

UNIVERSITÀ DEGLI STUDI DI  
MODENA E REGGIO EMILIA

DIPARTIMENTO DI SCIENZE  
FISICHE, INFORMATICHE, MATEMATICHE

Tesi di Laurea in

INFORMATICA

PROGETTAZIONE, REALIZZAZIONE ED  
ACCESSIBILITÀ DI UN DATABASE  
BIOMOLECOLARE SULLE SEQUENZE  
ULTRACONSERVATE DEL GENOMA UMANO

Relatore

Prof. Riccardo Martoglia

Candidato

Vincenzo Lomonaco

Correlatori

Prof. Federica Mandreoli  
Dott. Cristian Taccioli

Anno Accademico 2012/2013



UNIVERSITÀ DEGLI STUDI DI  
MODENA E REGGIO EMILIA

DIPARTIMENTO DI SCIENZE  
FISICHE, INFORMATICHE, MATEMATICHE

Tesi di Laurea in

INFORMATICA

PROGETTAZIONE, REALIZZAZIONE ED  
ACCESSIBILITÀ DI UN DATABASE  
BIOMOLECOLARE SULLE SEQUENZE  
ULTRACONSERVATE DEL GENOMA UMANO

Relatore

Prof. Riccardo Martoglia

Candidato

Vincenzo Lomonaco

Correlatori

Prof. Federica Mandreoli  
Dott. Cristian Taccioli

Anno Accademico 2012/2013



*Dedicata a Elisabetta Giubilato e Francesco Lomonaco,  
dall'amore sincero e riconoscente di un figlio  
verso i propri genitori.*

*Ringrazio tutti coloro i quali mi sono stati vicino o mi hanno supportato durante  
gli studi, arricchendone inevitabilmente il percorso.*

*Ringraziamenti particolari vanno al Relatore Prof. Martoglia ed ai Correlatori  
Prof. Mandreoli e Dott. Taccioli, per la correttezza, responsabilità e capacità con  
la quale hanno condotto il loro ruolo.*



# Introduzione

*L'emergere di computer ad alte prestazioni a basso costo e le tecnologie omiche ad alto throughput sono alla base di diversi progetti finalizzati alla costruzione di database bio-tecnologici, quali **RCSB Protein Data Bank**, **Gene Ontology**, **CMap**, **GEO NCBI**. Attualmente, tali banche dati contengono informazioni di vitale importanza e in una grande varietà di forme (sequenze genetiche, annotazioni testuali, dati in forma di grafo, etc) ma possono essere interrogate solo in modo indipendente e incompleto, rendendo estremamente difficile per un ricercatore poter rilevare anche le relazioni più elementari tra i dati.* [1]

La tesi di ricerca, dunque, è volta a soddisfare un nuovo paradigma di studio, progettazione e sviluppo di tecniche che consentano l'interoperabilità delle banche dati biologiche, dirette, inoltre, a sostenere modalità efficaci ed efficienti per la loro interrogazione.

Nello specifico essa ruota attorno la progettazione e realizzazione di una banca dati sulle *sequenze nucleotidiche ultraconservate*<sup>1</sup> del genoma umano e dati biologici relativi già esistenti.

Nella voluta risoluzione di noti problemi d'indagine come la sin troppo accentuata entropia delle informazioni inerenti, il progetto offre, dunque, innovativi *environments* d'investigazione ed interrogazione dei dati che rappresentano, sotto certi punti di vista, un'avanguardia nel mondo biologico della ricerca.

Basata su soluzioni portabili ed efficienti per il recupero di dati, come su sorprendenti compromessi di semplicità e capacità espressiva nell'interrogazione e rappresentazione degli stessi, la tesi di ricerca, infatti, è volta a segnare con decisione un nuovo traguardo nel campo delle sequenze nucleotidiche ultraconservate sia in termini di interrogabilità dei dati quanto di integrità e completezza.

---

<sup>1</sup>vd. Appendice **Elementi ultraconservati nel genoma umano**

## Premessa

Si tiene a precisare che lo scritto rappresenta il risultato di un lavoro frutto di sudate sfide ed intuizioni, demoralizzazioni e rallegramenti, nella collaborazione di due mondi e discipline, quali l'Informatica e la Biologia, così diverse ma così intrensicamente legate. Sia il lettore, dunque, così gentile da superare, in virtù di quanto pocanzi espresso, eventuali forme di imprecisione o imperfezioni sfuggite ora all'uno ora all'altro punto di vista.

## Obiettivi della Tesi

Entrando più nel dettaglio, la tesi può essere articolata attorno alla realizzazione di quattro obiettivi principali:

- **Recupero dei dati** d'interesse online mediante *Web Services*<sup>2</sup> e loro razionalizzazione e filtraggio in locale.
- **Studio, progettazione e Realizzazione** di una base di dati, o meglio degli *script*<sup>3</sup> SQL [2] portabili per la realizzazione di una base di dati sulle sequenze nucleotidiche ultraconservate del genoma umano.
- **Recupero e collegamento di un ontologia biomedica** sulle patologie umane per fornire una più completa gamma di possibilità nelle operazioni di *Data Mining*<sup>4</sup>
- **Studio di accessibilità ed interrogabilità** dei dati **mediante un portale web** che supporti diverse modalità di investigazione e visualizzazione e **relativa realizzazione**.

Nel corso dell'elaborato ogni obiettivo verrà approfondito con accuratezza seguendo l'ordine temporale di sviluppo come rammentato nella sezione seguente: **Struttura della Tesi**.

Verranno inoltre evidenziati, con meticolosa attenzione ai particolari, gli strumenti *Open Source*<sup>5</sup> utilizzati ed attraverso quale modalità. Verranno, altresì sottolineati punti di debolezza e le relative e possibili migliorie effettuabili.

---

<sup>2</sup>Sistema software progettato per supportare l'interoperabilità tra diversi elaboratori attraverso il web.

<sup>3</sup>In informatica, il termine script designa un tipo particolare di programma contrassegno da particolari caratteristiche di linearità e semplicità.

<sup>4</sup>Insieme di tecniche e metodologie che hanno per oggetto l'estrazione di un sapere o di una conoscenza a partire da grandi quantità di dati.

<sup>5</sup>Termine inglese che significa codice sorgente aperto.

## Struttura della tesi

Se da un lato, nel corso dell'elaborato, si sviluppano gli obiettivi enucleati logicamente nel paragrafo precedente, dall'altro l'organizzazione dei capitoli sottolinea maggiormente una suddivisione di tipo *temporale* per enfatizzare con maggior rilievo il tortuoso ed affascinante percorso effettuato per giungere ad una tesi di ricerca completa.

**Capitolo 1 - Background.** Lo stato dell'arte. Si risponderà ad alcune semplici domande: Quali sono i dati di partenza? Quali sono le tecnologie adoperate per interrogarli? Con quali risultati? È opportuno un'operazione di miglioria? In quali termini?

**Capitolo 2 - Sviluppo.** Lo sviluppo centrale della Tesi. Metodi, tecniche e sistemi sono valutati ed eventualmente progettati in questo capitolo: dai dati grezzi ai dati in forma normale, dalle problematiche di ridondanza a quelle di linking, dall'analisi visuale manuale a quella avanzata. Dall'analisi e progettazione si passa, poi, alla realizzazione concreta con le relative problematiche, le scelte di campo e nuove possibilità.

**Capitolo 3 - Risultati.** Presentazione dei risultati. Il prodotto finito e le ultime modifiche, eventuali rifiniture, possibilità di utilizzo.

**Capitolo 4 - Conclusioni e sviluppi futuri.** Conclusioni e possibili migliorie. Difetti e valutazioni tecnico-personali.

**Appendice A - Glossario di base di biologia molecolare.** Informazioni e termini di base riguardo il background biologico essenziale per affrontare con linearità l'argomentazione della tesi.

**Appendice B - Elementi ultraconservati nel genoma umano.** Dettaglio tecnico biologico sugli elementi ultraconservati e nello specifico sulle sequenze nucleotidiche ultraconservate del genoma umano.

**Appendice C - Archivio dei Codici.** Totalità del codice scritto nell'ambito del progetto di tesi, riportato in forma completa ed all'ultima versione utile per chiarimenti ed approfondimenti.



# Indice

<b>1</b>	<b>Background</b>	<b>11</b>
1.1	Tabelle xls e script in R . . . . .	11
1.2	Problematiche dei dati grezzi . . . . .	15
1.2.1	Ridondanza e rumore . . . . .	15
1.2.2	Incoerenza . . . . .	16
1.2.3	Analizzabilità dei dati . . . . .	16
1.3	Nuove necessità . . . . .	17
<b>2</b>	<b>Sviluppo</b>	<b>19</b>
2.1	Progettazione e realizzazione del DB . . . . .	20
2.1.1	diagramma E/R . . . . .	20
2.1.2	Progetto logico . . . . .	23
2.1.3	Script SQL . . . . .	24
2.1.4	Creazione DB in Mysql . . . . .	26
2.2	Script Java per il recupero dati e Insert SQL . . . . .	27
2.2.1	Perchè Java? . . . . .	27
2.2.2	Il codice . . . . .	27
2.2.3	Risultati . . . . .	33
2.3	Inserimento e linking della HDO . . . . .	34
2.3.1	Recupero e riduzione dell'ontologia . . . . .	34
2.3.2	Adattamento della struttura del Database . . . . .	40
2.3.3	Adattamento e creazione degli script Java . . . . .	42
2.4	Progettazione e realizzazione dell'interfaccia Web . . . . .	49
2.4.1	Struttura e modalità di accesso . . . . .	49
2.4.2	Cenni sulla realizzazione . . . . .	50
<b>3</b>	<b>Risultati</b>	<b>57</b>
<b>4</b>	<b>Conclusioni e sviluppi futuri</b>	<b>63</b>
<b>A</b>	<b>Glossario di base di biologia molecolare</b>	<b>65</b>
<b>B</b>	<b>Elementi ultraconservati nel genoma umano</b>	<b>68</b>

<b>C Archivio dei Codici</b>	<b>70</b>
C.1 Uc.biomaRt_pathology2.java . . . . .	70
C.2 Tabelle.sql . . . . .	71
C.3 BiomartMain.java . . . . .	73
C.4 GeneRecover.java . . . . .	73
C.5 SplicingRecover.java . . . . .	76
C.6 SnpRecover.java . . . . .	79
C.7 PathologyRecover.java . . . . .	81
C.8 TestCorrispondenze.java . . . . .	85
C.9 OboEditAllFilter.java . . . . .	88
C.10 PathologyObj.java . . . . .	91
C.11 NTree.java . . . . .	92
C.12 NTreeNode.java . . . . .	95
C.13 TreeMaker.java . . . . .	96
C.14 RiempipiPatTable.java . . . . .	98
C.15 Index.html . . . . .	99
C.16 Uc_data_mining.php . . . . .	102
C.17 Table.php . . . . .	106
C.18 Sequence.php . . . . .	109
C.19 Result_search_uc.php . . . . .	110
C.20 Result_search_pat.php . . . . .	110
C.21 Result_gene.php . . . . .	111
C.22 Result.php . . . . .	115
C.23 Related_works.html . . . . .	120
C.24 Blastpat.php . . . . .	123
C.25 Blastn.php . . . . .	129
C.26 About_us.html . . . . .	132
C.27 Style.css . . . . .	134
<b>Bibliografia</b>	<b>150</b>

# Capitolo 1

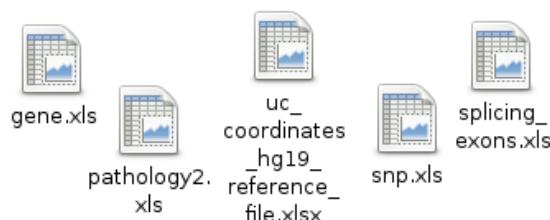
## Background

Al fine di rendere ben intellegibili le volontà che hanno spinto a profondere una tale quantità di energie, quali lo stesso studio di tesi o la redazione dell'elaborato in questione, si vuole illustrare con precisione lo stato tecnico operativo nel quale si versava inizialmente, prima dell'intervento di una vera collaborazione informatica nell'ambito biologico. Nei paragrafi successivi verrà dato ampio spazio, inoltre, alle problematiche ed alle incertezze introdotte da una scarsa accuratezza o formalità nell'archiviazione dei dati ed i danni che essa può produrre.

### 1.1 Tabelle xls e script in R

Nel seguente paragrafo si comincerà a fare riferimento oltre che a termini tecnici puramente informatici anche ad alcuni di matrice biologica. Si faccia riferimento, dunque, in caso di smarrimento, all'appendice **A - Glossario di base di biologia molecolare** per chiarimenti.

I dati grezzi, con i cui ci si poteva rapportare all'inizio del percorso, erano rappresentati da cinque tabelle nel formato *xls*<sup>1</sup>.



**Fig. 1.1:** Tabelle xls

---

<sup>1</sup>Estensione Microsoft Excel

Ognuno di questi file, seguendo per quanto possibile le intenzioni dell'ignoto progettista, avrebbe dovuto enucleare diversi concetti:

- **uc\_coordinates\_hg19\_reference\_file.xlsx**: File contenente le coordinate indispensabili ad ogni sequenza ultraconservata (identificatore della sequenza, numero cromosoma di appartenenza, coordinata d'inizio nel cromosoma, coordinata di fine, nome del filamento DNA).
- **gene.xls**: File contenente le informazioni circa i geni che contengono le sequenze ultraconservate (identificatore di sequenza coinvolta, nome del gene, identificatore univoco del gene, tipo di gene, nome del cromosoma di appartenenza, coordinata all'interno del cromosoma dell'inizio del gene, coordinata della fine del gene, banda cromosomica di appartenenza, nome del filamento di DNA in cui si trova il gene).
- **pathology2.xls**: File contenente le informazioni circa le patologie correlate ad un certo gene (identificatore ultraconservata coinvolta, identificatore gene correlato, descrizione *MIM*<sup>2</sup> circa la patologia genetica, nome della patologia, contenuto *GC*<sup>3</sup> del gene)
- **snp.xls**: File contenente informazioni circa le variazioni *SNP*<sup>4</sup> correlate ad una certa sequenza ultraconservata (identificatore di sequenza, id del polimorfismo, id dell'eventuale gene di appartenenza, *allele*<sup>5</sup> mutato, coordinata della variazione, nome del filamento di DNA sul quale è avvenuta la mutazione, Significatività clinica, descrizione fenotipica, eventuale validazione di un'ente).
- **plicing\_exons.xls**: File contenente informazioni circa gli *Splicing*<sup>6</sup> correlate ad una certo gene (identificatore di sequenza, id del gene correlato, tipo di Splicing, codice di Splicing e nome esteso).

---

<sup>2</sup>Mendelian Inheritance in Man: è una banca dati che cataloga tutte le patologie aventi una componente genetica.

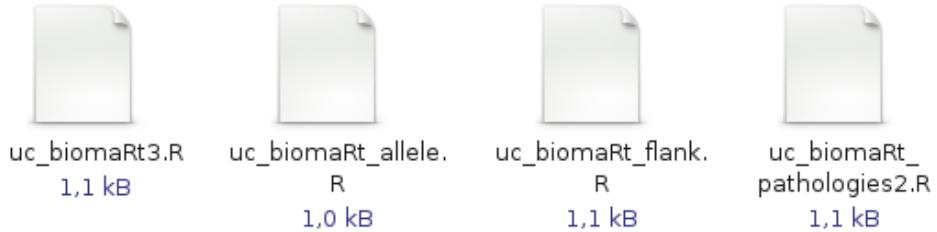
<sup>3</sup>Due delle cinque basi che compongono i nucleotidi del DNA e dell'RNA

<sup>4</sup>Single Nucleotide Polymorphism o SNP: polimorfismo a singolo nucleotide

<sup>5</sup>In genetica si definisce allele o fattore ogni variante di sequenza di un gene o di un locus genico

<sup>6</sup>In biologia molecolare e in genetica, splicing è una modifica del nascente pre-mRNA

I files in questione sono stati ottenuti mediante interrogazione del web service offerto da *BioMart*<sup>7</sup> [10]. Come? attraverso degli obsoleti script in *R*<sup>8</sup>:



**Fig. 1.2:** Script in R

Tuttavia come facilmente notificabile dai titoli e dalla disparità numerica rispetto al conto dei fai xls, questi script non combaciano con la creazione degli stessi, ma offrono, tuttavia, un ottimo spunto per capire come effettivamente il web service lavora e come è possibile interrogarlo. Di seguito il codice, da considerare come un esempio, del file **uc.biomaRt\_pathology2.R**:

```

1 library("biomaRt")
2 ensembl = useMart("ensembl",dataset="hsapiens_gene_ensembl")
3
4 options(width=120)
5 attributesR=c('ensembl_gene_id','mim_morbid_description',
6             'pathology','percentage_gc_content')
7 filtersR = c('chromosome_name','start','end')
```

Nella riga 1 importiamo il *package biomaRt*<sup>9</sup> [9]. Successivamente si costruiscono gli oggetti indispensabili all'interrogazione utilizzando le funzioni appropriate che il pacchetto biomaRt offre. Si sceglie il centro di raccolta dei database (*ensembl*<sup>10</sup>) e il database da interrogare (in questo caso il database genetico dell'*Homo Sapiens*). Infine gli *attributi* ossia le informazioni che vogliamo recuperare ed i *filtri* secondo i quali effettuare la ricerca.

---

<sup>7</sup>BioMart è un sistema software per rendere le informazioni biologiche più accessibili ai biologi

<sup>8</sup>R è un linguaggio di programmazione libero per l'analisi statistica dei dati e molto di più

<sup>9</sup>Un interfaccia scritta in R ai database offerti da BioMart attraverso il suo web service.

<sup>10</sup>Ensembl è una banca dati bioinformatica allestita con lo scopo di fornire informazioni aggiornate sui principali genomi eucariotici

```

8  out.list<-list()
9  for (i in 1:length(uc2009[,1])){
10    valuesR<-list(uc2009[i,2],uc2009[i,3],uc2009[i,4])
11    out<-getBM(attributes=attributesR,filters=filtersR,
12              values=valuesR, mart=ensembl,uniqueRows=F)
13    if (nrow(out)!=0){
14      out<-data.frame(cbind(uc2009[i,1],out))
15      colnames(out)[1]<-"uc"
16    }
17    else{
18      out.names<-colnames(out)
19      nulli<-as.data.frame(matrix(rep(NA,ncol(out)),nrow=1))
20      out<-data.frame(cbind(uc2009[i,1],nulli))
21      colnames(out)<-c("uc",out.names)
22    }
23
24    if(i==1)
25      write.table(out,"pathology.csv",append=FALSE,
26                  col.names=TRUE,quote=FALSE,row.names=FALSE,sep=";")
27    else
28      write.table(out,"pathology.csv",append=TRUE,
29                  col.names=FALSE,quote=FALSE,row.names=FALSE,sep=";")
30    out.list[[i]]<-out
31    print(uc2009[i,1])
32    print(out)
33  }

```

La seconda parte di codice cicla su ogni elemento dell'array uc2009 (ossia le sequenze ultraconservative), creato in precedenza nel workspace, ed utilizzando specifiche funzioni offerte da R scrive, infine, un output tabellare in formato *csv*<sup>11</sup>

Il risultato finale, se importiamo questi file csv in un software per l'elaborazione di fogli elettronici come ad esempio Open Office Calc, è modesto sia dal punto di vista estetico che da quello pragmatico d'indagine, ma questa intuizione verrà argomentata nel paragrafo successivo.

---

<sup>11</sup>comma-separated values (abbreviato in CSV) è un formato di file basato su file di testo utilizzato per l'importazione ed esportazione di una tabella di dati.

## 1.2 Problematiche dei dati grezzi

Appreso il funzionamento del web service e la grezza archiviazione dei dati apportata allo stato attuale, sono evidenziabili fin da subito problemi di elevata caratura per quanto concerne l'affidabilità e l'analizzabilità dei dati, fine ultimo per un biologo.

### 1.2.1 Ridondanza e rumore

La più evidente delle mancanze è senza dubbio la ridondanza dei dati presente sui più e più file.

A	B	C	D	
1	<b>Id</b>	<b>wikigene_name</b>	<b>ensembl_gene_id</b>	<b>gene_biotype</b>
2	uc.1	PEX14	ENSG00000142655	protein_coding
3	uc.2	CASZ1	ENSG00000130940	protein_coding
4	uc.3	CASZ1	ENSG00000130940	protein_coding
5	uc.4	CASZ1	ENSG00000130940	protein_coding
6	uc.5	CASZ1	ENSG00000130940	protein_coding
7	uc.6	CASZ1	ENSG00000130940	protein_coding
8	uc.7	CASZ1	ENSG00000130940	protein_coding
9	uc.8	CASZ1	ENSG00000130940	protein_coding

A	B	C	D	E	
1	<b>Id</b>	<b>refsnip_id</b>	<b>ensembl_gene_stable_id</b>	<b>allele</b>	<b>chrom_start</b>
2	uc.1	rs190053770	ENSG00000142655	G/A	10597738
3	uc.1	rs190053770	ENSG00000142655	G/A	10597738
4	uc.1	rs190053770	ENSG00000142655	G/A	10597738
5	uc.1	rs190053770	ENSG00000142655	G/A	10597738
6	uc.1	rs190053770	ENSG00000142655	G/A	10597738
7	uc.1	rs181426894	ENSG00000142655	A/C	10597822
8	uc.1	rs181426894	ENSG00000142655	A/C	10597822
9	uc.1	rs181426894	ENSG00000142655	A/C	10597822
10	uc.1	rs181426894	ENSG00000142655	A/C	10597822
11	uc.1	rs181426894	ENSG00000142655	A/C	10597822

A	B	C	
1	<b>Id</b>	<b>ensembl_gene_id</b>	<b>mim_morbid_description</b>
2	uc.1	ENSG00000142655	PEROXISOME BIOGENESIS FACTOR 14
3	uc.1	ENSG00000142655	PEROXISOME BIOGENESIS FACTOR 14
4	uc.1	ENSG00000142655	PEROXISOME BIOGENESIS FACTOR 14
5	uc.1	ENSG00000142655	PEROXISOME BIOGENESIS FACTOR 14
6	uc.1	ENSG00000142655	PEROXISOME BIOGENESIS FACTOR 14

**Fig. 1.3:** Ridondanza nel Formato tabellare

È subito ravvisabile nelle sue diverse forme: dalla ridondanza dovuta ai join tebellari effettuati lato server da BioMart a quelli introdotti dagli script lo-

cali in una mancata progettazione di archiviazione in forma normale dei dati stessi (si vedano le informazioni ripetute su più file).

Non meno evidenti le proporzioni di rumore (mancanza di dati informativi) presenti nei file.

310	uc.293	NA	NA
311	uc.294	NA	NA
312	uc.295	NA	NA
313	uc.296	NA	NA
314	uc.297	NA	NA
315	uc.298	NA	NA

**Fig. 1.4:** Rumore di fondo nel formato tabellare

Ridondanza e rumore accreccoscono le dimensioni dei file e diminuiscono la leggibilità dei dati ma introducono soprattutto i problemi elencati di seguito.

### 1.2.2 Incoerenza

Nel momento in cui si volesse aggiornare o modificare i dati sopracitati, ad esempio modificando l'identificatore di una sequenza ultraconservata, nulla potrebbe essere più drammatico: per ciascun file bisognerebbe manualmente trovare l'identificatore precedente e sostituirlo con il nuovo per ogni sua occorrenza. Non è molto difficile immaginare, come in un clima del genere, sia possibile incappare in grossolani errori di incoerenza a seguito di numerose modifiche o correzioni.

### 1.2.3 Analizzabilità dei dati

Quando ci si trova in situazioni analoghe la quantità di dati in presenza di rumore e ridondanze cresce vertiginosamente:

9947	uc.483	ENSG00000131374	NA
9948	uc.483	ENSG00000131374	NA
9949	uc.483	ENSG00000131374	NA
9950	uc.483	ENSG00000131374	NA
9951	uc.483	ENSG00000131374	NA
9952	uc.483	ENSG00000131374	NA

**Fig. 1.5:** Effetti ridondanza

Il danno è tutto a carico dell’analizzatore dei dati, il biologo in questo caso, che si ritrova costretto ad estrarre delle informazioni in condizioni sempre più difficilose.

### 1.3 Nuove necessità

Realizzate le problematiche evidenziate durante il capitolo, risulta dunque comprensibile la nuova volontà di un’evoluzione in una direzione tecnico-informatica per far fronte alle nuove necessità della biologia molecolare moderna: analizzare e gestire agevolmente grandi quantità di dati ma soprattutto estrarre **informazioni** in modo agevole, funzionale e creativo. In questa direzione sono sicuramente gli obiettivi tracciati e sicuramente raggiunti della tesi di ricerca: Un database per garantire le proprietà *ACID*<sup>12</sup> e tutte le cure necessarie per salvaguardare l’integrità e l’affidabilità dei dati; Un Portale web per rendere tali dati accessibili e navigabili agevolmente, senza nient’altro che un browser; La possibilità di analizzarli in senso creativo e trasversale: attraverso l’inserimento nel database di un’ontologia biomedica sulle patologie umane. Più avanti, nel capitolo successivo, sarà chiaro quanto a fronte dei dati grezzi tutte queste tecniche possano beneficiare ad un biologo, ricercatore o chiunque interessato a reperire informazioni sulle sequenze ultraconservate del genoma umano.

---

<sup>12</sup>Atomicità, Coerenza, Isolamento e Durabilità (ACID): Sono le proprietà logiche che devono avere le transazioni in un database moderno.



## **Capitolo 2**

### **Sviluppo**

Effettuata un'analisi accurata delle problematicità e necessità derivanti dall'uso di dati grezzi, analizzati nel capitolo precedente, è finalmente possibile passare alla fase creativa di sviluppo. Nel corso del capitolo verranno illustrate, seguendo un ordine temporale, le operazioni che porteranno dei dati sparpagliati, problematici e ridondanti a costituire una banca dati strutturata organicamente e funzionale all'interrogazione efficiente.

## 2.1 Progettazione e realizzazione del DB

Il primo e forse più importante passo per la buona riuscita del progetto di migliorie d'apportare è la progettazione del Database: Esso costituisce il punto nevralgico di organizzazione logica dei dati e il pilastro fondamentale per la garanzia di efficacia che il progetto vorrebbe rilanciare.

### 2.1.1 diagramma E/R

Si agisce innanzitutto isolando le entità che comporranno il nostro modello logico, ossia i concetti fondamentali che vorremmo rappresentare, e successivamente arricchendole con gli attribuiti, ossia le informazioni relative che vorremmo memorizzare per ciascuna entità:

- **Sequenza ultraconservata**

- **chr**: Identificativo del cromosoma di appartenenza
- **Start**: Numero di *bp*<sup>1</sup> dall'inizio del cromosoma che identifica l'inizio della sequenza.
- **End**: Numero di bp che identifica la fine della sequenza.
- **Sequence**: Sequenza reale delle basi azotate ATGC .
- **Strand**: Identificatore filamento di DNA dell'ultraconservata.
- **Uc\_name**: Nome mnemonico per identificare una sequenza.

- **Gene**

- **Gene\_start\_pos**: Numero bp identificante l'inizio del gene.
- **Gene\_end\_pos**: Numero bp identificante la fine del gene.
- **Ensembl\_gene\_id**: Identificativo univoco assegnato da *Ensembl*.
- **Band**: identificativo della banda cromosomica di appartenenza.
- **Strand**: Identificatore filamento di DNA di registrazione del gene.
- **Wikigene\_name**: Nome mnemonico per identificare un gene.
- **%G/C**: Percentuale delle basi azotate G e C all'interno del gene.
- **Gene\_biotype**: Tipo biologico di gene.

---

<sup>1</sup>Le coppie di basi o paia di basi (abbreviate come pb o, dall'inglese base pair, bp o bps) sono comunemente utilizzate come misura della lunghezza fisica di sequenze di acidi nucleici a doppio filamento

- **Patologia**

- **Id:** Identificativo univoco della patologia.
- **Nome:** Nome univoco esteso per definire la patologia.
- ...: Si vedrà nella sezione di linking dell'ontologia come migliorare gli attributi di questa entità.

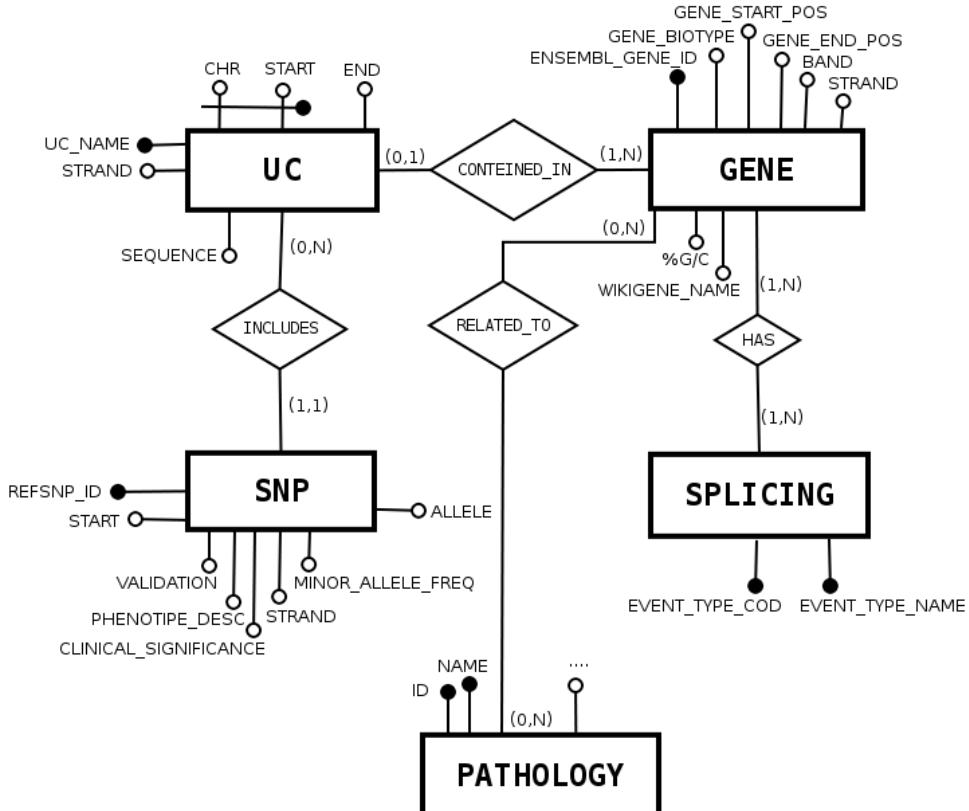
- **Splicing**

- **Event\_Type\_cod:** Codice univoco identificativo del tipo di splicing.
- **Event\_Type\_name:** Nome esteso identificativo del tipo di splicing.

- **SNP**

- **Refsnp\_id:** Codice univoco identificativo dell'SNP.
- **Start:** Identificativo numerico codificante la bp di inizio.
- **Allele:** Allele mutato.
- **Minor\_allele\_freq:** Frequenza di occorrenza dell'allele meno comune.
- **Phenotype\_desc:** Descrizione fenotipica.
- **Strand:** Identificativo del filamento di DNA su cui è avvenuta la mutazione.
- **Validation:** Ente che ha effettuato validazione.
- **Clinical\_Significance:** Significatività clinica.

Lo step successivo è il delineamento delle relazioni tra le sopracitate entità. Di seguito si propone il diagramma *E/R*<sup>2</sup> del Progetto:



**Fig. 2.1:** Diagramma E/R

Si presti particolare attenzione alle cardinalità delle relazioni e le chiavi delle entità. Come facilmente ravvisabile, sono già scomparse le ridondanze rilevate nei dati grezzi in formato tabellare. Nel prossimo paragrafo verrà sviluppato il primo passaggio per la traduzione del modello logico in un vero e proprio database: il progetto logico.

---

<sup>2</sup>Il modello entità-relazione in informatica è un modello per la rappresentazione concettuale dei dati ad un alto livello di astrazione.

### 2.1.2 Progetto logico

Il progetto logico rappresenta la traduzione schematica di quanto proposto sottoforma di E/R. Esso costituisce, in sunto, un valido canovaccio per la stesura del codice *SQL*<sup>3</sup> per la creazione del database. Di seguito, più esplicativo di mille parole, il risultato:

```
UC (UC_NAME, CHR, START,END, STRAND, SEQUENCE, ENSAMBL_GENE_ID)
AK: UC_NAME
FK: ENSAMBL_GENE_ID REFERENCES GENE

SNP (REFSNP_ID, START, ALLELE, VALIDATION, MINOR_ALLELE_FREQ,
      PHENOTYPE_DESC, CLINICAL_SIGNIFIANCE, STRAND, CHR, START)
FK: CHR, START REFERENCES UC NOT NULL

SPLICING (EVENT_TYPE_COD,EVENT_TYPE_NAME)
AK: EVENT_TYPE_NAME

HAS (EVENT_TYPE_COD, ENSEMBL_GENE_ID)
FK: EVENT_TYPE_COD REFERENCES SPLICING
FK: ENSEMBL_GENE_ID REFERENCES GENE

PATHOLOGY (ID, NAME, ...)
AK: NAME

RELATED_TO (ID, ENSAMBL_GENE_ID)
FK: ID REFERENCES PATHOLOGY
FK: ENSAMBL_GENE_ID REFERENCES GENE

GENE (ENSEMBL_GENE_ID, GENE_BYOTIPE, GENE_START_POS,
      GENE_END_POS, BAND, STRAND, WIKIGENE_NAME, G/C_PERC)
```

Effettuata la traduzione dal modello allo schema logico non resta nient'altro da fare che tradurre quanto proposto in codice SQL, che viene proposto nel paragrafo successivo.

---

<sup>3</sup>In informatica SQL (Structured Query Language) è un linguaggio standardizzato per database basati sul modello relazionale

### 2.1.3 Script SQL

Di seguito viene riportato il codice SQL per la creazione delle tabelle. Si rammenta che si tratta di un codice portabile in un qualsiasi *DBMS*<sup>4</sup> relazionale moderno.

```
34 CREATE TABLE GENE
35 (
36     ENSAMBL_GENE_ID VARCHAR(20),
37     GENE_BYOTYPE VARCHAR(30),
38     GENE_START_POS INTEGER,
39     GENE_END_POS INTEGER,
40     BAND VARCHAR(10),
41     STRAND INTEGER,
42     WIKIGENE_NAME VARCHAR(30),
43     G_C_PERC VARCHAR(10),
44     PRIMARY KEY (ENSAMBL_GENE_ID)
45 );
46
47 CREATE TABLE UC
48 (
49     UC_NAME VARCHAR(10) UNIQUE NOT NULL,
50     CHR VARCHAR(5),
51     START INTEGER,
52     ENDING INTEGER,
53     STRAND INTEGER,
54     SEQUENCE VARCHAR(1000) NOT NULL,
55     ENSAMBL_GENE_ID VARCHAR(20),
56     PRIMARY KEY (CHR,START),
57     FOREIGN KEY (ENSAMBL_GENE_ID) REFERENCES GENE
58         ON DELETE SET NULL ON UPDATE CASCADE
59 );
60
61 CREATE TABLE SNP
62 (
63     REFSNP_ID VARCHAR(30),
64     START INTEGER,
65     ALLELE VARCHAR(10),
66     VALIDATION VARCHAR(30),
67     MINOR_ALLELE_FREQ FLOAT(10),
68     PHENOTYPE_DESC VARCHAR(30),
69     CLINICAL_SIGNIFIANCE VARCHAR(30),
70     STRAND INTEGER,
71     CHR VARCHAR(5),
72     START_UC INTEGER,
73     PRIMARY KEY (REFSNP_ID),
```

---

<sup>4</sup>In informatica, un Database Management System (abbreviato in DBMS) è un sistema software progettato per consentire la creazione e la manipolazione (da parte di un amministratore) e l'interrogazione efficiente (da parte di uno o più utenti) di database.

```

74   FOREIGN KEY(CHR, START_UC) REFERENCES UC (CHR,START)
75   ON DELETE CASCADE ON UPDATE CASCADE
76 );
77
78
79 CREATE TABLE SPLICING
80 (
81   EVENT_TYPE_COD VARCHAR(8),
82   EVENT_TYPE_NAME VARCHAR(30) UNIQUE NOT NULL,
83   PRIMARY KEY(EVENT_TYPE_COD)
84 );
85
86 CREATE TABLE HAS
87 (
88   EVENT_TYPE_COD VARCHAR(8),
89   ENSAMBL_GENE_ID VARCHAR(20),
90   PRIMARY KEY(EVENT_TYPE_COD, ENSAMBL_GENE_ID),
91   FOREIGN KEY(EVENT_TYPE_COD) REFERENCES SPLICING
92   ON DELETE CASCADE ON UPDATE CASCADE,
93   FOREIGN KEY(ENSAMBL_GENE_ID) REFERENCES GENE
94   ON DELETE CASCADE ON UPDATE CASCADE
95 );
96
97 CREATE TABLE PATHOLOGY
98 (
99   ID VARCHAR(15),
100  NAME VARCHAR(100) NOT NULL,
101  PRIMARY KEY(ID)
102 );
103
104 CREATE TABLE RELATED_TO
105 (
106   ID VARCHAR(15),
107   ENSAMBL_GENE_ID VARCHAR(20),
108   PRIMARY KEY(ID, ENSAMBL_GENE_ID),
109   FOREIGN KEY(ID) REFERENCES PATHOLOGY
110   ON DELETE CASCADE ON UPDATE CASCADE,
111   FOREIGN KEY(ENSAMBL_GENE_ID) REFERENCES GENE
112   ON DELETE CASCADE ON UPDATE CASCADE
113
114 );

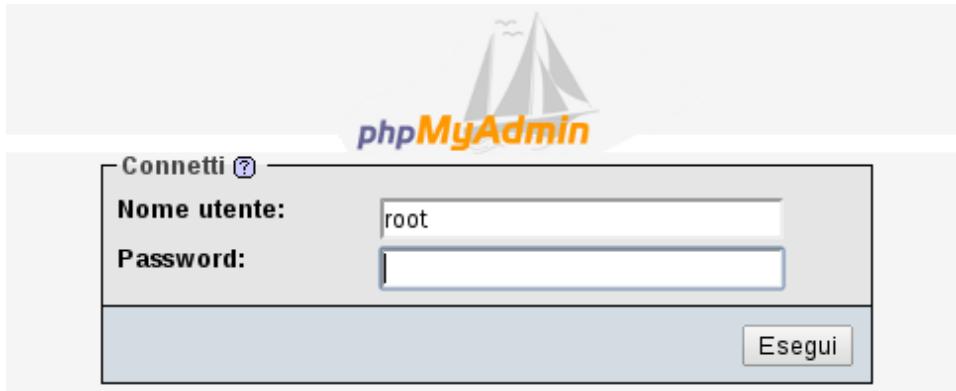
```

Si rammenta che per la tabella Pathology verranno apportate ulteriori migliorie nel paragrafo relativo all'inserimento e linking dell'ontologia biomedica.

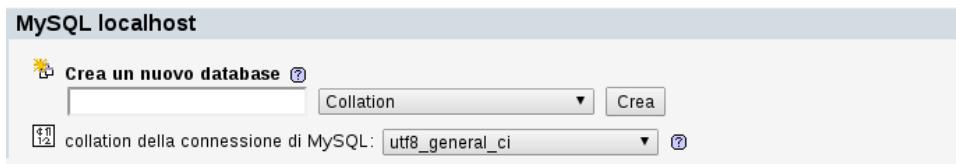
## 2.1.4 Creazione DB in Mysql

Ultimo step per la creazione effettiva del database è l'importazione nel DBMS dello script SQL visto nel paragrafo precedente: di seguito viene illustrata in massima sintesi la procedura adottata sul DBMS *Mysql* attraverso l'interfaccia scritta in *php*<sup>5</sup> *phpmyAdmin*.

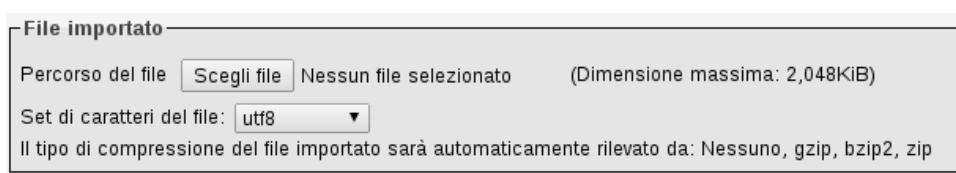
Innanzitutto si accede tramite “localhost/phpmyadmin” inserendo opportunamente le credenziali di accesso.



Succesivamente si accede alla procedura guidata per creare un nuovo database.



Infine si utilizza la scheda *importa* per importare un file sql da eseguire.



**Fig. 2.2:** Phpmyadmin

<sup>5</sup>PHP (acronimo ricorsivo di “PHP: Hypertext Preprocessor”) è un linguaggio di programmazione interpretato, originariamente concepito per la programmazione Web

## 2.2 Script Java per il recupero dati e Insert SQL

La sezione seguente si sviluppa nella direzione di realizzare gli script SQL di inserimento per il riempimento delle tabelle create in precedenza. Ci si rifà, in senso procedurale, a quanto detto riguardo il reperimento dei dati attraverso gli script in R e l'interrogazione del web service offerto da Bio-Mart. In questo caso, tuttavia, verrà utilizzato il linguaggio Java e gli script saranno impostati per ottenere i dati da biomart e per creare *direttamente* gli script SQL, in modo da enucleare al meglio, a livello concettuale, queste funzionalità.

### 2.2.1 Perchè Java?

La scelta di Java forse non risulta essere la più indicata per la prototipizzazione rapida non essendo un linguaggio prettamente dinamico. Tuttavia era stata presa in considerazione ed infine scelta perchè, stando alla documentazione biomart, erano offerti dei framework di lavoro per l'interrogazione del web service attraverso *SOAP API*<sup>6</sup> (ma che in realtà per scelte successive non verranno utilizzate). Inoltre risultava un linguaggio molto conosciuto e dalla indubbia portabilità, lasciando aperta ogni possibilità su future migliorie nel recupero *automatico* dei dati, ad esempio fornendo un'agevole *UI*<sup>7</sup>.

### 2.2.2 Il codice

In ambiente di sviluppo *Eclipse*<sup>8</sup> sono stati scritti 5 script principali:



**Fig. 2.3:** Script java

<sup>6</sup>In informatica SOAP (inizialmente acronimo di Simple Object Access Protocol) è un protocollo leggero per lo scambio di messaggi tra componenti software.

<sup>7</sup>L'interfaccia utente, anche conosciuta come UI (dall'inglese User Interface), è ciò che si frappone tra la macchina e l'utente, consentendo l'interazione tra i due.

<sup>8</sup>Eclipse è un ambiente di sviluppo integrato multi-linguaggio e multipiattaforma.

Nello specifico essi rappresentano:

- **BiomartMain.java**: Il main che inizializza le classi e avvia i metodi
- **GeneRecover.java**: Lo Script che recupera le informazioni inerenti le sequenze ultraconservate da file e quelle sui geni eventualmente correlati dal web service biomart. Quindi crea i uno script SQL per il riempimento delle tabelle UC e GENE.
- **SplicingRecover.java**: Lo Script che recupera le informazioni inerenti lo splicing dei geni e riempie le tabelle sql SPlicing ed HAS.
- **SNPRecover.java**: Lo Script che recupera le informazioni inerenti le SNP e quindi riempie la tabelle SNP.
- **PathologyRecover.java**: Lo Script che recupera le informazioni inerenti le Pathology e quindi riempie la tabelle PATHOLOGY e RELATED\_TO (lo script verrà rivisto poi alla luce delle migliorie introdotte a seguito del linking dell'ontologia biomedica).

A seguire si analizzerà solo uno di questi script, poichè essi agiscono tutti pressochè in modo analogo<sup>9</sup>. Si prende in analisi lo script **GeneRecover.java** analizzandone solo gli stralci principali, privi di controlli di robustezza per semplificarne la comprensione.

Nella primissima parte vengono importati gli oggetti di base per lo sviluppo. Successivamente si descrive la classe e gli attributi indispensabili al processo.

```
1 import java.io.*;
2 import java.net.URL;
3 import java.net.URLEncoder;
4 import java.util.StringTokenizer;
5
6 public class GeneRecover {
7     //dati per recupero info uc
8     private final int num_lines = 481;
9     private String[] chr = new String[num_lines];
10    private String[] bp_start = new String[num_lines];
11    private String[] bp_end = new String[num_lines];
12    private String[] strand_uc = new String[num_lines];
13    private String[] sequence = new String[num_lines];
14    strand_uc[lineNumber] = st.nextToken();
15    //dati per tabelle
16    private String ensambl_gene_id, gene_byotipe, gene_start_pos,
17        gene_end_pos, band, strand, wikigene_name,
18        g_c_perc, sql_gene, sql_uc;
```

<sup>9</sup>Si rimanda all'appendice **Archivio dei Codici** per approfondimenti

Tutto il lavoro viene svolto nel costruttore in modo da non dover chiamare nessun metodo e rendere il main il più semplice possibile. Vengono creati ed indirizzati opportunamente i descrittori di input ed output. Il file **uc\_coordinates\_hg19\_reference\_file.csv** rappresenta i metadati relativi a ciascuna sequenza nel formato:

```
"Id","Chromosome","Start","End","Strand"
"uc.1","1",10597697,10597903,1
"uc.2","1",10732543,10732749,1
```

```
1 public GeneRecover()throws Exception{
2     //leggo dal file le info
3     try{
4         //tsv file containing data
5         String strFile_uc =
6             "/home/uc_coordinates_hg19_reference_file.csv";
7         String strfile_seq =
8             "/home/ultraconservative/UC.txt";
9
10        /* Si creano i descrittori per scrivere i file sql
11           delle tabelle "UC" e "GENE". */
12        BufferedReader br =
13            new BufferedReader( new FileReader(strFile_uc));
14        BufferedReader br2 =
15            new BufferedReader( new FileReader(strfile_seq));
16        FileOutputStream sql_file =
17            new FileOutputStream("/home/uc_gene_insert.sql");
18        PrintStream out =
19            new PrintStream(sql_file);
```

Vengono recuperate le sequenze reali dal file UC.txt formattato come quanto di seguito:

```
>uc.1
TCCACCGACAATGACCAGTTAGTCCTCATTCTCTCAA...
>uc.2
GCCCGCCCCCCCCTCCCCGGGCCAATCTGTTTCAAAGTG...
```

```
1     //recupero sequenze uc
2     int i = 1, j=0, k=0;
3     while(i <= num_lines*2 ){
4         strLine = br2.readLine();
5         if(i%2 == 0) {sequence[j] = strLine; j++;}
6         else {uc_name[k] = strLine.substring(1); k++;}
7         i++;
8     }
```

Si da avvio al ciclo principale che itera sul file delle uc recuperando i metadati per l'interrogazione del web service e inserendoli negli attributi della classe secondo l'indice *lineNumber* incrementato ad ogni ciclo.

```

1 StringTokenizer st = null;
2 int lineNumber = 0;
3 //la prima riga contenente meta-info viene scartata
4 br.readLine();
5 /* ciclo sulle sequenze
6 nel file uc_coordinates\hg19\reference\file.csv */
7 while( (strLine = br.readLine()) != null ){
8     //break comma separated line using ","
9     st = new StringTokenizer(strLine, ",");
10    String uc = st.nextToken();
11    chr[lineNumber] = st.nextToken();
12    //togliamo fastidiosi apici
13    chr[lineNumber] = chr[lineNumber].replace("\\"", "");
14    bp_start[lineNumber] = st.nextToken();
15    bp_end[lineNumber] = st.nextToken();
16    strand_uc[lineNumber] = st.nextToken();

```

Si definisce, finalmente, il metodo di interrogazione: Se nel secondo capitolo, si era utilizzato il pacchetto *biomaRt* che metteva a disposizione i suoi metodi per l'interrogazione del web service, in questo caso si è scelto di effettuare la richiesta nel modo più semplice e naturale possibile, ossia tramite *REST*<sup>10</sup> API. Viene definito, dunque, un file xml per la strutturazione della query: Si scelgono in sequenza, il dataset su cui effettuare l'indagine, i filtri per la richiesta e gli attributi desiderati. Ad ogni ciclo, ovviamente, i filtri cambiano producendo il risultato atteso.

```

17 String myxml = "<Query virtualSchemaName = \"default\""
18         formatter = \"CSV\" header = \"1\" "
19         +"uniqueRows = \"0\" count = \"\""
20         datasetConfigVersion = \"0.6\" >"
21         +<Dataset name = \"hsapiens_gene_ensembl\""
22             interface = \"default\" >
23         +<Filter name = \"chromosome_name\""
24             value = \"\"+chr[lineNumber]+\"\"/>"
25         +<Filter name = \"start\""
26             value = \"\"+bp_start[lineNumber]+\"\"/>"
27         +<Filter name = \"end\""
28             value = \"\"+bp_end[lineNumber]+\"\"/>"
29         +<Attribute name = \"wikigene_name\" />"
30         +<Attribute name = \"ensembl_gene_id\" />"

```

---

<sup>10</sup>Representational state transfer (REST) è un tipo di architettura software per i sistemi di ipertesto distribuiti come il World Wide Web.

```

31    +"<Attribute name = \"gene_biotype\" />"  

32    +"<Attribute name = \"start_position\" />"  

33    +"<Attribute name = \"end_position\" />"  

34    +"<Attribute name = \"band\" />"  

35    +"<Attribute name = \"strand\" />"  

36    +"<Attribute name = \"percentage_gc_content\" />"  

37    +"</Dataset>"  

38    +"</Query>";

```

Una volta definita la stringa xml, bastano poche righe di codice per inoltrare la richiesta al web service ed ottenere risposta nel BufferedReader.

```

1  String encoded = URLEncoder.encode(myxml, "utf-8");  

2  URL url = new URL("http://www.biomart.org/biomart/martservice?  

   query="+encoded);  

3  InputStream response = url.openStream();  

4  BufferedReader reader = new BufferedReader(new InputStreamReader  

   (response));

```

All'interno del ciclo principale viene nidificato un'altro ciclo per l'analisi iterativa dei diversi record restituitici da biomart. Se non viene restituito nulla, si scrive su file solo il codice SQL per effettuare l'inserimento della sequenza ultraconservata, altrimenti, lo si accompagna al codice per l'inserimento del gene e le sue informazioni relative.

```

1  String line;  

2  /*nel caso non ci sia nessuna risposta  

   devo comunque inserire l'uc*/  

3  while (True){  

4      if((line = reader.readLine()) == null){  

5          sql_uc = "INSERT IGNORE INTO UC VALUES('"+  

6              uc_name[lineNumber]+",'"+ chr[lineNumber]+  

7              "',"+ bp_start[lineNumber]+  

8              "',"+ bp_end[lineNumber]+",  

9              "+strand_uc[lineNumber]+  

10             ", '"+sequence[lineNumber]+',null);";  

11         out.println(sql_uc);  

12         break;  

13     }  

14  }  

15 else{  

16     //estrapoliamo le info dal record restituitoci  

17     st = new StringTokenizer(line, ",");  

18     wikigene_name = st.nextToken();  

19     if(!wikigene_name.equals("null"))  

20         wikigene_name = "','" +wikigene_name+',';  

21     ensambl_gene_id = st.nextToken();

```

```

22     if(!ensambl_gene_id.equals("null"))
23         ensambl_gene_id = ""+ensambl_gene_id+"";
24     gene_byotipe = st.nextToken();
25     if(!gene_byotipe.equals("null"))
26         gene_byotipe = ""+gene_byotipe+"";
27     gene_start_pos = st.nextToken();
28     gene_end_pos = st.nextToken();
29     band = st.nextToken();
30     if(!band.equals("null"))
31         band = ""+band+"";
32     strand = st.nextToken();
33     if(!strand.equals("null"))
34         strand = ""+strand+"";
35     g_c_perc = st.nextToken();
36
37     sql_gene = "INSERT IGNORE INTO GENE VALUES("+
38             ensambl_gene_id+", "+gene_byotipe+", "+
39             gene_start_pos+", "+gene_end_pos+", "+
40             band+", "+strand+", "+wikigene_name+", "+
41             +g_c_perc+");";
42     //scriviamo su file
43     out.println(sql_gene);
44
45     sql_uc = "INSERT IGNORE INTO UC VALUES('"+
46             uc_name[lineNumber]+', '+ chr[lineNumber]+', '+
47             bp_start[lineNumber]+
48             ", "+ bp_end[lineNumber]+", "+strand_uc[lineNumber]+
49             ", "+sequence[lineNumber]+", "
50             + ensambl_gene_id +");";
51     //scriviamo su file
52     out.println(sql_uc);
53 }
54 }
```

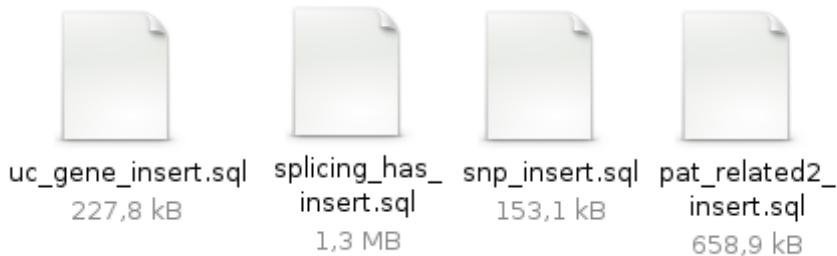
Nel main bisognerà semplicemente stanziare un oggetto della classe descritta, così come per le altre costruite seguendo lo stesso scheletro logico, per ottenere tutti gli script SQL per il riempimento della banca dati.

```

1 public class BiomartMain {
2     public static void main(String[] args) throws Exception {
3
4         GeneRecover gene = new GeneRecover();
5         SnpRecover.snp = new SnpRecover();
6         SplicingRecover spli = new SplicingRecover();
7         PathologyRecover path = new PathologyRecover();
8     }
```

### 2.2.3 Risultati

Eseguendo il main, otteniamo 4 script SQL per il riempimento automatico delle tabelle:



**Fig. 2.4:** Script SQL

Ciascuno dei quali avrà, all'incirca, lo stesso formato:

```
1 INSERT IGNORE INTO GENE VALUES('ENSG00000142655',  
                                'protein_coding',10532345,10690815,'p36.22','1','PEX14',45.03);  
2 INSERT IGNORE INTO UC VALUES('uc.1','1',10597697,10597903,1,  
                               'TCCACCGA...', 'ENSG00000142655');  
3 INSERT IGNORE INTO GENE VALUES('ENSG00000130940',  
                                'protein_coding',10696661,10856707,'p36.22','1','CASZ1',56.73);  
4 INSERT IGNORE INTO UC VALUES('uc.2','1',10732543,10732749,1,  
                               'GCCCGCCCCCTC...', 'ENSG00000130940');  
5 INSERT IGNORE INTO GENE VALUES('ENSG00000130940',  
                                'protein_coding',10696661,10856707,'p36.22','1','CASZ1',56.73);  
6 .....
```

Notare che sarà lo stesso DB (oltre che lo script nella sua interezza), grazie alla sua struttura, a filtrare le ridondanze ed il rumore nei dati. La clausola IGNORE viene qui utilizzata per trasformare gli *errors* dovuti alla presenza di record duplicati in *warnings* offrendoci un rendiconto generale sui problemi riscontrati ma senza interrompere il processo di inserimento.

Successivamente per ottenere un Database funzionante ed interrogabile basterà caricare gli script sopracitati allo stesso modo di quanto visto per l'importazione degli script per la creazione delle tabelle.

## 2.3 Inserimento e linking della HDO

Allo stato delle cose, il database risulta completo e funzionante. Tuttavia, si vorrebbe ottenere informazioni più dettagliate circa le patologie e soprattutto poter interrogare la basi di dati effettuando query strategiche che sappiano recuperare anche contenuti meno apparenti. L'idea è quella di introdurre nel Database una vera e propria ontologia medica sulle patologie umane, o quantomeno parte di essa, completa nelle informazioni gerarchiche tra le patologie. In questo modo si renderanno possibili interrogativi del tipo: *ottenere tutte le sequenze ultraconservative che hanno una qualche relazione con il cancro o sue sottopatologie.*

### 2.3.1 Recupero e riduzione dell'ontologia

Tra le tante ontologie accessibili, come la celeberrima *SNOMED* [14], la scelta è ricaduta sulla *Human disease ontology* consultabile e scaricabile gratuitamente grazie all'esperimento collaborativo dell'*OBO Foundry* [15]. Il file in formato *obo*<sup>11</sup> segue la formattazione seguente:

```
[Term]
id: DOID:7479
name: duodenal somatostatinoma
synonym: "duodenal delta cell somatostatin producing tumor" EXACT []
xref: NCI:C27407
xref: UMLS_CUI:C1333320
is_a: DOID:10021 ! duodenum cancer

[Term]
id: DOID:7480
name: large cell carcinoma with rhabdoid phenotype
synonym: "large cell carcinoma with rhabdoid phenotype
(morphologic abnormality)"
EXACT [SNOMEDCT_2005_07_31:128629005]
synonym: "large cell lung carcinoma with Rhabdoid Phenotype"
EXACT [NCI2004_11_17:C6876]
xref: NCI:C6876
xref: SNOMEDCT_2010_1_31:128629005
xref: UMLS_CUI:C1265997
is_a: DOID:4556 ! lung large cell carcinoma
```

Le enormi dimensioni dell'ontologia rappresentano, tuttavia, più un'ostacolo al progetto che un valore aggiunto, poiché solo una piccola parte delle

---

<sup>11</sup>Open Biomedical Ontology: formato adoperato per rappresentare esplicitamente significato e semantica di termini con vocabolari e relazioni tra gli stessi.

patologie presenti risulta essere effettivamente interessante allo scopo. È opportuno, dunque, estrarre solamente le patologie suggeriteci da Biomart in relazione a ciascun gene e **tutte quelle antenate sino alla radice**. Quello che si estrae, quindi, è un sotto-grafo minimale di **201** patologie a fronte **6393** elementi validi ed una centinaia di elementi considerati obsoleti del grafo originale.

**Come è stato possibile compiere questa meticolosa estrazione?** Innanzitutto tramite un nuovo script java, che per brevità e per la sua funzione strumentale non verrà mostrato nel dettaglio, si esaminano tutte le patologie recuperate da biomart e le si confronta con tutti i termini dell'ontologia nel formato precedentemente illustato. Se riusulta esserci un match tra in nome della patologia ed il nome del termine od un suo sinonimo allora quel termine viene rilevato dall'ontologia. Il secondo passo è l'indispensabile linking manuale delle patologie che biomart correla a determinati geni ma che non sono stati rilevati nell'ontologia o perchè troppo generici o perchè troppo particolarizzanti. Quello che dal punto informatico è possibile implementare è una serie di *suggerimenti* per il biologo tenuto ad effettuare il linking<sup>12</sup>. I file dei suggerimenti è ottenuto mediante l'algoritmo per calcolare la *Levenshtein distance*<sup>13</sup> di cui si è recuperata ed utilizzata la versione in java:

```

1 public static int computeLevenshteinDistance(CharSequence str1,
2                                     CharSequence str2) {
3     int[][] distance =
4         new int[str1.length() + 1][str2.length() + 1];
5
6     for (int i = 0; i <= str1.length(); i++)
7         distance[i][0] = i;
8     for (int j = 1; j <= str2.length(); j++)
9         distance[0][j] = j;
10
11    for (int i = 1; i <= str1.length(); i++)
12        for (int j = 1; j <= str2.length(); j++)
13            distance[i][j] = minimum(
14                distance[i - 1][j] + 1,
15                distance[i][j - 1] + 1,
16                distance[i - 1][j - 1]
17                + ((str1.charAt(i - 1) == str2.charAt(j - 1)) ? 0 : 1));
18
19    return distance[str1.length()][str2.length()];
20 }
```

---

<sup>12</sup>Si veda *TestCorrispondenze.java* nell'archivio dei codici per approfondimenti

<sup>13</sup>la distanza di Levenshtein, o distanza di edit, è una misura per la differenza fra due stringhe.

Per ottenere l'elenco dei suggerimenti basterà, in maniera piuttosto semplice, ciclare su tutte le patologie di cui ancora non è stato trovato il corrispettivo nell'ontologia effettuando un controllo del tipo:

```

1  ris= computeLevenshteinDistance(do_name.get(k), patmim_name.get(i));
2  //se le differenze non sono piu' della metà' dei caratteri
3  if(ris < patmim_name.get(i).length()/2)
4      out_file.println
5      (
6          patmim_name.get(i) +
7          " somiglia a: " + do_name.get(k)
8      );

```

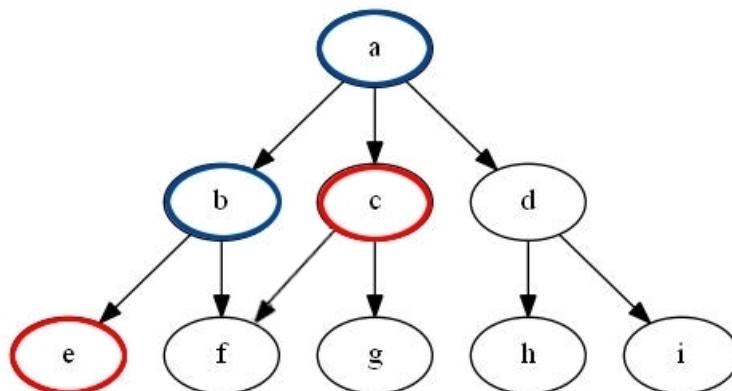
Come risultato si otterrà un file di testo con il quale un biologo esperto potrà agevolmente completare il lavoro di linking:

```

leukaemia somiglia a: leukemia
leukaemia somiglia a: leukopenia
lymphocytic somiglia a: lymphocelle
lymphocytic somiglia a: lymphopenia
goitre somiglia a: goiter
.....

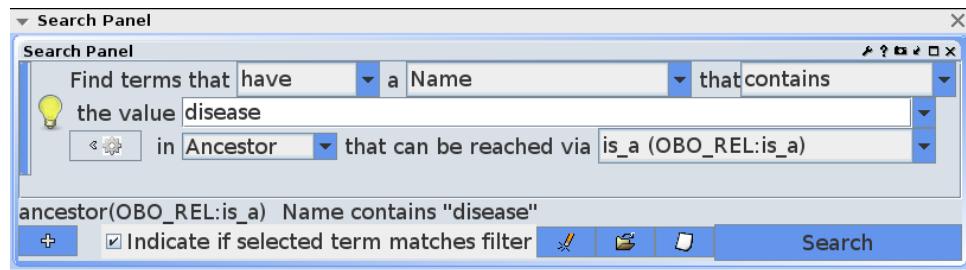
```

Ma il lavoro per l'estrapolazione del sub-graph non è ancora concluso: bisognerà trovare un algoritmo o delle tecniche efficienti per venire a conoscenza di tutte quelle patologie coinvolte nei cammini (anche più di uno) che collegano quelle di cui già si è a conoscenza alla radice dell'ontologia stessa. Solo in questo modo si potrà mantenere inalterata l'informazione gerarchica dell'ontologia. Si consideri l'esempio sottostante:



**Fig. 2.5:** Grafo d'esempio

Nell'immagine in figura si consideri in rosso (*c*, *e*) le patologie suggerite da biomart di cui si è ricondotto il termine corrispondente nell'ontologia. Ovviamente non si vuole prendere in considerazione il ramo del nodo *d* poichè non costituisce informazione rilevante ed in una futura interrogazione restituirebbe solo risultati vuoti circa le possibili sequenze o geni corelati. Saranno solo nodi cerchiati in blu (*a*, *b*) ad essere inseriti addizionalmente nel sotto-grafo, invece, perchè utili per generalizzare interrogazioni ottenendo informazioni sempre maggiori. Nell'esplorazione delle tecniche migliori per ottenere con il minor sforzo possibile gli effetti desiderati, risulta essere opportuno l'utilizzo del software open source *OBO-Edit* [19], un editor e visualizzatore di ontologie scritto in Java. Nello specifico sarà utilizzato il suo motore di ricerca ed il suo navigatore visuale per navigare il sottografo finale ed assicurarsene la correttezza. Nella figura sottostante viene riportata la modalità semi-strutturata di interrogazione che verrà utilizzata: si specifica, ad esempio, che si vuole ottenere tutti i nodi dell'ontologia che hanno come antenato il nodo root *disease*.



**Fig. 2.6:** Obo-Edit: Search panel

Tuttavia, compiere a mano questo lavoro, risulterebbe lungo e tedioso, ma fortunatamente, anche se non contemplata da OBO-edit nella sua documentazione è possibile un'altra soluzione: Il software sopracitato mette a disposizione un metodo per il salvataggio ed il caricamento delle query sotomesse al suo motore di ricerca; fortunatamente esse sono salvate in un intellegibile xml che segue pressocchè la seguente struttura:

```

1 <?xml version="1.0" encoding="UTF-8"?>
2 <java version="1.6.0_18" class="java.beans.XMLDecoder">
3   <object class="org.obo.filters.CompoundFilterImpl">
4     <void property="booleanOperation">
5       <int>1</int>
6     </void>
7     <void property="filters">
8       <void method="add">
9         <object class="org.obo.filters.ObjectFilterImpl">
10           <void property="aspect">
```

```

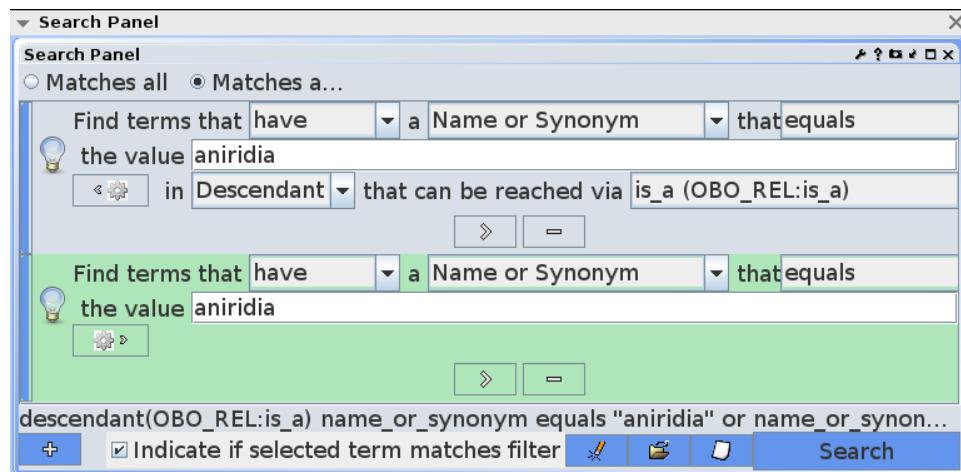
11      <object class="org.obo.filters.DescendantSearchAspect"/>
12    </void>
13    <void property="comparison">
14      <object id="EqualsComparison0" class="org.obo.filters.
15        EqualsComparison"/>
16    </void>
17    <void property="criterion">
18      <object id="NameSynonymSearchCriterion0" class="org.obo.
19        filters.NameSynonymSearchCriterion"/>
20    </void>
21    <void property="reasoner">
22      <object id="RuleBasedReasoner0" class="org.obo.reasoner.rbr.
23        RuleBasedReasoner">
24        <void property="properties">
25          <object class="java.util.HashSet"/>
26        </void>
27        </object>
28      </void>
29      <void property="traversalFilter">
30        <object class="org.obo.filters.LinkFilterImpl">
31          <void property="filter">
32            <void property="comparison">
33              <object class="org.obo.filters.EqualsComparison"/>
34            </void>
35            <void property="criterion">
36              <object class="org.obo.filters.IDSearchCriterion"/>
37            </void>
38            <void property="value">
39              <string>OBO_REL:is_a</string>
40            </void>
41            </object>
42          </void>
43        </object>
44      </void>
45      <void method="add">
46        <object class="org.obo.filters.ObjectFilterImpl">
47          <void property="comparison">
48            <object idref="EqualsComparison0"/>
49          </void>
50          <void property="criterion">
51            <object idref="NameSynonymSearchCriterion0"/>
52          </void>
53          <void property="reasoner">
54            <object idref="RuleBasedReasoner0"/>
55          </void>
56        </object>

```

```

57      <void property="value">
58          <string>aniridia</string>
59      </void>
60      </object>
61      </void>
62      </void>
63  </object>
64 </java>
```

Il file xml come intuibile rappresenta la seguente query:

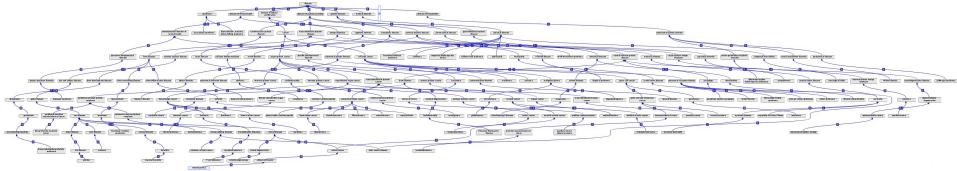


**Fig. 2.7:** Obo-Edit: Search panel

Ossia chiede che vengano restituiti tutti quei nodi che hanno tra i discendenti il nodo con nome e sinonimo “aniridia” o siano il nodo stesso. Anche se non propriamente strutturata per l’intellegibilità e modificabilità risulta chiaro l’utilizzo dei filtri e come sarebbe possibile aggiungere al file nuovi criteri di ricerca. Basterà dunque, attraverso un semplice script Java<sup>14</sup>, modificare il file in modo che, per ogni patologia tra quelle recuperate e collegate all’ontologia, venga introdotta una struttura analoga.

Il risultato costituirà il sottografo che si vorrebbe introdurre nel database, di cui si propone la rappresentazione grafica per coglierne l’ingenza: Costruito il sottografo di interesse viene esportata da OBO-Edit la lista completa dei nodi che lo costituiscono (nome e ID).

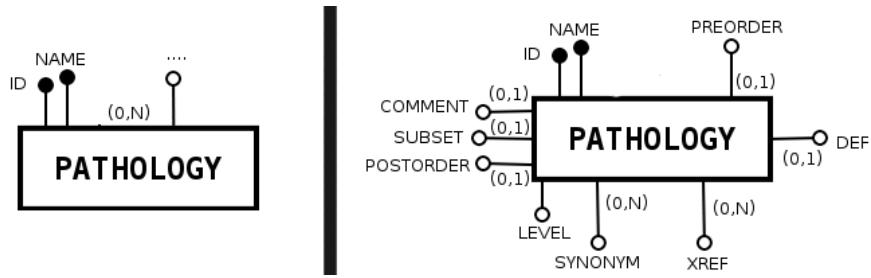
<sup>14</sup>Si veda *OboEditAllFilter.java* nell’alrchivio codici per approfondimenti



**Fig. 2.8:** Grafo della patologie estratte

### 2.3.2 Adattamento della struttura del Database

A seguito della scelta dell'ontologia è finalmente definibile con chiarezza la struttura di quanto lasciato precedentemente incompleto riguardo il modello logico e quindi l'implementazione del database, in particolare la definizione dell'entità pathology e dei suoi attributi. Di seguito la modifica del diagramma E/R ed del relativo progetto logico:



**Fig. 2.9:** Confronto migliorie

Oltre alle informazioni reperite dall'ontologia tradotte in attributi nel modello, si noti la modellazione della informazione gerarchica<sup>15</sup> espressa negli attributi: level, preorder e postorder.

**Perchè la scelta di questi attributi aggiuntivi?**

*“Given a node  $v$  [...] with  $\text{pre}(v)$  and  $\text{post}(v)$  ranks, the following properties are important towards our objectives: - all nodes  $x$  with  $\text{pre}(x) < \text{pre}(v)$  are the ancestors or preceding nodes of  $v$ ;*

- all nodes  $x$  with  $\text{pre}(x) > \text{pre}(v)$  are the descendants or following nodes of  $v$ ;*
- all nodes  $x$  with  $\text{post}(x) < \text{post}(v)$  are the descendants or preceding nodes of  $v$ ;*
- all nodes  $x$  with  $\text{post}(x) > \text{post}(v)$  are the ancestors or following nodes of  $v$ ;*

---

<sup>15</sup>Si rammenta che, per brevità, nel diagramma E/R l'entità pathology è direttamente riportata a seguito del collasso verso l'alto delle entità figlie.

- for any  $v$   $[..]$ , we have  $\text{pre}(v) - \text{post}(v) + \text{size}(v) = \text{level}(v)$ ;
- if  $\text{pre}(v) = 1$ ,  $v$  is the root, if  $\text{pre}(v) = n$ ,  $v$  is a leaf.  $[..]^{[20]}$

Ossia, per riuscire a sapere se una patologia  $y$  è una sottopatologia di una patologia  $x$  basterà verificare che l'attributo preorder di  $y$  sia maggiore di quello di  $x$  e che **contemporaneamente** l'attributo di postorder di  $y$  sia minore di quello di  $x$ . Tuttavia questa logica è applicabile solo sui nodi di un *albero*; Si spiegherà più in avanti in che modo nel nostro caso. Momentaneamente ci si cocentri sulla struttura del database la cui modifica del progetto logico, dunque, risulta:

```
.....
PATHOLOGY (ID, NAME, DEF, SUBSET, LEVEL, PREORDER,POSTORDER, COMMENT )
    AK: NAME

SYNONYM (ID, NAME)
    FK: ID REFERENCES PATHOLOGY

XREF (ID, REF)
    FK: ID REFERENCES PATHOLOGY
.....
```

Infine tradotto in SQL:

```
115 .....
116
117 CREATE TABLE PATHOLOGY
118 (
119     ID VARCHAR(15),
120     NAME VARCHAR(100) NOT NULL,
121     COMMENT VARCHAR(200),
122     DEFINITION VARCHAR(500),
123     SUBSET VARCHAR(100),
124     PREORDER INTEGER,
125     POSTORDER INTEGER,
126     LEVEL INTEGER,
127     PRIMARY KEY(ID)
128 );
129
130 CREATE TABLE SYNONYM
131 (
132     ID VARCHAR(15),
133     NAME VARCHAR(100),
134     PRIMARY KEY(ID, NAME),
135     FOREIGN KEY(ID) REFERENCES PATHOLOGY
```

```

136      ON DELETE CASCADE ON UPDATE CASCADE
137  );
138
139 CREATE TABLE XREF
140 (
141   ID VARCHAR(15),
142   REF VARCHAR(100),
143   PRIMARY KEY(ID, REF),
144   FOREIGN KEY(ID) REFERENCES PATHOLOGY
145     ON DELETE CASCADE ON UPDATE CASCADE
146 );
147 .....

```

### 2.3.3 Adattamento e creazione degli script Java

L’obiettivo successivo è quello di modificare gli script Java per la creazione delle INSERT SQL relativi alla nuova struttura dell’entità pathology pocanzi illustrata e crearne di nuovi. Nello specifico si dovrà:

1. Recuperare tutte le informazioni circa ogni singolo nodo del sub-graph, di cui attualmente si conosce solo nome ed id.
2. Strutturare le patologie a formare un *albero* per recuperare gli attributi di preorder, postorder e level, e successivamente inserire le stesse nel database.
3. Collegare ed inserire nel database le pathology indicate da biomart con particolare attenzione alla creazione, nel database, delle patologie di cui non si è trovato il corrispettivo nell’ontologia.

Effettuati questi punti si potrebbe dire concluso il procedimento di inserimento dell’ontologia. Si riportano i tratti salienti di ciascun punto:

1. Per il recupero delle informazioni circa una patologia viene creata una classe che prende in ingresso un id e cercando opportunamente nel file testuale costituente l’ontologia (la cui struttura è mostrata nel paragrafo sulla riduzione dell’ontologia) restituisce un oggetto pathology completo degli attributi di interesse. Si estrae di seguito il core principale<sup>16</sup>.

```

1 public class PathologyObj {
2
3   public String id="";
4   public String name="";
5   public String def="";
6   public String subset="";

```

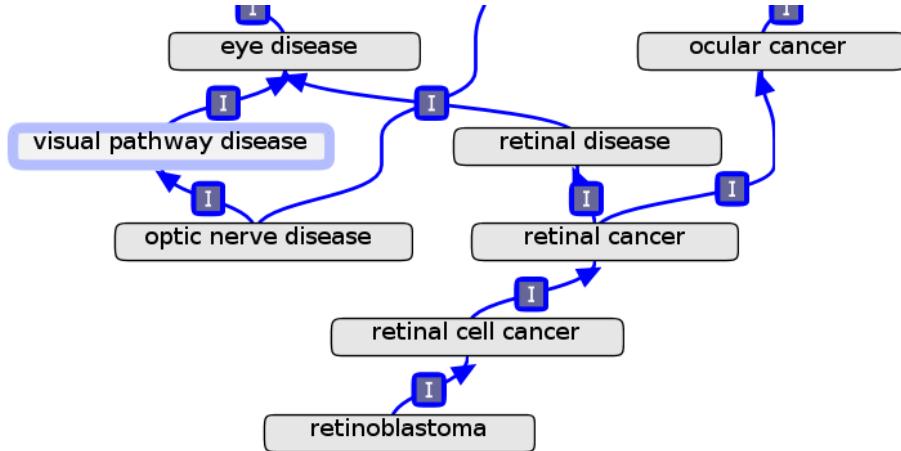
---

<sup>16</sup>Si rimanda al codice completo nell’appendice *PathologyObj.java*

```

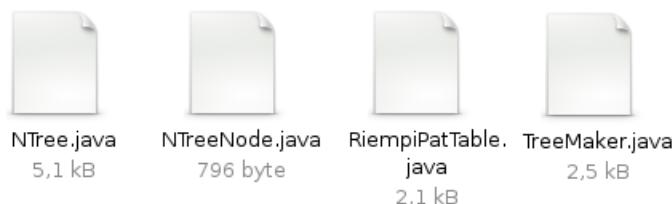
7  public String comment="";
8  public ArrayList<String> synonym = new ArrayList<String>();
9  public ArrayList<String> xref = new ArrayList<String>();
10 public String is_a="";
11 public int pre = 0;
12 public int post = 0;
13 public int level = 0;
14
15 public PathologyObj(String id) throws IOException{
16     BufferedReader br =
17         new BufferedReader( new FileReader("/home/HumanDO.txt"));
18     String strLine;
19     while( (strLine = br.readLine()) != null ){
20         if(strLine.startsWith("id: "+id)){
21             this.id = strLine;
22             while( !(strLine = br.readLine()).startsWith("[Term]")){
23                 if(strLine.startsWith("name:"))
24                     this.name = strLine;
25                 else if(strLine.startsWith("def:"))
26                     this.def = strLine;
27                 else if(strLine.startsWith("subset:"))
28                     this.subset = strLine;
29                 else if(strLine.startsWith("synonym:"))
30                     this.synonym.add(strLine);
31                 else if(strLine.startsWith("xref:"))
32                     this.xref.add(strLine);
33                 else if(strLine.startsWith("is_a:")) {
34                     StringTokenizer st;
35                     st = new StringTokenizer(strLine, "!");
36                     String tmp = st.nextToken();
37                     is_a = tmp.substring(0,tmp.length()-1);
38                 }
39             }
40             break;
41         }
42     }
43 }
```

2. Questo punto rappresenta la parte più laboriosa e delicata, anzitutto perché non si sta trattando con un albero bensì con un grafo. Poiché, però, solamente 2 dei 201 nodi d'interesse detiene più di una relazione di discendenza, si è deciso di poter introdurre una lieve ridondanza all'interno del database a fronte di un'indubbia semplificazione nella gestione delle informazioni gerarchiche considerando il grafo come un'albero. In questo caso, quindi, come è possibile verificare nello scorso di grafo riportato in figura, la ridondanza è circoscritta ad i due nodi con due padri ed ad i relativi sotto-grafi: in tutto solo 4 nodi vengono ripetuti.



**Fig. 2.10:** Nodi con due padri

Per recuperare gli attributi, in ultima analisi, sarà sufficiente creare un programma che utilizzi la classe PathologyObj per il recupero delle informazioni circa ogni patologia di nostro interesse e che costruisca l'albero per poi effettuare le *visite in preordine e postordine*<sup>17</sup>. Nella pratica sono stati creati i seguenti file:



**Fig. 2.11:** Script Java per la gestione dell'albero

- **NTree.java:** Costituisce la classe dell'albero di patologie con le relative funzioni di visita.
- **NTreeNode.java:** Costituisce la classe per la definizione degli oggetti che costituiranno i nodi dell'albero.
- **TreeMaker.java:** Costituisce il Main per l'utilizzo dell'albero definito nelle classi precedenti (Si rammenta anche l'utilizzo della classe PathologyObj definita in precedenza).

---

<sup>17</sup>Preorder e Postorder sono due diversi algoritmi di attraversamento e visita dei nodi di un'albero.

- **RiempPatTable.java** : Costituisce lo script per il riempimento, nel database, delle tabelle Pathology, Xref e Synonim.

Brevemente, si propongono alcuni scambi di codice di ciascuna classe per capirne meglio la funzione<sup>18</sup>:

Il file **NTreeNode.java** costituisce solo la struttura di ogni si nodo dell'albero memorizzandone l'elemento associato ma anche il riferimento al primo figlio ed al fratello.

```

1 public class NTreeNode
2 {
3     protected NTreeNode firstSon; //primo figlio
4     protected NTreeNode brother; //fratello
5     protected PathologyObj element; //elemento del nodo
6     public NTreeNode() //costruttore
7     {
8         this(null, null, null);
9     }
10    public NTreeNode(PathologyObj element, NTreeNode firstSon)
11    {
12        this(element, firstSon, null);
13    }
14
15    .....

```

Mentre, la classe descritta in **NTree.java** realizza l'albero con tutti i metodi indispensabili ad inserimenti, cancellazioni e modifiche oltre che le visite in preordine e postordine. Si evita di riproporre l'interezza del codice costituisce la normale gestione di un'albero e si ripropone solo l'implementazione dei metodi *set\_preorder* e *set\_postorder* che rappresentano l'unico elemento di novità assegnando a ciascun oggetto PathologyObj il relativo attributi di preordine e postordine ossia il numero progressivo di ordine di visita.

```

1 public class NTreeNode
2 {
3     public class NTree //classe pubblica albero
4     {
5         protected NTreeNode root; //radice dell'albero
6         public ArrayList<PathologyObj> array = new ArrayList<
7             PathologyObj>(); //per dopo.. lista
8         public NTree() //costruttore
9         {
10             root=null;
11         }

```

---

<sup>18</sup>Si rimanda all'**appendice C** per approfondimenti sul codice

```

11
12     public void setRoot(NTreeNode root){
13         this.root = root;
14     }
15
16     .....
17
18     public void set_preorder(){
19         set_preorder(root,1);
20     }
21     public int set_preorder(NTreeNode p, int num){
22         if(p!=null){
23             p.element.pre = num;
24             NTreeNode t=p.firstSon;
25             while(t!=null){
26                 num = set_preorder(t, ++num);
27                 t=t.brother;
28             }
29         }
30         return num;
31     }
32
33     public void set_postorder(){
34         set_postorder(root,1);
35     }
36     public int set_postorder(NTreeNode p, int num){
37         if(p!=null){
38             NTreeNode t=p.firstSon;
39             while(t!=null){
40                 num = set_postorder(t, num);
41                 num++;
42                 t=t.brother;
43             }
44             p.element.post = num;
45         }
46         return num;
47     }

```

Nel file **TreeMaker.java**, invece, non si fa altro che utilizzare le classi messe a disposizione in precedenza creando l'albero ed inserendo ciclicamente tutti le patologie (ogni ciclo aggiunge al minimo un livello) recuperandole dal file *allnodesfromobo.txt* che rappresenta l'insieme estratto precedentemente da Obo-Edit (Si rammenta che si è evitato di mostrare, ma è presente nello stesso file l'inserimento “manuale” delle patologie “duplicate” per la resa della gerarchia in forma d’albero):

```

1 public class TreeMaker {
2     NTree tree;
3
4     public TreeMaker() throws IOException {
5         BufferedReader br = new BufferedReader(
6             new FileReader("/home/allnodesfromobo.txt"));
7         ArrayList<PathologyObj> list =
8             new ArrayList<PathologyObj>();
9         StringTokenizer st = null;
10        String strLine,id;
11
12        //while per la creazione della lista di oggetti
13        while((strLine = br.readLine()) != null ){
14            st = new StringTokenizer(strLine, " ");
15            id = st.nextToken();
16            //System.out.println(id);
17            list.add(new PathologyObj(id));
18        }
19
20        //creiamo l'albero e colleghiamo i nodi
21        tree = new NTree();
22        for(int i=0; i<list.size(); i++){
23            if(list.get(i).id.equals("id: DOID:4")){
24                tree.setRoot(
25                    new NTreeNode(list.get(i),null,null));
26                list.remove(i);
27            }
28
29            while(!list.isEmpty()){
30                for(int i=0; i<list.size(); i++){
31                    //attenzione agli obsoleti! sono 3!
32                    if(list.get(i).is_a.equals(""))
33                        list.remove(i);
34                    String parent_id = "id: "+ list.get(i).is_a.substring(6);
35                    if(tree.insert(parent_id, new NTreeNode(list.get(i),null,
36                                null)))
37                        list.remove(i);
38
39                }
40            }
}

```

Infine, in **RiempipatTable.java**, di cui si riporta solo la parte fondamentale, dopo l'aggiunta “manuale” delle patologie obsolete (estranee alla gerarchia), si cicla sulla lista delle patologie recuperata tramite uno specifico metodo della classe NTree, ed dopo avere effettuato specifici controlli di tipo e formato circa i sigoli attributi si scrive il file SQL per l'inserimento nel database:

```

1
2     for(int i=0; i<list.size(); i++){
3         //System.out.println(list.size());
4         tmp = list.get(i);
5
6         .....
7
8         String sql = "INSERT INTO PATHOLOGY VALUES(
9             '"+tmp.id.substring(4)+"', "+
10            "'"+tmp.name.substring(6).replace("'", "\\'")+
11            "'"+comment+def+subset+pre+post+"','"+tmp.level+");";
12         out.println(sql);
13
14         for(int j=0; j<tmp.synonym.size(); j++){
15             sql = "INSERT INTO SYNONYM VALUES (
16                 '"+tmp.id.substring(4)+"', "+
17                 tmp.synonym.get(j).substring(9).replace("'", "\\'")+"');";
18             out.println(sql);
19         }
20
21         for(int j=0; j<tmp.xref.size(); j++){
22             sql = "INSERT INTO XREF VALUES (
23                 '"+tmp.id.substring(4)+"', "+
24                 tmp.xref.get(j).substring(6).replace("'", "\\'")+"');";
25             out.println(sql);
26         }
27
28     }

```

3. Per concludere, si attuano le modifice opportune al file **PathologyRecover.java** Introdotto nella sezione **Progettazione e realizzazione del DB**. Il codice, che, poichè concettualmente semplice, non viene mostrato, si occupa della già affrontata interrogazione del web service offerto da Biomart e nella scrittura, a differenza dello script originale, delle sole insert per la tabella RELATED\_TO<sup>19</sup>. Solo nel caso in cui si trattasse di patologie di cui non si è trovata corrispondenza nell'ontologia ma di cui comunque si vuol tener traccia (es Unclassifiable Pathology) allora si preoccupa anche di inserire il loro record nella tabella PATHOLOGY. Inoltre, da non dimenticare, vi è anche la gestione a parte dei nodi “replicati” per i quali ancora “manualmente” è stato effettuata l’aggiunta nella tabella RELATED\_TO. Si ritiene, dunque, conclusa la progettazione e sviluppo del database nella sua interezza e funzionalità. Nella prossima sezione verranno ampiamente discusse, invece, le modalità di accesso allo stesso mediante il web ed il loro relativo sviluppo.

---

<sup>19</sup>Si rimanda all’Appendice C per approfondimenti

## 2.4 Progettazione e realizzazione dell’interfaccia Web

Se da un lato lo sviluppo del database realizza e concretizza la ricerca efficiente dei dati, da non sottovalutare è l’importanza della componente funzionale ma soprattutto di agevolezza introdotta senz’altro da un’interfaccia grafica. Tra le tante e possibili interfacce che potrebbero essere accompagnate allo sviluppo di un database la più facile ed accessibile al giorno d’oggi è rappresentata da quella web. La scelta effettuata per il progetto, infatti, è ricaduta proprio su quest’ultima, in virtù della sua più semplice accessibilità remota mediante browser, dell’assente necessità di librerie o componenti esterni preinstallati, e poichè ricade su paradigmi noti e consolidati in termini di scelte grafiche e modalità di interazione con l’utente. Inoltre non vi sono particolari problematiche legate alla privacy del database e/o di sicurezza poichè in conflitto con lo scopo stesso del progetto: rendere possibile un più semplice ed agevole accesso ai dati e favorire l’interoperabilità.

### 2.4.1 Struttura e modalità di accesso

Per garantire un contesto all’interrogazione dei dati vera e propria, si è scelto di accennare la creazione un vero e proprio portale Web, per eventuali e possibili sviluppi futuri. Il portale sarà costituito da diverse sezioni elencate di seguito:

- **Home:** La sezione che offre la visione iniziale del portale con un’introduzione ed una lista di suggerimenti e nozioni fondamentali per l’utente più spaesato.
- **UC Data Mining:** La sezione che offre le diverse modalità di interrogazione ed un’introduzione sommaria sulla struttura del database. Il *core* del portale nonché centro degli sforzi ad interesse esclusivo del progetto di tesi.
- **Related Works:** La sezione che sottopone l’insieme navigabile dei lavori correlati circa la sequenze e gli elementi ultraconservati.
- **About Us:** La sezione che enuclea le informazioni circa i gruppi di ricerca coinvolti nel progetto, ed altre informazioni di reperibilità.

Nel proseguo dello scritto non si scenderà nel dettaglio per descrivere ogni singola pagina del portale ed i tecnicismi informatici grazie i quali è stato possibile ottenere certi risultati poichè non strettamente inerenti e rilevanti al progetto di tesi. Si approfondiscono, invece, la struttura e le modalità scelte per l’interrogazione nella sezione **UC Data Mining**. La sezione, come pocanzi accennato, introdurrà innanzitutto la struttura del database per rendere il più chiara possibile l’interazione col database e l’indagine più seria

e conscia, successivamente offrirà due gradi categorie di interrogazione: la prima, per gli utenti meno esperti, mediante una serie di *query prestrutturate*, la seconda secondo un'accesso completo al database mediante la possibilità di somministrare *query libere* al sistema. Per la prima categoria si è pensato di introdurre le seguenti *query prestrutturate* e di quindi altrettanti *forms*:

1. Possibilità di conoscere tutte le informazioni relative una sequeza uc specificata.
2. Possibilità di conoscere tutte le informazioni relative un gene specificato.
3. Possibilità di conoscere tutte sequenze uc rilevate su di uno specifico cromosoma e comprese tra diverse *bp* specificate.
4. Possibilità di conoscere tutte patologie correlate ad una sequeza uc e **tutti i suoi sottotipi** (Si notino le potenzialità introdotte dall'inserimento dell'ontologia).
5. Possibilità di confrontare secondo *matching approssimato*<sup>20</sup> una sequenza nucleotidica con quelle delle uc presenti nel database ed ottenere una lista di possibili risultati elencati per *scores*.
6. Possibilità di confrontare secondo *matching approssimato* una sequenza nucleotidica con quelle delle uc presenti nel database ed ottenere una lista di possibili risultati elencati per *scores* e filtrati secondo una patologia (ossia che siano necessariamente correlati a quella specifica patologia o ad un suo sottotipo).

Per la seconda categoria verrà invece introdotto un unico *form* la somministrazione delle query libere in SQL .

#### 2.4.2 Cenni sulla realizzazione

Per la realizzazione del progetto di accessibilità esposto nel paragrafo precedente, si rende opportuna la conoscenza e l'utilizzo del linguaggio **html** per la strutturazione dei contenuti, **css** per la gestione del *layout*, e del **php** per la gestione degli eventi lato server e la comunicazione con il database. Senza scendere troppo nel dettaglio si mettono in evidenza solo le strategie fondamentali per la realizzazione di alcuni dei punti elencati precedentemente.

Si prenda in considerazione, ad esempio, la possibilità di conoscere tutte patologie correlate ad una sequeza uc e tutti i suoi sottotipi. Sfruttando gli attributi delle patologie che rendono possibile la gerarchizzazione, come già spiegato nei paragrafi inerenti il collegamento della patologia (nell'esempio “disease”), si potrebbe giungere al seguente SQL:

---

<sup>20</sup>Per matching approssimato si intende con possibilità di errori

```

1
2 SELECT UC_NAME, NAME, MIN( LEVEL ) AS LEVEL
3 FROM UC
4 INNER JOIN
5 (
6
7   SELECT ENSAMBL_GENE_ID, NAME,
8   LEVEL FROM RELATED_TO
9   INNER JOIN
10  (
11
12    SELECT ID, NAME, LEVEL
13    FROM PATHOLOGY
14    WHERE PREORDER >=
15      (
16        SELECT PREORDER
17        FROM PATHOLOGY
18        WHERE NAME = 'disease'
19      )
20    AND POSTORDER <=
21      (
22        SELECT POSTORDER
23        FROM PATHOLOGY
24        WHERE NAME = 'disease'
25      )
26    ORDER BY LEVEL
27
28  ) AS B
29  WHERE RELATED_TO.ID = B.ID
30 ) AS C
31 WHERE UC.ENSAMBL_GENE_ID = C.ENSAMBL_GENE_ID
32 GROUP BY uc_name, name
33 ORDER BY LEVEL

```

In seguito alla somministrazione nel form da parte dell'utente del nome della patologia quindi, verrà inoltrata una richiesta *GET*<sup>21</sup> con i relativi parametri ad una pagina di risposta che li recupererà e connettendosi al database fornirà le risposte desiderate. Si prenda visione della sequenza indispensabile dei comandi in php presenti nella pagina di risposta:

```

1
2 <?php
3 //recupero del parametro pat presente nella richiesta GET
4 $pat = filter_input(INPUT_GET, "pat", FILTER_SANITIZE_STRING);
5

```

---

<sup>21</sup>Il metodo GET è un metodo del protocollo HTTP usato per ottenere il contenuto della risorsa indicata come URI

```

6   .... //codice di controllo, html etc..
7
8 //connessione al database
9 $con = mysql_connect("localhost", "client", "password");
10 if (!$con)
11 {
12   die('Could not connect: ' . mysql_error());
13 }
14 mysql_select_db("UCbase", \$con);
15
16 .... //codice di controllo, html etc..
17
18 //si stanzia ed esegue la query
19 $query = "SELECT UC_NAME,NAME,MIN(LEVEL)AS LEVEL
20 FROM UC
21 INNER JOIN (SELECT ENSAMBL_GENE_ID,NAME,LEVEL
22 FROM RELATED_TO INNER JOIN
23 (SELECT ID,NAME,LEVEL FROM PATHOLOGY WHERE PREORDER >=
24 (SELECT PREORDER FROM PATHOLOGY WHERE NAME='".$\$pat."')
25 AND POSTORDER
26 <=(SELECT POSTORDER FROM PATHOLOGY WHERE NAME='".$\$pat."')
27 ORDER BY LEVEL)AS B
28 WHERE RELATED_TO.ID = B.ID)AS C
29 WHERE UC.ENSAMBL_GENE_ID = C.ENSAMBL_GENE_ID
30 group by uc_name, name ORDER BY LEVEL";
31
32 $result = mysql_query($query);
33
34 //se ci sono risultati si crea la tabella
35 while($row = mysql_fetch_array($result, MYSQL_ASSOC))
36 {
37   if($i%2 == 0) echo "<tr class ='alt'>";
38   else echo "<tr>";
39
40   if($i==0){
41     foreach($row as $key => $value)
42     {
43       echo "<th>";
44       echo $key;
45       echo "</th>";
46     }
47     echo "</tr>";
48   }
49
50   if($i==0) echo "<tr>";
51   foreach($row as $x=>$x_value)
52   {
53     echo "<td>";
54     if($x == "UC_NAME"){

```

```

55     $uc_name = $x_value;
56     echo "<a href='result.php?uc_id=". 
57         $uc_name."'>".$x_value."</a>";
58 }
59 else if($x == "UC_NAME"){
60     $uc_name = $x_value;
61     echo $x_value;
62 }
63 else if($x == "DEFINITION" && $x_value != ""){
64     $x_value = preg_replace(
65         '!^(http|ftp|scp)(s)?:[\\/][a-zA-Z0-9.?_=~/]+!', 
66         "<a target='_blank' href=\"$\\1$\">$\\1</a>", $x_value);
67     echo $x_value;
68 }
69 else if($x == "SEQUENCE") {
70     echo strtolower(substr($x_value, 0, 13)."...");
71     echo "<a class='button' href='sequence.php?uc_id=". 
72         $uc_name."'> Get Sequence!</a>"; }
73 else
74     echo $x_value;
75 echo "</td>";
76 }
77 echo "</tr>";
78 $i++;
79 }
80 .... //codice finale, di controllo, html etc..
82 ?>

```

Con lo stesso meccanismo, potranno essere realizzati anche i primi 3 punti. Per quanto riguarda invece gli ultimi due, non si implemeterà un algoritmo di confronto ma si utilizzerà il collaudato e rinomato *BLAST*<sup>22</sup>[28] nello specifico la sua versione specializzata per il confronto di sequenze nucleotidiche *Blastn*.

Prima di ogni altra cosa bisogna creare ed indicizzare il piccolo database delle sequenze a partire dal file nel formato *FASTA*<sup>23</sup> usando il comando *makeblastdb* offerto sempre dal pacchetto scaricabile dal sito dell'NCBI [27].

```
./makeblastdb -in UC.fasta -dbtype nucl
-parse_seqids -out UCdb -title "Human UC"
```

---

<sup>22</sup>BLAST (Basic Local Alignment Search Tool, ovvero strumento di ricerca di allineamento locale) è un algoritmo usato per comparare le informazioni contenute nelle strutture biologiche primarie

<sup>23</sup>In bioinformatica, il formato FASTA è un formato testuale per la rappresentazione di sequenze nucleotidiche e peptidiche.

Successivamente si potrà utilizzare direttamente il comando blastn indicando il file contenente la query sempre nel formato FASTA ed il file di output:

```
./blastn -db UCdb -query query.fasta  
-task blastn-short -out results.out
```

Il file di output potrà ad esempio essere:

```
Database: Human UC  
          481 sequences; 126,007 total letters  
  
Query= prova  
  
Length=29  
Sequences producing significant alignments:  
Score      E  
           (Bits) Value  
  
lcl|uc.0          58.0   8e-12  
lcl|uc.23         20.3    1.8  
.....  
  
>lcl|uc.0  
Length=207  
  
Score = 58.0 bits (29),  Expect = 8e-12  
Identities = 29/29 (100%), Gaps = 0/29 (0%)  
Strand=Plus/Plus  
  
Query 1  TCCACCGACAATGACCAGTTAGTCCTCAT  29  
        |||||||  
Sbjct  1  TCCACCGACAATGACCAGTTAGTCCTCAT  29  
  
>lcl|uc.23  
Length=235  
  
Score = 20.3 bits (10),  Expect = 1.8  
Identities = 16/18 (89%), Gaps = 0/18 (0%)  
Strand=Plus/Minus  
  
Query 12  TGACCAGTTAGTCCTCAT  29  
        ||||| | |||||  
Sbjct  64  TGACCAGATGGTCCTCAT  47  
  
.....
```

Sia per la realizzazione del punto **4.** che del punto **5.** basterà recuperare per eventualmente poi rielaborare questi contenuti in php con una semplice riga di codice che inserisce nella variabile output un array di stringhe che costituiscono l'output visto in precedenza:

```
1 exec( "
2   cd /var/www/UCbase/ncbi-blast-2.2.27+/bin/;
3   ./blastn -db UCdb -query query.fasta
4   -task blastn-short 2>&1
5   ", $output
6 );
```

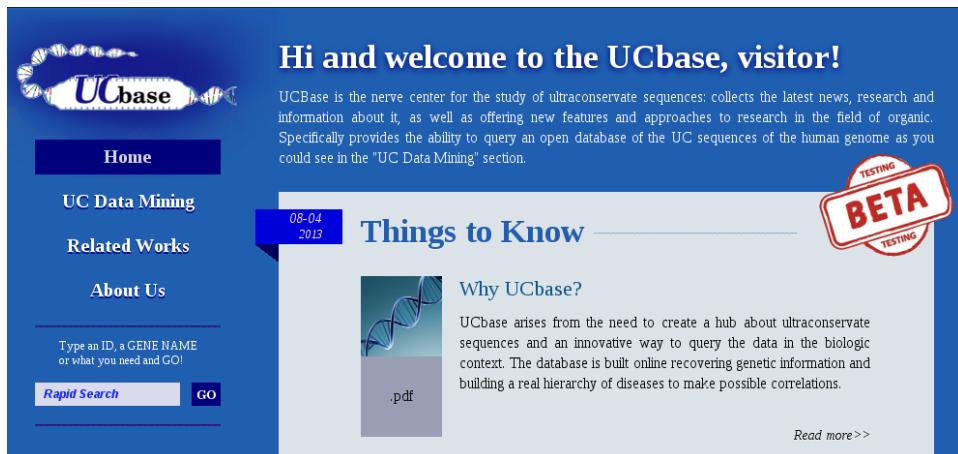
Si rammenta che, ovviamente, per il punto **5.** la rielaborazione dell'output non sarà solamente in termini grafici, ma si dovrà effettuare il filtraggio e l'ordinamento delle uc in funzione di quelle restituite da una query analoga a quanto visto per il risolvimento del punto **3.** : meccanicamente si eliminerà quelle sequenze uc non correlate a nessuna patologia e laddove presente la si indicherà; A parità di score (indicato da blastn), inoltre, la sequenza uc correlata alla patologia il cui livello è più affine a quello della patologia indicata dovrà risultare prioritaria nell'ordinamento dei risultati. (Data la linearità concettuale, anche in questo caso, si rimanda all'archivio dei codici per approfondimenti)



# Capitolo 3

## Risultati

Concluso e raffinato il codice per la definizione del layout di ogni pagina, ed a seguito di ulteriori migliorie ed aggiunte applicate al modello di accessibilità descritto nel capitolo precedente, si presentano i risultati raggiunti. Si tratta di un portale dal design semplice, leggero e molto asciutto. In primis la homepage, accessibile attualmente all'indirizzo  
<http://www.dsb.unimo.it/uibase>:



**Fig. 3.1:** UCbase: Pagina principale

Presenti in alto a sinistra i link fondamentali per il proseguo della navigazione. Si prenda visione delle altre pagine:

**Related works**

In this section you will find all information relating to the UC sequences even outside of our research laboratory. Paper and news about it will be posted below.

**Articles and Papers**

**UCbase & MiRfunc: A Database Of Ultraconserved Sequences And MicroRNA Function**

One hundred and eighty-one ultraconserved sequences (UCRs) longer than 200 bases were discovered in the genomes of human, mouse and rat. These are DNA sequences showing 100% identity among the three species...

*Read more >>*

**The BioMart Project**

BioMart is a unique open source data federation technology that provides unified access to distributed databases storing a wide range of data. This DATABASE issue recognizes BioMart's outstanding contributions to bioinformatics and documents the achievements of the BioMart community, which has grown impressively over the last ten years to become what it is today...

**Fig. 3.2:** UCbase: Related works

**About Us**

**INFORMATION SYSTEMS GROUP - UNIVERSITY OF MODENA AND REGGIO EMILIA**

The work of the ISGroup, here at the Computer Engineering Department (DII) of the University of Modena and Reggio Emilia, mainly focuses on the design and development of new systems, algorithms and data structures for the access and management of Information.

**JOIN OUR WORK!**

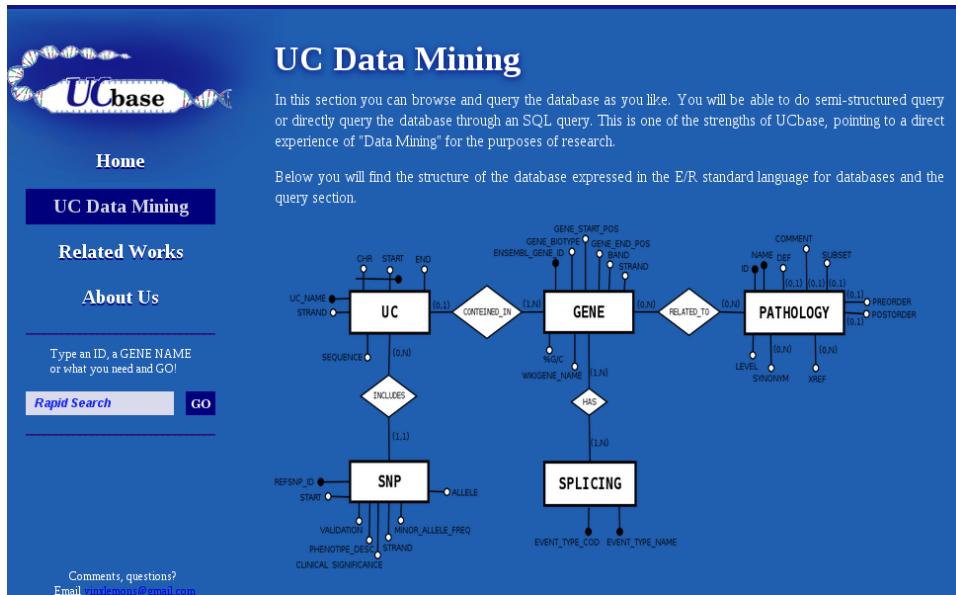
If you're experiencing issues and concerns about our work, please email us or help us to create a better work working together with our research group.

Design version: beta.  
Code version: beta.

**CONTACT US**  
**Telephone Nos. :** -----, -----  
**Email :** vinxelmons@gmail.com  
**Street Address :** ----- Modena, Italy

**Useful links**

**Fig. 3.3:** UCbase: About us



**Fig. 3.4:** UCbase: Uc data mining

### Pre-formed Query

All info about  ?

All info about  ?

All UC correled to  or its subtypes?

All UC in  starting between  and  ?

Blast sequence:

Blast sequence:  with pathology:

### Type your own Query!

**Fig. 3.5:** UCbase: Query

## Search Result

```
SELECT UC_NAME,NAME,MIN(LEVEL)AS LEVEL FROM UC INNER JOIN (SELECT  
ENSAMBL_GENE_ID,NAME,LEVEL FROM RELATED_TO INNER JOIN (SELECT  
ID,NAME,LEVEL FROM PATHOLOGY WHERE PREORDER >= (SELECT PREORDER  
FROM PATHOLOGY WHERE NAME='disease') AND POSTORDER <=(SELECT  
POSTORDER FROM PATHOLOGY WHERE NAME='disease') ORDER BY LEVEL)AS B  
WHERE RELATED_TO.ID = B.ID)AS C WHERE UC.ENSAMBL_GENE_ID =  
C.ENSAMBL_GENE_ID group by uc_name, name ORDER BY LEVEL
```

UC_NAME	NAME	LEVEL
<a href="#">uc.284</a>	disease of cellular proliferation	1
<a href="#">uc.76</a>	disease of cellular proliferation	1
<a href="#">uc.321</a>	disease of cellular proliferation	1
<a href="#">uc.323</a>	disease of cellular proliferation	1
<a href="#">uc.3</a>	disease of cellular proliferation	1
<a href="#">uc.63</a>	disease of cellular proliferation	1
<a href="#">uc.108</a>	disease of cellular proliferation	1
<a href="#">uc.144</a>	disease of cellular proliferation	1
<a href="#">uc.34</a>	disease of cellular proliferation	1

Fig. 3.6: UCbase: Risultati in formato tabellare

## Search Result

uc.0

<b>Uc_name</b>	uc.0
<b>Chr</b>	1
<b>Start</b>	10597697
<b>Ending</b>	10597903
<b>Strand</b>	1
<b>Sequence</b>	TCCACCGACAATGACCAGTT... <a href="#">Get Sequence!</a>
<b>Ensambl_gene_id</b>	ENSG00000142655

Fig. 3.7: UCbase: Risultati in formato schedario

**Blast: ACGTACAGTACGATC**

Database: Human UC (481 sequences; 126,007 total letters)

Click on the uc link to read specific Info.

Query= ACGTACAGTACGATC  
Length=15

Sequences alignments:	Score (Bits)	E Value	Score_pat (Level)	Pathology Name
lcl  <a href="#">uc.190</a>	18.3	2.2	1	disease of cellular proliferation
lcl  <a href="#">uc.58</a>	18.3	2.2	1	disease of cellular proliferation
lcl  <a href="#">uc.478</a>	16.4	8.7	1	disease of cellular proliferation
lcl  <a href="#">uc.456</a>	16.4	8.7	1	disease of cellular proliferation
lcl  <a href="#">uc.282</a>	16.4	8.7	1	disease of cellular proliferation
lcl  <a href="#">uc.151</a>	16.4	8.7	1	disease of cellular proliferation
lcl  <a href="#">uc.136</a>	16.4	8.7	1	disease of cellular proliferation

**Fig. 3.8:** UCbase: Risultati blastn

**Specific Info**

Click on the link for more info about the uc.

>lcl|[uc.190](#)  
Length=319

Score = 18.3 bits (9), Expect = 1.1  
Identities = 9/9 (100%), Gaps = 0/9 (0%)  
Strand=Plus/Plus

```
Query 3  GTACAGTAC  11
         |||||||||
Sbjct 260 GTACAGTAC  268
```

**Fig. 3.9:** UCbase: Informazioni specifiche

La pagina che raggruppa i contenuti relativi agli elementi ultraconservati *Related works* si presenta come mostrato in **Fig.3.2**. Segue lo stesso layout grafico anche la pagine *About Us* (**Fig.3.3**).

Si prenda visione, infine, della sezione più importante per la ricerca esplorativa sui dati: la pagina *UC Data Mining* (**Fig.3.4**). La pagina d'apprima illustra la struttura del database poi fornisce come una serie di form compilabili le opzioni di investigazione illustrate nel capitolo precedente (**Fig.3.5**). Sono stati implementate, poi, 3 grandi tipologie di layout per l'output dei dati: La prima in formato tabellare per rispondere ad una query generica come quelle provenienti dal form delle query libere oppure dal form 3. e 4. dal gruppo delle query prestrutturate. (**Fig.3.6**) La seconda, invece, in un formato un po' più da "schedario" per le informazioni circa i campi di un singolo record, quindi per rispondere alle query prestrutturate 1. e 2. (**Fig.3.7**)

La terza ed ultima categoria di visualizzazione utilizzate per le query prestrutturate 5. e 6., infine, in cui viene lasciato inalterato lo scheletro di output di blastn, ma aggiunta, come nel caso mostrato in figura (**Fig.3.8**), qualche informazione aggiuntiva.

Si rimanda nuovamente al sito <http://www.ds.unimo.it/ucbase> per scoprire ed approfondire appieno le potenzialità dell'interfaccia.

## Capitolo 4

# Conclusioni e sviluppi futuri

A conclusione dell’elaborato, successivamente alla realizzazione del database ed all’interfaccia web per la sua accessibilità, risulta quantomeno appropriato fornire sommarie conclusioni circa l’effettivo raggiungimento degli obiettivi che ci si era posti all’inizio del percorso di tesi, le difficoltà incontrate durante il percorso e consapevoli quanto umili commenti circa le riconosciute mancanze e grossolanità del progetto sicuramente migliorabili in futuro. In primis si premette che, a fronte di un risicato tempo di studio, si è realizzato un progetto dalle dimensioni rilevanti, le cui basi tematiche necessitano, perl’altro, di ulteriori e necessari approfondimenti. Per cui si ritenga l’elaborato frutto di un attento e laborioso lavoro di riassunto, che non lascia evincere le difficoltà con le quali ci si è scontrati, né i tentativi e scelte esplicative effettuate nelle più disparate direzioni e mai più realizzate. Tuttavia, la totalità degli obiettivi enunciati nel primo capitolo, sono stati portati a termine e dunque possono essere considerate risolte le problematiche che quest’ultimi erano tesi a dirimere. Nonostante questo, numerosi sono gli aspetti migliorabili del progetto di cui si è a consapevolezza e si rimanda in un futuro:

- **Il linguaggio java:** Una delle prime pecche risulta essere la scelta del linguaggio, per il quale si sarebbe dovuto prediligere un linguaggio per la prototipizzazione veloce come un linguaggio dinamico (ad esempio il Python).
- **Ridondanza nel database e linking semi-automatico delle patologie:** Un importante e nuovo livello di pulizia dei dati riguarderebbe la riprogettazione di un linking automatico con l’ontologia di riferimento (attualmente si ricorda che il processo richiede l’intervento di un esperto in materia per le patologie il cui nome non combacia direttamente con quello dell’ontologia) e di quella per evitare la ridondanza all’interno del database (introdotta attualmente per semplificare il processo di gestione delle gerarchie attraverso l’uso un’albero anzichè di un grafo)

- **Impossibilità di aggiornamento dei dati:** Attualmente il progetto non prevede la possibilità di aggiornare dati già esistenti ma solo di creare gli script indispensabili alla creazione ed all'inserimento di dati in un database ex-novo. Si suggerisce anche questo punto come una delle possibili migliorie effettuabili.
- **Semplicità dell'interfaccia Web:** Infinitamente migliorabile, infine, l'interfaccia web, sia in termini di design sia in termini di funzionalità offerte.

Sottolineate le mancanze e possibili migliorie, non rimane che riavvalorare, in antitesi, quanto, invece, risulta essere un punto di forza del progetto:

- **Open Database e linking dell'HDO:** Nell'ambito degli elementi e delle sequenze nucleotidiche ultraconservate, il servizio offerto per l'interrogazione del database, costituisce uno dei pochissimi esempi di libero accesso ed investigazione in ambito biologico e rappresenta un'avanguardia sia per modalità di interrogazione mediante SQL, sia per la struttura stessa del database che collega direttamente numerose informazioni (comprese quelle di tipo **gerarchico**) sulle patologie correlate ad un gene e quindi ad alcuni elementi ultraconservati.
- **L'interfaccia Web:** Offrire il sopracitato servizio mediante il Web risulta un'ulteriore agevolazione rivolta all'utenza, attraverso un'interfaccia semplice e veloce, funzionale ed efficiente, che non comporta nessun prerequisito o difficoltà iniziale.
- **Recupero automatico dei dati:** Con poche e semplici operazioni (e pochi file di supporto) è possibile, anche per i poco esperti, ottenere informazioni **aggiornate** sulle metainformazioni inerenti le sequenze ultraconservate e, quindi, costruire il database, senza il minimo sforzo ed in tempi brevi.

## Appendice A

# Glossario di base di biologia molecolare

Di seguito vengono riportate alcune informazioni utili (citazioni e definizioni) inerenti.

**Biologia Molecolare:** “ *La biologia molecolare è una branca della biologia che studia gli esseri viventi a livello dei meccanismi molecolari alla base della loro fisiologia, concentrandosi in particolare sulle interazioni tra le macromolecole, ovvero proteine e acidi nucleici (DNA e RNA). Per biologia molecolare si intendono spesso una serie di tecniche che consentono la rilevazione, l'analisi, la manipolazione, l'amplificazione (PCR) e la copia (clonaggio) degli acidi nucleici.* ” [23]

**DNA e RNA:** “ *L'acido deossiribonucleico (DNA) costituisce, con l'acido ribonucleico (RNA), la classe di polimeri informazionali definita 'acidi nucleici', componenti fondamentali delle strutture viventi. Tra tutte le molecole biologiche soltanto gli acidi nucleici possiedono la potenzialità di autoduplicazione che permette la replicazione e la trascrizione dell'informazione chimica in essi contenuta. Altri polimeri costituenti fondamentali delle strutture biologiche, quali le proteine e i polisaccaridi, sono privi di tali capacità e, pur dotati di comparabile complessità polimerica strutturale e informazionale, di ricchissima informazione chimica e di grande flessibilità funzionale, non sono in grado di rivestire un ruolo genetico. Questo ruolo è legato alla particolare organizzazione strutturale degli acidi nucleici e alle proprietà chimiche, fisiche e topologiche sia dei suoi costituenti che delle macromolecole considerate nel loro insieme.* ” [24]

**Cromosoma:** “ *Nome dato da W. Waldeyer nel 1888 ai piccoli corpi intensamente colorabili, in genere di forma bastoncellare, visibili nel nucleo della cellula durante la mitosi. Secondo la teoria cromosomica , dimostrata*

*da T.H. Morgan con ricerche su Drosophila melanogaster, i cromosomi sono le strutture essenziali dell'eredità. ”* [24]

**Banda Cromosomica:** “ *La denaturazione o la digestione enzimatica della cromatina, seguite dall'incorporazione di un colorante specifico per il DNA, fa sì che nei cromosomi mitotici di organismi complessi si individui una serie di bande alternantesi in senso trasversale fino a un totale di 400, 550 o 850 visibili al microscopio ottico. Le si ottengono sottponendo i cromosomi a digestione con tripsina prima della colorazione con Giemsa; si formano bande scure (bande G) e chiare (bande G negative). ”* [24]

**Base azotata:** “ *In biochimica, per base azotata si intende una delle cinque basi che compongono i nucleotidi del DNA e dell'RNA. Si distinguono in purine e pirimidine. Purine: Adenina e Guanina Pirimidine: Citosina, Timina e Uracile Nel DNA si trovano entrambe le purine e, tra le pirimidine, citosina e timina. Nell'RNA la timina è sostituita dall'uracile. Nel DNA le basi si accoppiano a due a due con legami a idrogeno, mentre nell'RNA, essendo questo una catena singola, non sono legate tra loro. DNA: adenina-timina e citosina-guanina. RNA: adenina-uracile e citosina-guanina. In chimica, una base azotata è un qualsiasi composto che manifesta proprietà basiche per via della presenza di un doppietto elettronico non condiviso su un atomo di azoto, ad esempio l'ammoniaca e le ammine. ”* [23]

**Sequenza Nucleotidica:** *In biologia molecolare, la successione ordinata di nucleotidi del DNA. Ogni nucleotide è formato da una molecola di desossiribosio, un gruppo fosfato e una base azotata. In particolare, sono le basi azotate a determinare la sequenza nucleotidica specifica dell'informazione genetica.”* [24]

**Gene:** “ *Nome proposto da W. Johannsen (1909) per indicare l'unità ereditaria scoperta da G. Mendel; è materializzato nella cellula da un segmento di una molecola di DNA, in cui l'ordine di successione dei quattro tipi di nucleotidi determina, secondo la legge di corrispondenza espressa dal codice genetico, l'ordine di successione dei diversi tipi di amminoacido nel corrispondente polipeptide. Mentre nelle cellule procariotiche ogni gene è espresso da un segmento continuo di DNA, di regola i geni delle cellule eucariotiche sono in pezzi: segmenti di DNA portatori di informazione (esoni) si alternano a segmenti non informazionali (introni); la sintesi del polipeptide avviene in tal caso mediante la trascrizione di un tratto di molecola di DNA che comprende tutti gli esoni e i relativi introni. La lunga molecola di RNA così sintetizzata subisce quindi un processo di maturazione, che consiste nell'escissione degli introni, con la formazione della molecola di mRNA necessaria per l'ordinario processo di traduzione dell'informazione.[...] ”* [24]

**Base Pair:** “Le coppie di basi o paia di basi (abbreviate come pb o, dall’inglese base pair, bp o bps) sono comunemente utilizzate come misura della lunghezza fisica di sequenze di acidi nucleici a doppio filamento. Spesso, data la grande dimensione dei genomi, viene utilizzata l’unità kbp (k sta per kilo), pari a mille paia di basi. Il numero di paia di basi rappresenta il numero di coppie di basi azotate che contiene il filamento in analisi. Ogni paio di basi è tenuto insieme attraverso legami idrogeno che si instaurano tra le molecole, ricche di azoto. ” [23]

**SNP:** “Un polimorfismo a singolo nucleotide (spesso definito in inglese Single Nucleotide Polymorphism o SNP, pronunciato snip) è un polimorfismo, cioè una variazione, del materiale genico a carico di un unico nucleotide, tale per cui l’allele polimorfico risulta presente nella popolazione in una proporzione superiore all’1%. Al di sotto di tale soglia si è soliti parlare di mutazione. Ad esempio, se le sequenze individuate in due pazienti sono AAGCCTA e AAGCTTA, è presente uno SNP che differenzia i due alleli C e T. ” [23]

**Splicing:** “In biologia molecolare e in genetica, splicing è una modificazione del nascente pre-mRNA che avviene insieme o dopo la trascrizione (biologia), nella quale gli introni sono rimossi e gli esoni vengono uniti. La cosa è necessaria per il tipico RNA messaggero prima che possa essere usato per produrre una corretta proteina tramite la traduzione (biologia) o sintesi proteica. ” [23]

**Allele:** “In genetica si definisce allele o fattore ogni variante di sequenza di un gene o di un locus genico (generalmente un gruppo di geni). Il genotipo di un individuo relativamente ad un gene è il corredo di alleli che egli si trova a possedere. In un organismo diploide, in cui sono presenti due copie di ogni cromosoma (i cosiddetti cromosomi omologhi), il genotipo è dunque costituito da due alleli. Se un organismo possiede per lo stesso gene due alleli identici allora si definisce omozigote per quel locus genico, mentre se sono differenti viene detto eterozigote. In realtà per dati geni si possono riscontrare diverse varianti alleliche tanto che per molti geni noti è stata messa in discussione la terminologia allele mutante perché le varietà alleliche non necessariamente portano ad uno svantaggio selettivo. ” [23]

**Strand:** “In genetica, ciascuna delle due catene di una molecola di DNA, normalmente costituita da una coppia di filamenti fra loro complementari. Nel processo di sintesi del DNA si chiama filamento guida (leading strand) quello che è sintetizzato in maniera continua, in direzione 5’3’, e filamento tardivo (lagging strand) quello sintetizzato in direzione opposta 3’5’, in maniera discontinua sotto forma di brevi frammenti successivamente legati da legami covalenti. ” [24]

## Appendice B

# Elementi ultraconservati nel genoma umano

Di seguito vengono riportate alcune informazioni utili (citazioni e definizioni) inerenti.

### Elemento ultraconservato:

*Un elemento ultraconservato (UCE: ultraconserved element) è una regione di DNA che è identica in almeno due specie. Uno dei primi studi sugli UCES mostra che sequenze della lunghezza di 200 nucleotidi o più del DNA umano si sono conservate interamente (sequenze identiche di acidi nucleici) in entrambi ratti e topi. Rispetto alla maggioranza di DNA non codificante, alcuni elementi ultraconservati sono stati rilevati come attivi a livello trascrizionale, originando molecole di RNA non codificante. [23]*

### Estratto tradotto dall'articolo originale Ultraconserved Elements in the Human Genome [22] :

Ci sono 481 segmenti più lunghi di 200 bp che sono conservati (100% di identità senza inserimenti né cancellazioni) tra regioni ortologhe del genoma umano, del topo e del ratto. Quasi tutti questi segmenti sono anche *ultraconservati* nel genoma del pollo e del cane con un 95% e 99% di identità, rispettivamente. Molti sono anche significativamente conservati nel pesce. Questi elementi **ultraconservati** del genoma umano sono spesso situati o sovrapposti agli esoni nei geni coinvolti nel processamento dell'RNA o in introni o vicino a geni coinvolti nella regolazione della trascrizione e sviluppo. Insieme a più di 5.000 sequenze di oltre 100bp che sono conservative tra i tre mammiferi sequenziati, questi rappresentano una classe di elementi genetici le cui funzioni ed origini evolutive sono ancora da determinare, ma che sono più altamente conservate tra queste specie che le proteine, e sembrano essere essenziali per l'ontogenesi dei mammiferi e altri vertebrati.

Anche se solo circa l'1,2% del genoma umano appare codificare proteine, è

stato stimato che circa il 5% è più conservato di quanto previsto dall'evoluzione naturale dopo la divisione con i roditori. Studi hanno evidenziato segmenti non codificanti specifici nel genoma umano che sembrano essere in fase di selezione, con una conservazione del 70% o 80% di identità con il topo su più di 100bp. Uno studio di questi elementi sul cromosoma 21 dell'uomo ha scoperto che quelli che erano altamente conservati in più specie contenevano un numero significante di elementi non codificanti. Risultati simili sono stati trovati confrontando l'uomo, il topo ed il ratto in uno studio sulla regione cromosomica contenente il gene CFTR (1.8 Mbp) e in uno studio funzionale sul locus SIM2 in molti mammiferi.

# Appendice C

## Archivio dei Codici

Di seguito, vengono riportati, in ordine di apparizione, tutti i codici di cui, a volte, si è mostratosolo gli scorci più importanti. Si rammenta che per i codici modificati nel corso dell'elaborato verrà riportata solo l'ultima versione.

### C.1 Uc\_biomaRt\_pathology2.java

```
1 library("biomaRt")
2 ensembl = useMart("ensembl",dataset="hsapiens_gene_ensembl")
3
4 options(width=120)
5 attributesR=c('ensembl_gene_id','mim_morbid_description','pathology',
6   percentage_gc_content')
7 filtersR = c('chromosome_name','start','end')
8
9 out=list<-list()
10 for (i in 1:length(uc2009[,1])){
11   valuesR<-list(uc2009[i,2],uc2009[i,3],uc2009[i,4])
12   out<-getBM(attributes=attributesR,filters=filtersR,values=valuesR,
13   mart=ensembl,uniqueRows=F)
14   if (nrow(out)!=0){
15     out<-data.frame(cbind(uc2009[i,1],out))
16     colnames(out)[1]<-"uc"
17   }
18   else{
19     out.names<-colnames(out)
20     nulli<-as.data.frame(matrix(rep(NA, ncol(out)),nrow=1))
21     out<-data.frame(cbind(uc2009[i,1],nulli))
22     colnames(out)<-c("uc",out.names)
23   }
24   if(i==1)
25     write.table(out,"pathology.csv",append=FALSE,col.names=TRUE,
26       quote=FALSE,row.names=FALSE,sep=";")
27   else
28     write.table(out,"pathology.csv",append=TRUE,col.names=FALSE,
29       quote=FALSE,row.names=FALSE,sep=";")
```

```

27     out.list[[i]]<-out
28     print(uc2009[i,1])
29     print(out)
30 }

```

## C.2 Tabelle.sql

```

1  CREATE TABLE GENE
2  (
3      ENSAMBL_GENE_ID VARCHAR(20),
4      GENE_BYOTYPE VARCHAR(30),
5      GENE_START_POS INTEGER,
6      GENE_END_POS INTEGER,
7      BAND VARCHAR(10),
8      STRAND INTEGER,
9      WIKIGENE_NAME VARCHAR(30),
10     G_C_PERC VARCHAR(10),
11     PRIMARY KEY(ENSAMBL_GENE_ID)
12 );
13
14 CREATE TABLE UC
15 (
16     UC_NAME VARCHAR(10) UNIQUE NOT NULL,
17     CHR VARCHAR(5),
18     START INTEGER,
19     ENDING INTEGER,
20     STRAND INTEGER,
21     SEQUENCE VARCHAR(1000) NOT NULL,
22     ENSAMBL_GENE_ID VARCHAR(20),
23     PRIMARY KEY (CHR,START),
24     FOREIGN KEY (ENSAMBL_GENE_ID) REFERENCES GENE
25         ON DELETE SET NULL ON UPDATE CASCADE
26 );
27
28 CREATE TABLE SNP
29 (
30     REFSNP_ID VARCHAR(30),
31     START INTEGER,
32     ALLELE VARCHAR(10),
33     VALIDATION VARCHAR(30),
34     MINOR_ALLELE_FREQ FLOAT(10),
35     PHENOTYPE_DESC VARCHAR(30),
36     CLINICAL_SIGNIFIANCE VARCHAR(30),
37     STRAND INTEGER,
38     CHR VARCHAR(5),
39     START_UC INTEGER,
40     PRIMARY KEY (REFSNP_ID),
41     FOREIGN KEY(CHR, START_UC) REFERENCES UC (CHR,START)
42         ON DELETE CASCADE ON UPDATE CASCADE
43 );
44
45
46 CREATE TABLE SPLICING

```

```

47  (
48    EVENT_TYPE_COD VARCHAR(8),
49    EVENT_TYPE_NAME VARCHAR(30) UNIQUE NOT NULL,
50    PRIMARY KEY(EVENT_TYPE_COD)
51  );
52
53 CREATE TABLE HAS
54  (
55    EVENT_TYPE_COD VARCHAR(8),
56    ENSAMBL_GENE_ID VARCHAR(20),
57    PRIMARY KEY(EVENT_TYPE_COD, ENSAMBL_GENE_ID),
58    FOREIGN KEY(EVENT_TYPE_COD) REFERENCES SPlicing
59      ON DELETE CASCADE ON UPDATE CASCADE,
60    FOREIGN KEY(ENSAMBL_GENE_ID) REFERENCES GENE
61      ON DELETE CASCADE ON UPDATE CASCADE
62  );
63
64 CREATE TABLE PATHOLOGY
65  (
66    ID VARCHAR(15),
67    NAME VARCHAR(100) NOT NULL,
68    COMMENT VARCHAR(200),
69    DEFINITION VARCHAR(500),
70    SUBSET VARCHAR(100),
71    PREORDER INTEGER,
72    POSTORDER INTEGER,
73    LEVEL INTEGER,
74    PRIMARY KEY(ID)
75  );
76
77 CREATE TABLE RELATED_TO
78  (
79    ID VARCHAR(15),
80    ENSAMBL_GENE_ID VARCHAR(20),
81    PRIMARY KEY(ID, ENSAMBL_GENE_ID),
82    FOREIGN KEY(ID) REFERENCES PATHOLOGY
83      ON DELETE CASCADE ON UPDATE CASCADE,
84    FOREIGN KEY(ENSAMBL_GENE_ID) REFERENCES GENE
85      ON DELETE CASCADE ON UPDATE CASCADE
86  );
87
88 CREATE TABLE SYNONYM
89  (
90    ID VARCHAR(15),
91    NAME VARCHAR(100),
92    PRIMARY KEY(ID, NAME),
93    FOREIGN KEY(ID) REFERENCES PATHOLOGY
94      ON DELETE CASCADE ON UPDATE CASCADE
95  );
96
97 CREATE TABLE XREF
98  (
99    ID VARCHAR(15),
100

```

```

101    REF VARCHAR(100),
102    PRIMARY KEY(ID, REF),
103    FOREIGN KEY(ID) REFERENCES PATHOLOGY
104        ON DELETE CASCADE ON UPDATE CASCADE
105 );

```

### C.3 BiomartMain.java

```

1 import java.net.*;
2 import java.io.*;
3
4 public class BiomartMain {
5     public static void main(String[] args) throws Exception {
6
7         GeneRecover gene = new GeneRecover();
8         SnpRecover.snp = new SnpRecover();
9         SplicingRecover spli = new SplicingRecover();
10        PathologyRecover path = new PathologyRecover();
11    }
12
13 }

```

### C.4 GeneRecover.java

```

1 import java.io.FileOutputStream;
2 import java.io.FileReader;
3 import java.io.InputStream;
4 import java.io.InputStreamReader;
5 import java.io.PrintStream;
6 import java.net.URL;
7 import java.net.URLEncoder;
8 import java.util.StringTokenizer;
9
10 public class GeneRecover {
11     //dati per recupero info uc
12     private final int num_lines = 481;
13     private String[] chr = new String[num_lines];
14     private String[] bp_start = new String[num_lines];
15     private String[] bp_end = new String[num_lines];
16     private String[] strand_uc = new String[num_lines];
17     private String[] sequence = new String[num_lines];
18     private String[] uc_name = new String[num_lines];
19     //dati per tabella
20     private String ensambl_gene_id, gene_byotipe, gene_start_pos,
21             gene_end_pos, band, strand, wikigene_name, g_c_perc, sql_gene,
22             sql_uc;
23
24     public GeneRecover()throws Exception{
25         //leggo dal file le info
26         try
27             {

```

//tsv file containing data

```

28     String strFile_uc = "/home/uc_coordinates_hg19_reference_file.
29             csv";
30     String strFile_seq = "/home/ultraconservate/UC.txt";
31
32     //create BufferedReader to read tsv file e per scrivere i file
33             sql
34     //delle tabelle "UC" e "GENE".
35     BufferedReader br = new BufferedReader( new FileReader(
36             strFile_uc));
37     BufferedReader br2 = new BufferedReader( new FileReader(
38             strFile_seq));
39     FileOutputStream sql_file = new FileOutputStream("/home/kira/
40             Scrivania/tesi/uc_gene_insert.sql");
41     PrintStream out = new PrintStream(sql_file);
42
43     String strLine = "";
44     //recupero sequenze uc
45     int i = 1, j=0, k=0;
46     while(i <= num_lines*2 ){
47         strLine = br2.readLine();
48         if(i%2 == 0) {sequence[j] = strLine; j++;}
49         else {uc_name[k] = strLine.substring(1); k++;}
50         i++;
51     }
52
53     StringTokenizer st = null;
54     int lineNumber = 0;
55     br.readLine(); //la prima riga contenente meta-info viene
56             scartata cosÃ¢n
57     //read comma separated file line by line
58     while( (strLine = br.readLine()) != null )
59     {
60         //break comma separated line using ","
61         st = new StringTokenizer(strLine, ",");
62         String uc = st.nextToken();
63         chr[lineNumber] = st.nextToken();
64         //togliamo fastidiosi apici
65         chr[lineNumber] = chr[lineNumber].replace("\",\"");
66         bp_start[lineNumber] = st.nextToken();
67         bp_end[lineNumber] = st.nextToken();
68         strand_uc[lineNumber] = st.nextToken();
69         System.out.println("Sto recuperando "+uc+": "+chr[lineNumber]
70             +" "+bp_start[lineNumber]+ " "+bp_end[lineNumber]);
71
72
73     String myxml = "<Query virtualSchemaName = \"default\""
74             formatter = \"CSV\" header = \"1\" "
75             +"uniqueRows = \"0\" count = \"\" datasetConfigVersion =
76                 \"0.6\" >"
77             +"<Dataset name = \"hsapiens_gene_ensembl\" interface = \""
78                 default\" >"
79             +"<Filter name = \"chromosome_name\" value = \"\"+chr[
80                 lineNumber]+\"\"/>"
```

```

71      +"<Filter name = \"start\" value = \\""+bp_start[lineNumber]+"
72          \"/>"
73      +"<Filter name = \"end\" value = \\""+bp_end[lineNumber]+"\\"/>
74          "
75      +"<Attribute name = \"wikigene_name\" />"
76      +"<Attribute name = \"ensembl_gene_id\" />"
77      +"<Attribute name = \"gene_biotype\" />"
78      +"<Attribute name = \"start_position\" />"
79      +"<Attribute name = \"end_position\" />"
80      +"<Attribute name = \"band\" />"
81      +"<Attribute name = \"strand\" />"
82      +"<Attribute name = \"percentage_gc_content\" />"
83      +"</Dataset>"
84      +"</Query>";
85      //System.out.println(myxml);
86      String encoded = URLEncoder.encode(myxml, "utf-8");
87      //System.out.println(encoded.length());
88      URL url = new URL("http://www.biomart.org/biomart/
89                      martservice?query="+encoded);
90      InputStream response = url.openStream();
91      BufferedReader reader = new BufferedReader(new
92          InputStreamReader(response));
93      i=0;
94      for (String line; ;i++) { //nel caso non ci sia nessuna
95          risposta devo comunque inserire l'uc
96          if((line = reader.readLine()) == null && (i == 1)) {
97              sql_uc = "INSERT IGNORE INTO UC VALUES('"+uc_name[
98                  lineNumber]+','"+ chr[lineNumber]+",'"+ bp_start[
99                  lineNumber]+
100                 ","+ bp_end[lineNumber]+","+strand_uc[lineNumber]+",'"+
101                 +sequence[lineNumber]+",'null');";
102              System.out.println(sql_uc);
103              out.println(sql_uc);
104              break;
105          }
106          else if((i!=0) && (line != null)){//se Ã¹ la prima riga ->
107              scarto, altrimenti scrivo su file sql
108              line.replace(",",","+null+",");
109              if(line.charAt(0) == ',') line = "null"+ line;
110              System.out.println(line);
111              st = new StringTokenizer(line, ",");
112              wikigene_name = st.nextToken();
113              if(!wikigene_name.equals("null")) wikigene_name = ""+
114                  wikigene_name+"";
115              ensambl_gene_id = st.nextToken();
116              if(!ensambl_gene_id.equals("null")) ensambl_gene_id = ""+
117                  +ensambl_gene_id+"";
118              gene_byotipe = st.nextToken();
119              if(!gene_byotipe.equals("null")) gene_byotipe = ""+
120                  gene_byotipe+"";
121              gene_start_pos = st.nextToken();
122              gene_end_pos = st.nextToken();
123              band = st.nextToken();

```

```

113         if(!band.equals("null")) band = ""+band+"";
114         strand = st.nextToken();
115         if(!strand.equals("null")) strand = ""+strand+"";
116         g_c_perc = st.nextToken();
117
118         sql_gene = "INSERT IGNORE INTO GENE VALUES("+
119             ensambl_gene_id+", "+gene_byotype+", "+
120             gene_start_pos+", "+gene_end_pos+", "+band+", "+strand+", "+
121             wikigene_name+", "+
122             +g_c_perc+");";
123         System.out.println(sql_gene);
124         out.println(sql_gene);
125
126         sql_uc = "INSERT IGNORE INTO UC VALUES('"+uc_name[
127             lineNumber]+"' ,'"'+ chr[lineNumber]+"' ,'"'+ bp_start[
128             lineNumber]+
129             ", '"+ bp_end[lineNumber]+", "+strand_uc[lineNumber]+", '"
130             +sequence[lineNumber]+"' ,"
131             + ensambl_gene_id +");";
132         System.out.println(sql_uc);
133         out.println(sql_uc);
134     }
135     reader.close();
136     lineNumber++;
137   }
138   catch(Exception e)
139   {
140     System.out.println("Exception: " + e);
141   }
142 }
143
144 }
```

## C.5 SplicingRecover.java

```

1 import java.io.BufferedReader;
2 import java.io.FileOutputStream;
3 import java.io.FileReader;
4 import java.io.InputStream;
5 import java.io.InputStreamReader;
6 import java.io.PrintStream;
7 import java.net.URL;
8 import java.net.URLEncoder;
9 import java.util.StringTokenizer;
10
11 public class SplicingRecover {
12   private final int num_lines = 481;
13   private String[] chr = new String[num_lines];
```

```

14 private String[] bp_start = new String[num_lines];
15 private String[] bp_end = new String[num_lines];
16
17 private String splicing_cod, splicing_name, ensambl_gene_id, sql_has,
18     sql_splicing;
19
20 public SplicingRecover()throws Exception{
21     //leggo dal file le info
22     try
23     {
24         //tsv file containing data
25         String strFile = "/home/kira/Scrivania/tesi/
26             uc_coordinates_hg19_reference_file.csv";
27         //sql output file
28         FileOutputStream sql_file = new FileOutputStream("/home/
29             splicing_has_insert.sql");
30         PrintStream out = new PrintStream(sql_file);
31
32         //create BufferedReader to read tsv file
33         BufferedReader br = new BufferedReader( new FileReader(strFile
34             ));
35         String strLine = "";
36         StringTokenizer st = null;
37         int lineNumber = 0;
38         br.readLine(); //la prima riga contenente meta-info viene
39             scartata cosÃš
40         //read comma separated file line by line
41         while( (strLine = br.readLine()) != null )
42         {
43             //break comma separated line using ","
44             st = new StringTokenizer(strLine, ",");
45             String uc = st.nextToken();
46             chr[lineNumber] = st.nextToken();
47             chr[lineNumber] = chr[lineNumber].replace("\\", "");//"
48                 togliamo fastidiosi apici
49             bp_start[lineNumber] = st.nextToken();
50             bp_end[lineNumber] = st.nextToken();
51             System.out.println("Sto recuperando "+uc+": "+chr[lineNumber]
52                 +" "+bp_start[lineNumber]+" "+bp_end[lineNumber]);
53
54             String myxml = "<Query virtualSchemaName = \"default\" "
55                 formatter = "\"CSV\" header = \"1\" "
56                 +"uniqueRows = \"0\" count = \"\" datasetConfigVersion =
57                     \"0.6\" >"
58                 +"<Dataset name = \"hsapiens_gene_ensembl\" interface = \""
59                     default\" >"
60                 +"<Filter name = \"chromosome_name\" value = \""
61                     +chr[
62                         lineNumber]+"/>"
63                 +"<Filter name = \"start\" value = \""
64                     +bp_start[lineNumber]+"/>"
65                 +"<Filter name = \"end\" value = \""
66                     +bp_end[lineNumber]+"/>"
67
68                 +"<Attribute name = \"ensembl_gene_id\" />"
```

```

55      +"<Attribute name = \"splicing_event_type\" />"  

56      +"<Attribute name = \"name_106\" />"  

57      +"</Dataset>"  

58      +"</Query>"  

59      //System.out.println(myxml);  

60      String encoded = URLEncoder.encode(myxml, "utf-8");  

61      //System.out.println(encoded.length());  

62      URL url = new URL("http://www.biomart.org/biomart/  

       martservice?query="+encoded);  

63      InputStream response = url.openStream();  

64      BufferedReader reader = new BufferedReader(new  

       InputStreamReader(response));  

65      int i=0;  

66      for (String line; (line = reader.readLine()) != null;) {  

67          if(i!=0){//se Ã¹ la prima riga -> scarto  

68              while(line != (line = line.replace(", , ,null,")));  

69              line = line.replace(",,", "\r\n");  

70              if(line.charAt(0) == ',') line = "null"+ line;  

71              if(line.charAt(line.length()-1) == ',') line = line + "  

               null";  

72              st = new StringTokenizer(line, ",");  

73              System.out.println(line);  

74  

75              ensambl_gene_id = st.nextToken();  

76              if(!ensambl_gene_id.equals("null")) ensambl_gene_id = ""+  

               ensambl_gene_id+",";  

77              splicing_cod = st.nextToken();  

78              if(!splicing_cod.equals("null")) splicing_cod = ""+  

               splicing_cod+",";  

79              splicing_name = st.nextToken();  

80              if(!splicing_name.equals("null")) splicing_name = ""+  

               splicing_name+",";  

81  

82              //se sono nulli non li inserisco per nulla  

83              if(!splicing_cod.equals("null") && !splicing_name.equals(""  

               null")) ){  

84                  sql_splicing = "INSERT INTO SPlicing VALUES("+  

                   splicing_cod+","+splicing_name+") ON DUPLICATE KEY  

                   UPDATE EVENT_TYPE_COD = EVENT_TYPE_COD;"  

85                  System.out.println(sql_splicing);  

86                  out.println(sql_splicing);  

87  

88                  sql_has = "INSERT INTO HAS VALUES("+splicing_cod+", "+  

                   ensambl_gene_id+") ON DUPLICATE KEY UPDATE  

                   EVENT_TYPE_COD = EVENT_TYPE_COD;"  

89                  System.out.println(sql_has);  

90                  out.println(sql_has);  

91              }  

92          }  

93          i++;  

94      }  

95      lineNumber++;  

96      reader.close();  

97  }

```

```

98     }
99     catch(Exception e)
100    {
101        System.out.println("Exception while reading csv file: " + e);
102    }
103 }
104
105 }
```

## C.6 SnpRecover.java

```

1 import java.io.BufferedReader;
2 import java.io.FileOutputStream;
3 import java.io.FileReader;
4 import java.io.InputStream;
5 import java.io.InputStreamReader;
6 import java.io.PrintStream;
7 import java.net.URL;
8 import java.net.URLEncoder;
9 import java.util.StringTokenizer;
10
11
12 public class SnpRecover {
13     private final int num_lines = 481;
14     private String[] chr = new String[num_lines];
15     private String[] bp_start = new String[num_lines];
16     private String[] bp_end = new String[num_lines];
17
18     private String refsnp_id, allele, start, strand, clinical, phenotype,
19         validated, allele_freq, sql;
20
21     public SnpRecover()throws Exception{
22
23         try
24         {
25             //tsv file containing data
26             String strFile = "/home/uc_coordinates_hg19_reference_file.csv";
27             //per output sql
28             FileOutputStream sql_file = new FileOutputStream("/home/
29                         snp_insert.sql");
30             PrintStream out = new PrintStream(sql_file);
31
32             //create BufferedReader to read tsv file
33             BufferedReader br = new BufferedReader( new FileReader(strFile
34                                         ));
35             String strLine = "";
36             StringTokenizer st = null;
37             int lineNumber = 0;
38             br.readLine(); //la prima riga contenente meta-info viene
39                         scartata cosÃ
40             //read comma separated file line by line
41             while( (strLine = br.readLine()) != null )
42             {
```

```

39 //break comma separated line using ","
40 st = new StringTokenizer(strLine, ",");
41 String uc = st.nextToken();
42 chr[lineNumber] = st.nextToken();
43 chr[lineNumber] = chr[lineNumber].replace("\", \"");//  

   togliamo fastidiosi apici
44 bp_start[lineNumber] = st.nextToken();
45 bp_end[lineNumber] = st.nextToken();
46 System.out.println("Sto recuperando "+uc+": "+chr[lineNumber]
    +" "+bp_start[lineNumber]+ " "+bp_end[lineNumber]);
47
48
49 String myxml = "<Query virtualSchemaName = \"default\""
   formatter = \"CSV\" header = \"1\" "
50 +"uniqueRows = \"0\" count = \"\" datasetConfigVersion =
   \"0.6\" >
51 +<Dataset name = \"hsapiens_snp\" interface = \"default\" >
52 +<Filter name = \"chr_name\" value = \""+chr[lineNumber]+"
   \"/>
53 +<Filter name = \"chrom_start\" value = \""+bp_start[
   lineNumber]+"\"/>
54 +<Filter name = \"chrom_end\" value = \""+bp_end[lineNumber]
   +"\"/>
55 +<Attribute name = \"refsnip_id\" />
56 +<Attribute name = \"allele\" />
57 +<Attribute name = \"chrom_start\" />
58 +<Attribute name = \"chrom_strand\" />
59 +<Attribute name = \"clinical_significance\" />
60 +<Attribute name = \"phenotype_description\" />
61 +<Attribute name = \"validated\" />
62 +<Attribute name = \"minor_allele_freq\" />
63 +</Dataset>
64 +</Query>";
65 //System.out.println(myxml);
66 String encoded = URLEncoder.encode(myxml, "utf-8");
67 //System.out.println(encoded.length());
68 URL url = new URL("http://www.biomart.org/biomart/
   martservice?query="+encoded);
69 InputStream response = url.openStream();
70 BufferedReader reader = new BufferedReader(new
   InputStreamReader(response));
71 int i=0;
72 for (String line; (line = reader.readLine()) != null;) {
73 //se Ã¹ la prima riga -> scarto
74 if(i!=0) {
75     int first = line.indexOf("\");
76     int last = line.lastIndexOf("\");
77     if(first != -1){
78         String substring = line.substring(first,last);
79         String substring_replaced = substring.replace(", ", "-")
           ;
80         line = line.replace(substring, substring_replaced);
81         line = line.replace("\", \"");
82     }

```

```

83     while(line != (line = line.replace(", , ,null,")));
84     if(line.charAt(0) == ',') line = "null"+ line;
85     if(line.charAt(line.length()-1) == ',') line = line + "
86         null";
87     st = new StringTokenizer(line, ",");
88     System.out.println(line);
89
90     refsnp_id = st.nextToken();
91     if(!refsnp_id.equals("null")) refsnp_id = ""+refsnp_id+","
92         ;
93     allele = st.nextToken();
94     if(!allele.equals("null")) allele = "" + allele + ",";
95     start = st.nextToken();
96     if(!start.equals("null")) start = "" + start + ",";
97     strand = st.nextToken();
98     clinical = st.nextToken();
99     if(!clinical.equals("null")) clinical = ""+clinical+","
100        ;
101    phenotype = st.nextToken();
102    if(!phenotype.equals("null")) phenotype = ""+phenotype+","
103        ;
104    validated = st.nextToken();
105    if(!validated.equals("null")) validated = ""+validated+","
106        ;
107    allele_freq = st.nextToken();
108
109    sql = "INSERT IGNORE INTO SNP VALUES("+ refsnp_id +","
110        +
111        start+","
112        +allele+","
113        +validated+
114        ,
115        +allele_freq+","
116        +phenotype+","
117        +clinical+","
118        +strand+","
119        +
120        +chr[lineNumber]+","
121        +
122        ,bp_start[lineNumber]+");"
123
124
125    reader.close();
126    lineNumber++;
127
128 }
129 catch(Exception e)
130 {
131     System.out.println("Exception while reading csv file: " + e);
132 }
133
134 }
```

## C.7 PathologyRecover.java

```

1 import java.io.BufferedReader;
2 import java.io.FileOutputStream;
```

```

3 import java.io.FileReader;
4 import java.io.InputStream;
5 import java.io.InputStreamReader;
6 import java.io.PrintStream;
7 import java.net.URL;
8 import java.net.URLEncoder;
9 import java.util.StringTokenizer;
10
11 public class PathologyRecover {
12     private final int num_lines = 481;
13     private String[] chr = new String[num_lines];
14     private String[] bp_start = new String[num_lines];
15     private String[] bp_end = new String[num_lines];
16
17     public PathologyRecover()throws Exception{
18         //leggo dal file le info
19         //tsv file containing data
20         String strFile = "/home/uc_coordinates_hg19_reference_file.csv";
21
22         //reader per controllare id patologia
23         BufferedReader read = new BufferedReader( new FileReader(
24             "/home/MIM+PATHOLOGY_PULITO_ORDINATO+cri.csv"));
25
26         //file sql
27         FileOutputStream sql_file = new FileOutputStream("/home/
28             pat_related2_insert.sql");
29         PrintStream out = new PrintStream(sql_file);
30
31         //create BufferedReader to read tsv file
32         BufferedReader br = new BufferedReader( new FileReader(strFile
33             ));
34         String strLine = "";
35         StringTokenizer st = null;
36         int lineNumber = 0;
37         br.readLine(); //la prima riga contenente meta-info viene
38             scartata cosÃ¬
39         //read comma separated file line by line
40         while( (strLine = br.readLine()) != null )
41         {
42             //break comma separated line using ","
43             st = new StringTokenizer(strLine, ",");
44             String uc = st.nextToken();
45             chr[lineNumber] = st.nextToken();
46             chr[lineNumber] = chr[lineNumber].replace("\", \"");//"
47                 togliamo fastidiosi apici
48             bp_start[lineNumber] = st.nextToken();
49             bp_end[lineNumber] = st.nextToken();
50             System.out.println("Sto recuperando "+uc+": "+chr[lineNumber]
51                 +" "+bp_start[lineNumber]+ " "+bp_end[lineNumber]);
52
53             String myxml = "<Query virtualSchemaName = \"default\""
54                 formatter = \"TSV\" header = \"1\" "
55                 +"uniqueRows = \"0\" count = \"\" datasetConfigVersion =
56                 \"0.6\" >"
```

```

48      + "<Dataset name = \"hsapiens_gene_ensembl\" interface = \""
49      + "default\" >"
50      + "<Filter name = \"chromosome_name\" value = \"\""+chr[
51          lineNumber]+"/>"
52      + "<Filter name = \"start\" value = \"\""+bp_start[lineNumber]+"
53          +"/>"
54      + "<Filter name = \"end\" value = \"\""+bp_end[lineNumber]+"/>
55          "
56      + "<Attribute name = \"ensembl_gene_id\" />"
57      + "<Attribute name = \"mim_morbid_description\" />"
58      + "<Attribute name = \"pathology\" />"
59      + "</Dataset>"
60      + "</Query>";
61      //System.out.println(myxml);
62      String encoded = URLEncoder.encode(myxml, "utf-8");
63      //System.out.println(encoded.length());
64      URL url = new URL("http://www.biomart.org/biomart/
65          martservice?query="+encoded);
66      InputStream response = url.openStream();
67      BufferedReader reader = new BufferedReader(new
68          InputStreamReader(response));
69
70      String line = reader.readLine(); //scarto indice colonne
71
72      while ( (line = reader.readLine()) != null) { //per ogni
73          riga che mi restituiscono
74          st = new StringTokenizer(line, " ");
75          String gene_id = st.nextToken();
76          String pat="", mim="";
77          if(st.hasMoreTokens())
78              mim = st.nextToken();
79          if(st.hasMoreTokens())
80              pat = st.nextToken();
81
82          //qui scopri id patologia o inventalo
83
84          if(mim.contains(";")){
85              st = new StringTokenizer(mim, ";");
86              mim = st.nextToken();
87          }
88          if(pat.contains(";")){
89              st = new StringTokenizer(pat, ";");
90              pat = st.nextToken();
91          }
92
93          if(!mim.equals("") && mim.startsWith(" "))
94              mim = mim.substring(1); //il famoso spazio iniziale

```

```

95         while ( (line = read.readLine()) != null){
96             if(!pat.equals("")){
97                 if (line.toLowerCase().startsWith(pat)){
98                     st = new StringTokenizer(line, ";");
99                     st.nextToken();
100                    id_pat = st.nextToken();
101                    System.out.println("trovato "+pat+" id:"+id_pat);
102                    break;
103                }
104            }
105            else
106                break;
107        }
108        //resetto read
109        //riporto all'inizio file
110        read =new BufferedReader( new FileReader("/home/MIM+
111                                     PATHOLOGY_PULITO_ORDINATO+cri.csv"));
112
113        while ( (line = read.readLine()) != null){
114            if(!mim.equals("")){
115                if (line.toLowerCase().startsWith(mim)){
116                    st = new StringTokenizer(line, ";");
117                    st.nextToken();
118                    id_mim = st.nextToken();
119                    System.out.println("trovato "+mim+" id:"+id_mim);
120                    break;
121                }
122            }
123            else
124                break;
125        }
126        //resetto read
127        //riporto all'inizio file
128        read =new BufferedReader( new FileReader("/home/MIM+
129                                     PATHOLOGY_PULITO_ORDINATO+cri.csv"));
130
131        if(mim.contains(",")){
132            mim = mim.replace(",", "\\" );
133        }
134        if(pat.contains(",")){
135            pat = pat.replace(",", "\\" );
136
137            if(!pat.equals("") && !id_pat.startsWith("DOID")){
138                String sql1 = "INSERT IGNORE INTO PATHOLOGY VALUES('"+
139                                id_pat+"','"+"+pat+
140                                "','null,null,null,null,null,-1');";
141                out.println(sql1);
142                System.out.println(sql1);
143            }
144            if(!mim.equals("") && !id_mim.startsWith("DOID")){
145                String sql1 = "INSERT IGNORE INTO PATHOLOGY VALUES('"+
146                                id_mim+"','"+"+mim+
147                                "','null,null,null,null,null,-1');";

```

```

145         out.println(sql1);
146         System.out.println(sql1);
147     }
148     //inserimento finale
149     if(!pat.equals("")){
150         String sql = "INSERT IGNORE INTO RELATED_TO VALUES('"+
151             id_pat+"','"+gene_id+"');";
152         out.println(sql);
153         System.out.println(sql);
154
155         if(id_pat.equals("DOID:1891") || id_pat.equals("DOID:4645"
156             ) || id_pat.equals("DOID:771") || id_pat.equals("DOID
157             :768") ){
158             String sql_ = "INSERT IGNORE INTO RELATED_TO VALUES('"+
159                 id_pat+"bis','"+gene_id+"');";
160             out.println(sql_);
161             break;
162         }
163     }
164     if(!mim.equals("")){
165         String sql = "INSERT IGNORE INTO RELATED_TO VALUES('"+
166             id_mim+"','"+gene_id+"');";
167         out.println(sql);
168         System.out.println(sql);
169
170         if(id_mim.equals("DOID:1891") || id_mim.equals("DOID:4645"
171             ) || id_mim.equals("DOID:771") || id_mim.equals("DOID
172             :768") ){
173             String sql_ = "INSERT IGNORE INTO RELATED_TO VALUES('"+
174                 id_mim+"bis','"+gene_id+"');";
175             out.println(sql_);
176         }
177     }

```

## C.8 TestCorrispondenze.java

```

1 import java.io.BufferedReader;
2 import java.io.FileNotFoundException;
3 import java.io.FileOutputStream;
4 import java.io.FileReader;
5 import java.io.PrintStream;
6 import java.util.ArrayList;
7
8
9 public class TestCorrispondenze {

```

```

10
11 //TESTA LE CORRISPONDENZE DIRETTE TRA LE PATHOLOGIE BIOMART E QUELLE DELLA
12 // ONTOLOGIA
13 public static void main(String[] args) throws Exception {
14
15     extractSourceFileFromDO();
16     String strFile = "/home/kira/Scrivania/AllpathologyfromDO.txt";
17     BufferedReader br = new BufferedReader( new FileReader(strFile));
18     String strFile2 = "/home/kira/Scrivania/pat_mim.txt";
19     BufferedReader br2 = new BufferedReader( new FileReader(strFile2));
20     FileOutputStream file = new FileOutputStream("/home/result.txt");
21     PrintStream out = new PrintStream(file);
22     FileOutputStream file2 = new FileOutputStream("/home/suggests.txt");
23     PrintStream out_sugg = new PrintStream(file2);
24
25     ArrayList<String> do_name = new ArrayList<String>();
26     ArrayList<String> patmim_name = new ArrayList<String>();
27     String line = "";
28
29     //riempio le patologie dall'ontologia
30     while( (line = br.readLine()) != null ){
31         do_name.add(line.toLowerCase());
32     }
33     //riempio le patologie di biomart
34     while( (line = br2.readLine()) != null ){
35         patmim_name.add(line.toLowerCase());
36     }
37
38     for(int i=0; i< patmim_name.size(); i++){
39         if(do_name.contains(patumim_name.get(i))){
40             System.out.println(patumim_name.get(i));
41             out.println(patumim_name.get(i));
42         }
43         else{
44             int ris;
45             for(int k=0; k<do_name.size(); k++){
46                 ris = computeLevenshteinDistance(do_name.get(k), patmim_name.
47                     get(i));
48                 if(ris < patmim_name.get(i).length()/2)//se le differenze non sono
49                     piÙ della metà dei caratteri
50                     out_sugg.println(patumim_name.get(i) + " somiglia a: " + do_name.
51                     get(k) );
52             }
53         }
54     }
55
56     private static void extractSourceFileFromDO() throws Exception{
57         String strFile = "/home/kira/Scrivania/HumanDO.txt";
58         BufferedReader br = new BufferedReader( new FileReader(strFile));
59         FileOutputStream file = new FileOutputStream("/home/AllpathologyfromDO
60             .txt");
61         PrintStream out = new PrintStream(file);

```

```

59     ArrayList<String> do_name = new ArrayList<String>();
60     String line = "";
61     while( (line = br.readLine()) != null ){
62         if(line.startsWith("name:")) || line.startsWith("synonym:")){
63             if(line.startsWith("name:"))
64                 line = line.substring(6);
65             else{
66                 line = line.substring(9);
67                 int firstI, secondI;
68                 firstI=line.indexOf("\\");
69                 secondI=line.lastIndexOf("\\");
70                 line = line.substring(firstI+1, secondI);
71             }
72             if(!do_name.contains(line)) do_name.add(line);
73         }
74         else if(line.startsWith("[Typedef]"))
75             break;
76     }
77
78     //stampa sul file i nomi puliti senza doppioni
79     for(int i=0; i< do_name.size(); i++)
80         out.println(do_name.get(i));
81
82 }
83
84 private static int minimum(int a, int b, int c) {
85     return Math.min(Math.min(a, b), c);
86 }
87
88 public static int computeLevenshteinDistance(CharSequence str1,
89         CharSequence str2) {
90     int[][] distance = new int[str1.length() + 1][str2.length() + 1];
91
92     for (int i = 0; i <= str1.length(); i++)
93         distance[i][0] = i;
94     for (int j = 1; j <= str2.length(); j++)
95         distance[0][j] = j;
96
97     for (int i = 1; i <= str1.length(); i++)
98         for (int j = 1; j <= str2.length(); j++)
99             distance[i][j] = minimum(
100                 distance[i - 1][j] + 1,
101                 distance[i][j - 1] + 1,
102                 distance[i - 1][j - 1]
103                     + ((str1.charAt(i - 1) ==
104                         str2.charAt(j - 1))
105                         ? 0
106                         : 1));
107
108     return distance[str1.length()][str2.length()];
109 }
```

## C.9 OboEditAllFilter.java

```
1 import java.io.BufferedReader;
2 import java.io.FileNotFoundException;
3 import java.io.FileOutputStream;
4 import java.io.FileReader;
5 import java.io.PrintStream;
6 import java.util.StringTokenizer;
7
8
9 public class OboEditAllFilter {
10     //classe per creare il file dei filtri per obo-edit per ottenere tutti i
11     //nodi
12     //dell'albero che ci interessa, anche quelli intermedi
13
14     public static void main(String[] args) throws Exception {
15         String srcFile = "/home/MIM+PATHOLOGY_PULITO_ORDINATO+cri.csv";
16         BufferedReader br = new BufferedReader( new FileReader(srcFile));
17         FileOutputStream out_file = new FileOutputStream("/home/OBOAllfilter.
18             xml");
19         PrintStream out = new PrintStream(out_file);
20
21         //parte iniziale
22         out.println("<?xml version=\"1.0\" encoding=\"UTF-8\"?>" +
23             "<java version=\"1.6.0_18\" class=\"java.beans.XMLDecoder\">" +
24             "<object class=\"org.obo.filters.CompoundFilterImpl\">" +
25             "<void property=\"booleanOperation\">" +
26             "<int>1</int>" +
27             "</void>" +
28             "<void property=\"filters\">");
```

```
29             StringTokenizer st = null;
30             String strLine = "";
31             while( (strLine = br.readLine()) != null ){
32                 st = new StringTokenizer(strLine, ";");
33                 String pat_name = st.nextToken();
34                 String id = st.nextToken();
35                 if(id.startsWith("DOID")){
36                     out.println("<void method=\"add\">" +
37                         "<object class=\"org.obo.filters.ObjectFilterImpl\">" +
38                         "<void property=\"comparison\">" +
39                         "<object id=\"EqualsComparison0\" class=\"org.obo.filters.
40                             EqualsComparison\"/>" +
41                         "</void>" +
42                         "<void property=\"criterion\">" +
43                         "<object id=\"IDSearchCriterion0\" class=\"org.obo.filters.
44                             IDSearchCriterion\"/>" +
45                         "</void>" +
46                         "<void property=\"reasoner\">" +
47                         "<object id=\"RuleBasedReasoner0\" class=\"org.obo.reasoner
48                             .rbr.RuleBasedReasoner\">" +
49                         "<void property=\"properties\">" +
50                         "<object class=\"java.util.HashSet\"/>" +
```

```

47         "</void>" +
48         "</object>" +
49         "</void>" +
50         "<void property=\"value\">" +
51         "<string>" + id + "</string>" +
52         "</void>" +
53         "</object>" +
54         "</void>" +
55         //parte per richiedere padri raggiungibili con is_a
56         "<void method=\"add\">" +
57         "<object class=\"org.obo.filters.ObjectFilterImpl\">" +
58         "<void property=\"aspect\">" +
59         "<object id=\"DescendantSearchAspect0\" class=\"org.obo.
               filters.DescendantSearchAspect\"/>" +
60         "</void>" +
61         "<void property=\"comparison\">" +
62         "<object id=\"EqualsComparison0\" class=\"org.obo.filters.
               EqualsComparison\"/>" +
63         "</void>" +
64         "<void property=\"criterion\">" +
65         "<object id=\"IDSearchCriterion0\" class=\"org.obo.filters.
               IDSearchCriterion\"/>" +
66         "</void>" +
67         "<void property=\"reasoner\">" +
68         "<object id=\"RuleBasedReasoner0\" class=\"org.obo.reasoner.
               rbr.RuleBasedReasoner\"/>" +
69         "<void property=\"properties\">" +
70         "<object class=\"java.util.HashSet\"/>" +
71         "</void>" +
72         "</object>" +
73         "</void>" +
74         "<void property=\"traversalFilter\">" +
75         "<object class=\"org.obo.filters.LinkFilterImpl\">" +
76         "<void property=\"filter\">" +
77         "<void property=\"comparison\">" +
78         "<object class=\"org.obo.filters.EqualsComparison\"/>" +
79         "</void>" +
80         "<void property=\"criterion\">" +
81         "<object class=\"org.obo.filters.IDSearchCriterion\"/>" +
82         "</void>" +
83         "<void property=\"value\">" +
84         "<string>OBO_REL:is_a</string>" +
85         "</void>" +
86         "</void>" +
87         "</object>" +
88         "</void>" +
89         "<void property=\"value\">" +
90         "<string>" + id + "</string>" +
91         "</void>" +
92         "</object>" +
93         "</void>);
```

94

95 }else if(id.startsWith("SI")){

96 out.println("<void method=\"add\">" +

```

97      "<object class=\"org.obo.filters.ObjectFilterImpl\">" +
98      "<void property=\"comparison\">" +
99      "<object id=\"EqualsComparison0\" class=\"org.obo.filters.
100         EqualsComparison\"/>" +
101     "</void>" +
102     "<void property=\"criterion\">" +
103     "<object class=\"org.obo.filters.NameSynonymSearchCriterion
104         \"/>" +
105     "</void>" +
106     "<void property=\"reasoner\">" +
107     "<object id=\"RuleBasedReasoner0\" class=\"org.obo.reasoner
108         .rbr.RuleBasedReasoner\">" +
109     "<void property=\"properties\">" +
110     "<object class=\"java.util.HashSet\"/>" +
111     "</void>" +
112     "</object>" +
113     "</void>" +
114     "<void property=\"value\">" +
115     "<string>" + pat_name + "</string>" +
116     "</void>" +
117     "</object>" +
118     "</void>" +
119     //parte per richiedere padri raggiungibili con is_a
120     "<void method=\"add\">" +
121     "<object class=\"org.obo.filters.ObjectFilterImpl\">" +
122     "<void property=\"aspect\">" +
123     "<object id=\"DescendantSearchAspect0\" class=\"org.obo.
124         filters.DescendantSearchAspect\"/>" +
125     "</void>" +
126     "<void property=\"comparison\">" +
127     "<object id=\"EqualsComparison0\" class=\"org.obo.filters.
128         EqualsComparison\"/>" +
129     "</void>" +
130     "<void property=\"criterion\">" +
131     "<object id=\"NameSynonymSearchCriterion0\" class=\"org.obo.
132         filters.NameSynonymSearchCriterion\"/>" +
133     "</void>" +
134     "<void property=\"reasoner\">" +
135     "<object id=\"RuleBasedReasoner0\" class=\"org.obo.reasoner.
136         rbr.RuleBasedReasoner\"/>" +
137     "<void property=\"properties\">" +
138     "<object class=\"java.util.HashSet\"/>" +
139     "</void>" +
140     "<void property=\"traversalFilter\">" +
141     "<object class=\"org.obo.filters.LinkFilterImpl\">" +
142     "<void property=\"filter\">" +
143     "<void property=\"comparison\">" +

```

```

144         "<void property=\"value\">" +
145         "<string>OBO_REL:is_a</string>" +
146         "</void>" +
147         "</void>" +
148         "</object>" +
149         "</void>" +
150         "<void property=\"value\">" +
151         "<string>" + pat_name + "</string>" +
152         "</void>" +
153         "</object>" +
154         "</void>");;
155     }
156 }
157
158 //fine
159 out.println("</void></object></java>");
160 }
161
162 }
```

## C.10 PathologyObj.java

```

1 package albero;
2
3 import java.io.BufferedReader;
4 import java.io.FileNotFoundException;
5 import java.io.FileOutputStream;
6 import java.io.FileReader;
7 import java.io.IOException;
8 import java.io.PrintStream;
9 import java.util.ArrayList;
10 import java.util StringTokenizer;
11
12 public class PathologyObj {
13     public String id="";
14     public String name="";
15     public String def="";
16     public String subset="";
17     public String comment="";
18     public ArrayList<String> synonym = new ArrayList<String>();
19     public ArrayList<String> xref = new ArrayList<String>();
20     public String is_a="";
21     public int pre = 0;
22     public int post = 0;
23     public int level = 0;
24
25     public PathologyObj(String id) throws IOException{
26         BufferedReader br = new BufferedReader( new FileReader("/home/
27             HumanDO.txt"));
28         String strLine;
29         while( (strLine = br.readLine()) != null ){
30             //System.out.println(strLine);
31             if(strLine.startsWith("id: "+id)){
```

```

31     //System.out.println(strLine);
32     this.id = strLine;
33     while( !(strLine = br.readLine()).startsWith("[Term]")){
34         if(strLine.startsWith("name:")) this.name = strLine;
35         else if(strLine.startsWith("def:")) this.def = strLine;
36         else if(strLine.startsWith("subset:")) this.subset = strLine;
37         else if(strLine.startsWith("synonym:")) this.synonym.add(strLine);
38         else if(strLine.startsWith("xref:")) this.xref.add(strLine);
39         else if(strLine.startsWith("is_a:")) {
40             StringTokenizer st;
41             st = new StringTokenizer(strLine, "!");
42             String tmp = st.nextToken();
43             is_a = tmp.substring(0,tmp.length()-1);
44         }
45     }
46     break;
47 }
48 }
49 }
50
51 public void tostring(){
52     System.out.println(id);
53     System.out.println(name);
54     System.out.println(def);
55     System.out.println(subset);
56     System.out.println(comment);
57     for(int i=0; i<synonym.size();i++)
58         System.out.println(synonym.get(i));
59     for(int i=0; i<xref.size();i++)
60         System.out.println(xref.get(i));
61     System.out.println(is_a);
62 }
63
64 public String toString(){
65     return id+" "+level;
66 }
67
68 }

```

## C.11 NTree.java

```

1 package albero;
2
3 import java.util.ArrayList;
4
5 public class NTree //classe pubblica albero
6 {
7     protected NTreeNode root; //radice dell'albero
8     public ArrayList<PathologyObj> array = new ArrayList<PathologyObj>
9         >(); //per dopo.. lista
10    public NTree() //costruttore
11    {
12        root=null;
13    }
14}

```

```

12
13
14     public void setRoot(NTreeNode root){
15         this.root = root;
16     }
17
18     public boolean insert(String parent_id, NTreeNode son) //metodo per
19         l'inserimento
20     {
21         if(this.search(root,parent_id)&&son!=null)
22         {
23             //System.out.println("trovato!");
24             NTreeNode temp=this.searchNode(root,parent_id);
25             NTreeNode temp2 = temp.firstSon;
26             temp.firstSon=son;
27             son.brother = temp2;
28             //settiamo level prima di inserire
29             son.element.level = temp.element.level+1;
30             return true;
31         }
32         else
33             System.out.println("Errore, nodo inesistente.
34                         Impossibile inserire i nodi dei figli");
35             //stampa errore
36             return false;
37     }
38     public boolean search(NTreeNode StartPoint, String node_id) //
39         metodo per la ricerca di un nodo
40     {
41         NTreeNode p = StartPoint;
42         if(p!=null){
43             boolean ris;
44             //System.out.println(p.element.id + " confronto "+node_id);
45             if(p.element.id.equals(node_id)) return true;
46             else
47             {
48                 NTreeNode t=p.firstSon;
49                 while(t!=null)
50                 {
51                     ris = search(t, node_id);
52                     if(ris) return true;
53                     t=t.brother;
54                 }
55             }
56         }
57         return false;
58     }
59     public NTreeNode searchNode(NTreeNode StartPoint, String node_id) //
60         metodo per la ricerca di un nodo
61     {
62         NTreeNode p = StartPoint;
63         NTreeNode ris;

```

```

61     if(p!=null){
62         if(p.element.id.equals(node_id)) return p;
63         else
64         {
65             NTreeNode t=p.firstSon;
66             while(t!=null)
67             {
68                 ris = searchNode(t, node_id);
69                 if(ris != null) return ris;
70                 t=t.brother;
71             }
72         }
73     }
74     return null;
75 }
76
77 public void preorder(){
78     preorder(root);
79 }
80 public void preorder(NTreeNode p){
81     if(p!=null){
82         System.out.println(p.toString());
83         NTreeNode t=p.firstSon;
84         while(t!=null){
85             preorder(t);
86             t=t.brother;
87         }
88     }
89 }
90
91 public void set_preorder(){
92     set_preorder(root,1);
93 }
94 public int set_preorder(NTreeNode p, int num){
95     if(p!=null){
96         p.element.pre = num;
97         System.out.println(p.toString()+" "+p.element.pre);
98         NTreeNode t=p.firstSon;
99         while(t!=null){
100             num = set_preorder(t, ++num);
101             t=t.brother;
102         }
103     }
104     return num;
105 }
106
107 public void postorder(){
108     postorder(root);
109 }
110 public void postorder(NTreeNode p){
111     if(p!=null){
112         NTreeNode t=p.firstSon;
113         while(t!=null){
114             postorder(t);

```

```

115             t=t.brother;
116         }
117         System.out.println(p.toString());
118     }
119 }
120
121 public void set_postorder(){
122     set_postorder(root,1);
123 }
124 public int set_postorder(NTreeNode p, int num){
125     if(p!=null){
126         NTreeNode t=p.firstSon;
127         while(t!=null){
128             num = set_postorder(t, num);
129             num++;
130             t=t.brother;
131         }
132         p.element.post = num;
133         System.out.println(p.toString()+" "+ p.element.post);
134     }
135     return num;
136 }
137
138 public void toArrayList(){
139     toArrayList(root);
140 }
141 public void toArrayList(NTreeNode p){
142     if(p!=null){
143         array.add(p.element);
144         NTreeNode t=p.firstSon;
145         while(t!=null){
146             toArrayList(t);
147             t=t.brother;
148         }
149     }
150 }
151 }
```

## C.12 NTreeNode.java

```

1 package albero;
2 public class NTreeNode //implementazione nodo albero con tipo generico E
3 {
4     protected NTreeNode firstSon; //primo figlio
5     protected NTreeNode brother; //fratello
6     protected PathologyObj element; //elemento del nodo
7     public NTreeNode() //costruttore
8     {
9         this(null,null, null);
10    }
11    public NTreeNode(PathologyObj element, NTreeNode firstSon)
12    {
13        this(element, firstSon, null);
```

```

14     }
15     public NTreeNode(PathologyObj element, NTreeNode firstSon, NTreeNode
16         brother)
16     {
17         this.element=element;
18         this.firstSon=firstSon;
19         this.brother=brother;
20     }
21     public String toString(){
22         return element.toString();
23     }
24 }
```

## C.13 TreeMaker.java

```

1 package albero;
2
3 import java.io.BufferedReader;
4 import java.io.FileNotFoundException;
5 import java.io.FileReader;
6 import java.io.IOException;
7 import java.util.ArrayList;
8 import java.util.StringTokenizer;
9
10 public class TreeMaker {
11     NTree tree;
12
13     public TreeMaker() throws IOException {
14         BufferedReader br = new BufferedReader(new FileReader("/home/
15             allnodesfromobo.txt"));
15         ArrayList<PathologyObj> list = new ArrayList<PathologyObj>();
16         StringTokenizer st = null;
17         String strLine,id;
18
19         //while per la creazione della lista di oggetti che ci interessano
20         while((strLine = br.readLine()) != null ){
21             st = new StringTokenizer(strLine, " ");
22             id = st.nextToken();
23             //System.out.println(id);
24             list.add(new PathologyObj(id));
25         }
26
27         //creiamo l'albero e colleghiamo i nodi
28         tree = new NTree();
29         for(int i=0; i<list.size(); i++)
30             if(list.get(i).id.equals("id: DOID:4")) {
31                 tree.setRoot(new NTreeNode(list.get(i),null,null));
32                 list.remove(i);
33             }
34
35         while(!list.isEmpty()){
36             for(int i=0; i<list.size(); i++){

```

```

37     if(i == 0) { System.out.println("Stampo albero:"); tree.preorder();
38         System.out.println();}
39     //attenzione agli obsoleti! sono 3!
40     if(list.get(i).is_a.equals("")){ System.out.println("obsolete: "+
41         list.get(i).id); list.remove(i); }
42     String parent_id = "id: "+list.get(i).is_a.substring(6);
43     System.out.println("io sono: "+list.get(i).id + " cerco: " +
44         parent_id);
45     if(tree.insert(parent_id, new NTreeNode(list.get(i),null,null))) {
46         list.remove(i);
47         System.out.println("nodo esistente!!!");}
48     }
49
50     //collego a mano gli ultimi con due is_a
51
52     //optical nerve disease
53     PathologyObj tmp = new PathologyObj("DOID:1891");
54     tmp.id = tmp.id+"bis";
55     tmp.is_a = "id: DOID:1393";
56     tree.insert(tmp.is_a, new NTreeNode(tmp,null,null));
57
58     //retinal cancer
59     tmp = new PathologyObj("DOID:4645");
60     tmp.id = tmp.id+"bis";
61     tmp.is_a = "id: DOID:2174";
62     tree.insert(tmp.is_a, new NTreeNode(tmp,null,null));
63
64     //retinal cell cancer
65     tmp = new PathologyObj("DOID:771");
66     tmp.id = tmp.id+"bis";
67     tmp.is_a = "id: DOID:4645bis";
68     tree.insert(tmp.is_a, new NTreeNode(tmp,null,null));
69
70     //retinoblastoma
71     tmp = new PathologyObj("DOID:768");
72     tmp.id = tmp.id+"bis";
73     tmp.is_a = "id: DOID:771bis";
74     tree.insert(tmp.is_a, new NTreeNode(tmp,null,null));
75
76     tree.set_preorder();
77     System.out.println();
78     tree.set_postorder();
79 }
80
81 public ArrayList<PathologyObj> array(){
82     tree.toArraylist();
83     return tree.array;
84 }
85
86 }
```

## C.14 RiempipiTable.java

```
1 package albero;
2
3 import java.io.FileOutputStream;
4 import java.io.IOException;
5 import java.io.PrintStream;
6 import java.util.ArrayList;
7
8 public class RiempipiTable {
9     //classe per creare script per riempire tabelle: PATHOLOGY, SYNONYM, XREF
10    //con quelle che conosco dalla ontologia
11    public static void main(String[] args) throws IOException {
12
13        FileOutputStream sql_file = new FileOutputStream("/home/Path-syn-xref.
14            sql");
14        PrintStream out = new PrintStream(sql_file);
15
16        TreeMaker op = new TreeMaker();
17        ArrayList<PathologyObj> list = op.array();
18
19        //riempio quelle nell'array
20
21        //aggiungo obsoleti
22        PathologyObj tmp = new PathologyObj("DOID:8600");
23        tmp.level= -1;
24        System.out.println("try:"+tmp.level);
25        list.add(tmp);
26
27        tmp = new PathologyObj("DOID:4625");
28        tmp.level= -1;
29        list.add(tmp);
30
31        tmp = new PathologyObj("DOID:410");
32        tmp.level= -1;
33        list.add(tmp);
34        System.out.println("try:"+tmp.level);
35
36        for(int i=0; i<list.size(); i++){
37            //System.out.println(list.size());
38            tmp = list.get(i);
39            System.out.println("try:"+tmp.level);
40            String def=",null";
41            if(!tmp.def.equals("")) def = ","+tmp.def.substring(5).replace("'", "
42                ''")+"'";
42            String subset=",null";
43            if(!tmp.subset.equals("")) subset = ","+tmp.subset.substring(8)+"'";
44            String comment=",null";
45            if(!tmp.comment.equals("")) comment = ","+tmp.comment.substring(9)+"",
46                ";
46            String pre=",null",post=",null";
47            if(tmp.pre != 0)
48                pre = ","+Integer.toString(tmp.pre);
```

```

49     if(tmp.post != 0)
50         post = ","+Integer.toString(tmp.post);
51
52     String sql = "INSERT INTO PATHOLOGY VALUES('"+tmp.id.substring(4)+"', "
53             +"','"+tmp.name.substring(6).replace("'", "\")+
54             "'"+comment+def+subset+pre+post+"','"+tmp.level+");";
55     out.println(sql);
56
57     for(int j=0; j<tmp.synonym.size(); j++){
58         sql = "INSERT INTO SYNONYM VALUES ('"+tmp.id.substring(4)+"', '"+tmp.
59             synonym.get(j).substring(9).replace("'", "\")+"');";
60         out.println(sql);
61     }
62
63     for(int j=0; j<tmp.xref.size(); j++){
64         sql = "INSERT INTO XREF VALUES ('"+tmp.id.substring(4)+"', '"+tmp.
65             xref.get(j).substring(6).replace("'", "\")+"');";
66         out.println(sql);
67     }
68 }
69 }
```

## C.15 Index.html

```

1  <!DOCTYPE html>
2  <html>
3  <head>
4      <meta charset="UTF-8">
5      <title>UCbase</title>
6      <link rel="icon" href="http://localhost/appswebositetemplate/favicon.ico" /
7          >
8      <link rel="stylesheet" href="css/style.css" type="text/css">
9  </head>
10 <body>
11     <div class="page">
12         <div class="sidebar">
13             <div id="logo">
14                 <a href="index.html"></a>
15             </div>
16             <ul class="navigation">
17                 <li class="selected">
18                     <a href="index.html">Home</a>
19                 </li>
20                 <li>
21                     <a href="uc_data_mining.php">UC Data Mining</a>
22                 </li>
23                 <li>
24                     <a href="related_works.html">Related Works</a>
25                 </li>
```

```

26         <a href="about_us.html">About Us</a>
27     </li>
28 </ul>
29 <p id="suggest"> Type an ID, a GENE NAME <br> or what you need and GO
30   !</p>
31 <form action="">
32   <input type="text" value="Rapid Search" onblur="this.value=!this.
33     value?'Rapid Search':this.value;" onfocus="this.select()"
34     onclick="this.value='';">
35   <input type="submit" value="" onclick="alert('Not implemented yet
36     !');">
37 </form>
38
39 <p id="mail">Comments, questions? <br>
40   Email <a href="mailto:vinxlemons@gmail.com">vinxlemons@gmail.com</a>
41   ></p>
42 </div>
43 <div class="content">
44   <div class="article">
45     <h2>Hi and welcome to the UCbase, visitor!</h2>
46     <p>UCBase is the nerve center for the study of ultraconservative
47       sequences: collects the latest news, research and information
48       about it, as well as offering new features and approaches to
49       research in the field of organic. Specifically provides the
50       ability to query an open database of the UC sequences of the
51       human genome as you could see in the "UC Data Mining" section.<
52     /p>
53   </div>
54   <div class="blog">
55     <div class="date">
56       <span>08-04</span>
57       <span>2013</span>
58     </div>
59     <div>
60       <h2>Things to Know</h2>
61       <div>
62         border
63       </div>
64     </div>
65     <ul>
66       <li>
67         <div>
68           <div>
69             <a href="#"></a>
70             <a onclick="alert('Not present!');">.pdf</a>
71           </div>
72         <div>
73           <h3>Why UCbase? </h3>
74           <p>UCbase arises from the need to create a hub about
75             ultraconservative sequences and an innovative way to query
76             the data in the biologic context. The database is built
77             online recovering genetic information and building a
78             real hierarchy of diseases to make possible correlations
79           .</p>

```

```

64          <a onclick="alert('Not implemented yet!');">Read more></a>
65      </div>
66  </div>
67 </li>
68 <li>
69   <div>
70     <div>
71       <a href="#"></a>
72       <a onclick="alert('Not present!');">.pdf</a>
73     </div>
74   <div>
75     <h3>What are the UC sequences?</h3>
76     <p>In biology, conserved sequences are similar or identical
77       sequences that occur within nucleic acid sequences (such
78       as RNA and DNA sequences), protein sequences, protein
79       structures or polymeric carbohydrates across species (
80       orthologous sequences) or within different molecules
81       produced by the same organism (paralogous sequences).</p>
82     <a onclick="alert('Not implemented yet!');">Read more></a>
83   </div>
84 </li>
85 <li>
86   <div>
87     <div>
88       <a href="#"></a>
89       <a onclick="alert('Not present!');">.pdf</a>
90     </div>
91     <div>
92       <h3>Database and technical information</h3>
93       <p>Cascading Style Sheets (CSS) were used for UCbase web
94         development to style web pages written in HTML and PHP.
95         The CSS specifications are maintained by the World Wide
96         Web Consortium (W3C).
97       UCbase was created using MySQL database under Debian Etch
98         Linux OS installed on a Quad Core Processor machine
99         with 32 GB RAM.</p>
100      <a onclick="alert('Not implemented yet!');">Read more></a>
101    </div>
102  </ul>
103 <div class="section">
104   <a href="#"></a>

```

**102 **

**103 </body>**

**104 </html>**

## C.16 Uc\_data\_mining.php

```
1  <!DOCTYPE html>
2  <html>
3  <head>
4  <head>
5      <meta charset="UTF-8">
6      <title>UCbase - UC Data Mining</title>
7      <link rel="icon" href="http://localhost/appswebstiteftemplate/favicon.ico"
8          />
9      <link rel="stylesheet" href="css/style.css" type="text/css">
10 </head>
11 <body>
12     <div class="page">
13         <div class="sidebar">
14             <div id="logo">
15                 <a href="index.html"></a>
16             </div>
17             <ul class="navigation">
18                 <li>
19                     <a href="index.html">Home</a>
20                 </li>
21                 <li class="selected">
22                     <a href="uc_data_mining.php">UC Data Mining</a>
23                 </li>
24                 <li>
25                     <a href="related_works.html">Related Works</a>
26                 </li>
27                 <li>
28                     <a href="about_us.html">About Us</a>
29                 </li>
30             </ul>
31             <p id="suggest"> Type an ID, a GENE NAME <br> or what you need and GO
32                 !</p>
33             <form action="">
34                 <input type="text" value="Rapid Search" onblur="this.value=this.
35                     value?'Rapid Search':this.value;" onfocus="this.select()"
36                     onclick="this.value='';">
37                 <input type="submit" value="" onclick="alert('Not yet implemented!')"
38                     ;">
39             </form>
40             <p id="mail">Comments, questions? <br>
41                 Email <a href="mailto:vinxlemons@gmail.com">vinxlemons@gmail.com</a>
42             </p>
43             <div class="content">
44                 <div>
45                     <h2>UC Data Mining</h2>
46                     <p>In this section you can browse and query the database as you like
47                         .
48                         You will be able to do semi-structured query or directly query the database
49                         through an SQL query.
```

```

44 This is one of the strengths of UCbase, pointing to a direct experience of
45 "Data Mining" for the purposes of research.
46 </p><p>Below you will find the structure of the database expressed in the E
47 /R standard language for databases and the query section.</p>
48 <a href="images/Diagramma1.png" target="_blank"></img></a>
50 </div>
51 <div class="data_mining">
52   <div>
53     <h2>Pre-formed Query</h2>
54   </div>
55   <ul>
56     <span>All info about </span>
57     <form action="result.php" method="GET">
58       <select name="uc.id" >
59         <?php
60           $con = mysql_connect("localhost", "client", "is1fs");
61           if (!$con)
62           {
63             die('Could not connect: ' . mysql_error());
64           }
65           mysql_select_db("UCbase", $con);
66           $query = "SELECT UC_NAME FROM UC ";
67           $result = mysql_query($query);
68
69           while($row = mysql_fetch_array($result, MYSQL_NUM))
70           {
71             if($row[0] == "uc.0")
72               echo '<option value="'. $row[0]. '" selected="selected">' .
73                 $row[0]. '</option>';
74             else
75               echo '<option value="'. $row[0]. '" >' . $row[0]. '</option>';
76           }
77         ?>
78       </select>
79       <input class="submit" type="submit" value="" >
80     </form>
81     <span> ?</span>
82   </ul>
83   <ul>
84     <span>All info about </span>
85     <form action="result_gene.php" method="GET">
86       <select name="gene_name" >
87         <?php
88           $query = "SELECT WIKIGENE_NAME FROM GENE ";
89           $result = mysql_query($query);
90
91
92
93

```

```

94     while($row = mysql_fetch_array($result, MYSQL_NUM))
95     {
96         if($row[0]!=""){
97             if($row[0] == "PEX14")
98                 echo '<option value="'. $row[0].'" selected="selected">' .
99                     $row[0].'</option>';
100            else
101                echo '<option value="'. $row[0].'" >' . $row[0].'</option>' ;
102        }
103    ?>
104    </select>
105    <input class="submit" type="submit" value="">
106    </form>
107    <span?</span>
108    </ul>
109    <ul>
110        <span id="third">All UC correled to</span>
111        <form action="result_search_pat.php">
112
113            <select id="tre" name="pat" >
114            <?php
115
116                $query = "SELECT NAME FROM PATHOLOGY ";
117
118                $result = mysql_query($query);
119
120                while($row = mysql_fetch_array($result, MYSQL_NUM))
121                {
122                    if($row[0]!=""){
123                        if($row[0] == "disease")
124                            echo '<option value="'. $row[0].'" selected="selected">' .
125                                $row[0].'</option>';
126                        else
127                            echo '<option value="'. $row[0].'" >' . $row[0].'</option>' ;
128                    }
129                ?>
130                </select>
131                <input class="submit" type="submit" value="">
132                </form>
133                <span style="width: 150px;">or its subtypes?</span>
134            </ul>
135            <ul>
136                <span id="third">All UC in </span>
137                <form action="result_search_uc.php" method="GET">
138
139                    <select name="chr_id" >
140                    <?php
141
142                        $query = "SELECT DISTINCT CHR FROM UC ";
143

```

```

144     $result = mysql_query($query);
145
146     while($row = mysql_fetch_array($result, MYSQL_NUM))
147     {
148         if($row[0]!=""){
149             if($row[0] == "1")
150                 echo '<option value="'. $row[0]. '" selected="selected">' .
151                     $row[0]. '</option>';
152             else
153                 echo '<option value="'. $row[0]. '">'. $row[0]. '</option>';
154             ;
155         }
156     ?>
157     </select>
158     <span> starting between </span>
159
160     <input type="text" style="width: 70px;" value="num_bp" name="num_bp1"
161         onblur="this.value=!this.value?'num_bp':this.value;
162             " onfocus="this.select()" onclick="this.value='';">
163
164     <span> and</span>
165
166     <input type="text" style="width: 70px;" value="num_bp" name="num_bp2"
167         onblur="this.value=!this.value?'num_bp':this.value;
168             " onfocus="this.select()" onclick="this.value='';">
169     <input class="submit" type="submit" value="">
170
171     </form>
172     <span> ?</span>
173     </ul>
174     <ul>
175         <span>Blast sequence:</span>
176         <form action="blastn.php" method="GET">
177
178             <input type="text" title="insert a nucleotide sequence! ex:
179                 ACGTACAGTACG" style='width:360px;' value="ACGTACAGTACG" name=
180                     "sequence" onblur="this.value=!this.value?'ACGTACAGTACG':
181                         this.value;" onfocus="this.select()" onclick="this.value='';
182                         ">
183             <input class="submit" type="submit" value="">
184
185         </form>
186     </ul>
187     <ul>
188         <span>Blast sequence:</span>
189         <form action="blastpat.php" method="GET">
190
191             <input type="text" title="insert a nucleotide sequence! ex:
192                 ACGTACAGTACG" value="ACGTACAGTACG" name="sequence" onblur=
193                     "this.value=!this.value?'ACGTACAGTACG':this.value;" onfocus=
194                         "this.select()" onclick="this.value='';">
195             <span> with pathology:</span>
196             <select name="pathology" >
197             <?php
198
199

```

```

185         $query = "SELECT NAME FROM PATHOLOGY ";
186
187         $result = mysql_query($query);
188
189         while($row = mysql_fetch_array($result, MYSQL_NUM))
190     {
191             if($row[0]!=""){
192                 if($row[0] == "disease")
193                     echo '<option value="'. $row[0]. '" selected="selected">' .
194                         $row[0]. '</option>';
195                 else
196                     echo '<option value="'. $row[0]. '">'. $row[0]. '</option>' ;
197             }
198         ?>
199         </select>
200         <input class ="submit" type="submit" value="">
201     </form>
202     </ul>
203     <div>
204         <h2 id="second">Type your own Query!</h2>
205         <ul>
206
207             <form action="table.php">
208
209                 <input id="second" type="text" name="query" value="SELECT *
210                     FROM UC" onblur="this.value!=this.value?'SELECT * FROM UC
211                     ':this.value;" onfocus="this.select()" onclick="this.value
212                     ='';">
213                 <input class="submit" type="submit" value="">
214             </form>
215         </ul>
216         <div class="section">
217             <a id='end' href="#"></a>
218         </div>
219     </div>
220     </div>
221 </body>
222 </html>
```

## C.17 Table.php

```

1
2 <html>
3 <body>
4
5 <!DOCTYPE html>
6 <html>
7 <head>
```

```

8   <meta charset="UTF-8">
9   <title>UCbase - UC Data Mining</title>
10  <link rel="icon" href="http://localhost/appswsitetemplate/favicon.ico"
11    />
12  <link rel="stylesheet" href="css/style.css" type="text/css">
13 </head>
14 <body>
15   <div class="page">
16     <div class="sidebar">
17       <div id="logo">
18         <a href="index.html"></a>
19       </div>
20       <ul class="navigation">
21         <li>
22           <a href="index.html">Home</a>
23         </li>
24         <li class="selected">
25           <a href="uc_data_mining.php">UC Data Mining</a>
26         </li>
27         <li>
28           <a href="related_works.html">Related Works</a>
29         </li>
30         <li>
31           <a href="about_us.html">About Us</a>
32         </li>
33       </ul>
34       <p id="suggest"> Type an ID, a GENE NAME <br> or what you need and GO
35         !</p>
36       <form action="">
37         <input type="text" value="Rapid Search" onblur="this.value!=this.
38           value?'Rapid Search':this.value;" onfocus="this.select()"
39           onclick="this.value='';">
40         <input type="submit" value="" onclick="alert('Not yet implemented!')"
41           ;">
42       </form>
43
44       <p id="mail">Comments, questions? <br>
45         Email <a href="mailto:vinxlemons@gmail.com">vinxlemons@gmail.com</a>
46       </p>
47     </div>
48     <div class="content">
49       <div>
50         <h2>Search Result</h2>
51       </div>
52       <div class="data_mining">
53         <div class="result">
54           <h2 style='font-size:15px;'>
55             <?php
56               $query = $_GET["query"];
57               $query = str_replace("<", "try", $query);
58               $query = filter_var($query, FILTER_SANITIZE_STRING);
59               $query = str_replace("try", "<", $query);
60               echo $query;
61             ?>

```

```

56         </h2>
57     </div>
58     <div class="section" id="result">
59     <?php
60 $con = mysql_connect("localhost", "client", "is1fs");
61 if (!$con)
62 {
63 die('Could not connect: ' . mysql_error());
64 }
65
66 mysql_select_db("UCbase", $con);
67 $query = $_GET["query"];
68 $query = str_replace("<", "try", $query);
69 $query = filter_var($query, FILTER_SANITIZE_STRING);
70 $query = str_replace("try", "<", $query);
71
72 $query = html_entity_decode($query, ENT_QUOTES);
73 $result = mysql_query($query);
74
75 if($result != false){
76     echo "<table id='big'>";
77     $i=0;
78     $nohits=1;
79     while($row = mysql_fetch_array($result, MYSQL_ASSOC))
80     {
81         $nohits =0;
82         if($i%2 == 0) echo "<tr class ='alt'>";
83         else echo "<tr>";
84
85         if($i==0){
86             foreach($row as $key => $value)
87             {
88                 echo "<th>";
89                 echo $key;
90                 echo "</th>";
91             }
92             echo "</tr>";
93         }
94
95         if($i==0) echo "<tr>";
96         foreach($row as $x=>$x_value)
97         {
98             echo "<td>";
99             if($x == "UC_NAME"){
100                 $uc_name = $x_value;
101                 echo "<a href='result.php?uc_id=". $uc_name ."'>". $x_value . "</a>";
102             }
103             else if($x == "UC_NAME"){
104                 $uc_name = $x_value;
105                 echo $x_value;
106             }
107             else if($x == "DEFINITION" && $x_value != ""){

```

```

108         $x_value = preg_replace('!(http|ftp|scp)(s)?:\//[a-zA-Z0-9.?&
109             _=/]+!', "<a target='_blank' href=\"\\0\">\\0</a>",
110             $x_value);
111         echo $x_value;
112     }
113     else if($x == "SEQUENCE") {
114         echo strtolower(substr($x_value, 0, 13)."...");
115         echo "<a class='button' href='sequence.php?uc_id=".$uc_name."'> Get
116             Sequence!</a>"; }
117     else
118         echo $x_value;
119     echo "</td>";
120     $i++;
121 }
122 if($nohits)
123     echo "<p style='text-align:center; color:black; font-size:24px;'>NO
124             HITS FOUND.</p>";
125 echo "</table>";
126 }
127 else{
128     echo "<p style='margin-left:5px; font-size:24px; color:black;'>Query
129             error.</p>";
130     $message = 'Invalid query: ' . mysql_error() . "\n";
131     $message .= 'Whole query: ' . $query;
132     echo $message;
133 }
134 mysql_close($con);
135     </div>
136     </div>
137     </div>
138     </div>
139 </body>
140 </html>
141 </body>
142 </html>

```

## C.18 Sequence.php

```

1
2 <?php $con = mysql_connect("localhost","client","is1fs");
3     if (!$con){
4         die('Could not connect: ' . mysql_error());
5     }
6
7     $uc = filter_input(INPUT_GET, "uc_id", FILTER_SANITIZE_STRING);
8     mysql_select_db("UCbase", $con);
9     $result = mysql_query("SELECT * FROM UC WHERE UC_NAME='".$uc."'");

```

```

10 $row = mysql_fetch_array($result);
11 echo "<html> <body> <p style='font-size:12px;'> double click on it to
12     select! <br><br>".$row["SEQUENCE"]."</p></body></html>";
13 mysql_close($con);
14 ?>

```

## C.19 Result\_search\_uc.php

```

1
2 <?php
3
4 $chr = filter_input(INPUT_GET, "chr_id", FILTER_SANITIZE_STRING);
5 $num_bp1 = filter_input(INPUT_GET, "num_bp1", FILTER_SANITIZE_NUMBER_INT);
6 $num_bp2 = filter_input(INPUT_GET, "num_bp2", FILTER_SANITIZE_NUMBER_INT);
7
8 $query = "SELECT * FROM UC WHERE UC.CHR= '".$chr."' AND UC.START BETWEEN
9     ".$num_bp1." AND ".$num_bp2."";
10 header("location: http://localhost/UCbase/table.php?query=".$query);
11 ?>

```

## C.20 Result\_search\_pat.php

```

1
2 <?php
3
4 $pat = filter_input(INPUT_GET, "pat", FILTER_SANITIZE_STRING);
5 if($pat == "retinal cancer")
6     $query = "SELECT UC_NAME,NAME,MIN(LEVEL)AS LEVEL FROM UC INNER JOIN (
7             SELECT ENSAMBL_GENE_ID,NAME,LEVEL FROM RELATED_TO INNER JOIN (
8                 SELECT ID,NAME,LEVEL FROM PATHOLOGY WHERE PREORDER >= (SELECT
9                     PREORDER FROM PATHOLOGY WHERE ID='DOID:4645') AND POSTORDER <=(SELECT
10                     POSTORDER FROM PATHOLOGY WHERE ID='DOID:4645'))AS B WHERE
11                     RELATED_TO.ID = B.ID)AS C WHERE UC.ENSAMBL_GENE_ID = C.
12                     ENSAMBL_GENE_ID group by uc_name, name ORDER BY LEVEL";
13 else if($pat == "optic nerve disease")
14     $query = "SELECT UC_NAME,NAME,MIN(LEVEL)AS LEVEL FROM UC INNER JOIN (
15             SELECT ENSAMBL_GENE_ID,NAME,LEVEL FROM RELATED_TO INNER JOIN (
16                 SELECT ID,NAME,LEVEL FROM PATHOLOGY WHERE PREORDER >= (SELECT
17                     PREORDER FROM PATHOLOGY WHERE ID='DOID:1891') AND POSTORDER <=(SELECT
18                     POSTORDER FROM PATHOLOGY WHERE ID='DOID:1891'))AS B WHERE
19                     RELATED_TO.ID = B.ID)AS C WHERE UC.ENSAMBL_GENE_ID = C.
20                     ENSAMBL_GENE_ID group by uc_name, name ORDER BY LEVEL";
21 else if($pat == "retinal cell cancer")
22     $query = "SELECT UC_NAME,NAME,MIN(LEVEL) AS LEVEL FROM UC INNER JOIN (
23             SELECT ENSAMBL_GENE_ID,NAME,LEVEL FROM RELATED_TO INNER JOIN (
24                 SELECT ID,NAME,LEVEL FROM PATHOLOGY WHERE PREORDER >= (SELECT
25                     PREORDER FROM PATHOLOGY WHERE ID='DOID:771') AND POSTORDER <=(SELECT
26                     POSTORDER FROM PATHOLOGY WHERE ID='DOID:771') )AS B WHERE
27                     RELATED_TO.ID = B.ID)AS C WHERE UC.ENSAMBL_GENE_ID = C.
28                     ENSAMBL_GENE_ID group by uc_name, name ORDER BY LEVEL";
29 else if($pat == "retinoblastoma")
30     $query = "SELECT UC_NAME,NAME, MIN(LEVEL) AS LEVEL FROM UC INNER JOIN (
31             SELECT ENSAMBL_GENE_ID,NAME,LEVEL FROM RELATED_TO INNER JOIN (

```

```

    SELECT ID,NAME,LEVEL FROM PATHOLOGY WHERE PREORDER >= (SELECT
    PREORDER FROM PATHOLOGY WHERE ID='DOID:768') AND POSTORDER <=(

    SELECT POSTORDER FROM PATHOLOGY WHERE ID='DOID:768') )AS B WHERE
    RELATED_TO.ID = B.ID)AS C WHERE UC.ENSAMBL_GENE_ID = C.
    ENSAMBL_GENE_ID group by uc_name, name ORDER BY LEVEL";
13   else $query = "SELECT UC_NAME,NAME,MIN(LEVEL)AS LEVEL FROM UC INNER JOIN
    (SELECT ENSAMBL_GENE_ID,NAME,LEVEL FROM RELATED_TO INNER JOIN (SELECT
    ID,NAME,LEVEL FROM PATHOLOGY WHERE PREORDER >= (SELECT PREORDER FROM
    PATHOLOGY WHERE NAME='".$pat."') AND POSTORDER <=(SELECT POSTORDER
    FROM PATHOLOGY WHERE NAME='".$pat."') ORDER BY LEVEL)AS B WHERE
    RELATED_TO.ID = B.ID)AS C WHERE UC.ENSAMBL_GENE_ID = C.
    ENSAMBL_GENE_ID group by uc_name, name ORDER BY LEVEL";
14
15   header("location: http://localhost/UCbase/table.php?query=".$query);
16 ?>

```

## C.21 Result\_gene.php

```

1 <html>
2 <body>
3
4
5 <!DOCTYPE html>
6 <!-- Website template by freewebsitetemplates.com -->
7 <html>
8 <head>
9   <meta charset="UTF-8">
10  <title>UCbase - UC Data Mining</title>
11  <link rel="icon" href="http://localhost/appswebsitetemplate/favicon.ico"
      />
12  <link rel="stylesheet" href="css/style.css" type="text/css">
13 </head>
14 <body>
15  <div class="page">
16    <div class="sidebar">
17      <div id="logo">
18        <a href="index.html"></a>
19      </div>
20      <ul class="navigation">
21        <li>
22          <a href="index.html">Home</a>
23        </li>
24        <li class="selected">
25          <a href="uc_data_mining.php">UC Data Mining</a>
26        </li>
27        <li>
28          <a href="related_works.html">Related Works</a>
29        </li>
30        <li>
31          <a href="about_us.html">About Us</a>
32        </li>
33      </ul>

```

```

34      <p id="suggest"> Type an ID, a GENE NAME <br> or what you need and GO
35      !</p>
36      <form action="">
37          <input type="text" value="Rapid Search" onblur="this.value!=this.
38              value?'Rapid Search':this.value;" onfocus="this.select()"
39              onclick="this.value='';">
40          <input type="submit" value="" onclick="alert('Not yet implemented!')"
41              ;">
42      </form>
43
44      <p id="mail">Comments, questions? <br>
45          Email <a href="mailto:vinxlemons@gmail.com">vinxlemons@gmail.com</a></p>
46
47      <div class="content">
48          <div>
49              <h2>Search Result</h2>
50          </div>
51          <div class="data_mining">
52              <div class="result">
53                  <h2> <?php
54                      $gene_name = filter_input(INPUT_GET, "gene_name",
55                          FILTER_SANITIZE_STRING);
56                      echo $gene_name;
57                      echo "</h2>" ;
58                      echo "</div>";
59                      $con = mysql_connect("localhost", "client", "is1fs");
60                      if (!$con){
61                          die('Could not connect: ' . mysql_error());
62                      }
63                      $gene = "";
64                      mysql_select_db("UCbase", $con);
65                      $result = mysql_query("SELECT * FROM GENE WHERE WIKIGENE_NAME='"
66                          . $gene_name . "'");
67
68                      $row = mysql_fetch_array($result); //solo una dovrebbe essere
69                      echo "<table id='result'>";
70                      if($row != null){
71                          foreach($row as $key => $value) //elimino indici numerici
72                          {
73                              if(is_int($key)){
74                                  unset($row[$key]);
75                              }
76                          }
77                          foreach($row as $x=>$x_value)
78                          {
79                              $i=0;
80                              if($i%2 == 0) echo "<tr class ='alt'>";
81                              else echo "<tr>";
82                              echo "<td class = 'one'>";
83                              echo ucfirst(strtolower($x));
84                              echo "</td>";
85                              echo "<td>";
86                              if($x == "SEQUENCE") echo substr($x_value, 0, 20)."...";
87                              else {if($x_value != "") echo $x_value; else echo "none";}

```

```

81         if($x == "SEQUENCE")
82             echo "<a class='button' href='sequence.php?uc_id=".$uc.">
83                 Get Sequence!</a>";
84             echo "</td>";
85             echo "</tr>";
86             $i++;
87         }
88         echo "</table>";
89         $gene = $row["ENSAMBL_GENE_ID"];
90     }
91     else
92         echo "<p style='margin-left:22px; margin-top: 20px; font-size:24
93 px; color:black;'>Not found.</p>";
94
95     $result = mysql_query("SELECT SPLICING.* FROM SPLICING,HAS WHERE HAS.
96         ENSAMBL_GENE_ID ='".$gene."' AND HAS.EVENT_TYPE_COD = SPLICING.
97         EVENT_TYPE_COD ");
98     $i =0;
99     while(($row = mysql_fetch_array($result, MYSQL_ASSOC)) != null){
100         if($row != null && $i == 0) echo "<div class='result'><h2> Gene Splicing
101             </h2></div>";
102         echo "<table id='result'>";
103         foreach($row as $x=>$x_value)
104             {
105                 $i=0;
106                 if($i%2 == 0) echo "<tr class ='alt'>";
107                 else echo "<tr>";
108                 echo "<td class = 'one'>";
109                 echo ucfirst(strtolower($x));
110                 echo "</td>";
111                 echo "<td>";
112                 echo "</td>";
113                 $i++;
114             }
115         echo "</table>";
116         $i++;
117     }
118     $result = mysql_query("SELECT UC.* FROM UC WHERE UC.ENSAMBL_GENE_ID ='".
119         $gene."'");
120     $i =0;
121     while(($row = mysql_fetch_array($result, MYSQL_ASSOC)) != null){
122         if($row != null && $i == 0) echo "<div class='result'><h2> Gene correled
123             UC</h2></div>";
124         echo "<table id='result'>";
125         $uc ="";
126         foreach($row as $x=>$x_value)
127             {
128                 $i=0;
129                 if($i%2 == 0) echo "<tr class ='alt'>";
130                 else echo "<tr>";
131                 echo "<td class = 'one'>";
132                 echo ucfirst(strtolower($x));

```

```

128         echo "</td>";
129         echo "<td>";
130         if($x == "UC_NAME") {$uc = $x_value; $x_value ="<a href='result
131             .php?uc.id=".$uc."'>".$uc."</a>"; }
132         if($x == "SEQUENCE") echo substr($x_value, 0, 20)."...";
133         else {if($x_value != "") echo $x_value; else echo "none";}
134         if($x == "SEQUENCE")
135             echo "<a class='button' href='sequence.php?uc_id=".$uc."'>
136                 Get Sequence!</a>";
137             echo "</td>";
138             echo "</tr>";
139             $i++;
140         }
141     }
142
143     $result = mysql_query("SELECT PATHOLOGY.* FROM PATHOLOGY,RELATED_TO WHERE
144         RELATED_TO.ENSAMBL_GENE_ID ='".$gene."' AND RELATED_TO.ID = PATHOLOGY.
145         ID");
146
147     $i =0;
148     while(($row = mysql_fetch_array($result, MYSQL_ASSOC)) != null){
149         if($row != null && $i == 0) echo "<div class='result'><h2> Gene correled
150             pathology</h2></div>";
151         echo "<table id='result'>";
152         foreach($row as $x=>$x_value)
153         {
154             if($i%2 == 0) echo "<tr class ='alt'>";
155             else echo "<tr>";
156             echo "<td class = 'one'>";
157             echo ucfirst(strtolower($x));
158             echo "</td>";
159             echo "<td>";
160             if($x == "ID"){ $pat = $x_value; }
161             if($x == "DEFINITION" && $x_value != ""){
162                 $x_value = preg_replace('!(http|ftp|scp)(s)?://[a-zA-Z0-9.?&
163                     _=-]+!', "<a target='_blank' href=\"$\\0\"$\\1</a>",
164                     $x_value);
165                 echo $x_value;
166             }
167             $i=0;
168             $result2 = mysql_query("SELECT SYNONYM.NAME FROM SYNONYM WHERE
169                 SYNONYM.ID ='".$pat."'");
170             while(($row2 = mysql_fetch_array($result2, MYSQL_ASSOC)) != null){
171                 foreach($row2 as $x2=>$x_value2)
172                 {
173                     if($i%2 == 0) echo "<tr class ='alt'>";
174                     else echo "<tr>";

```

```

174         echo "<td class = 'one'>";
175         echo "Synonym".$i;
176         echo "</td>";
177         echo "<td>";
178         if($x_value2 != "") echo $x_value2; else echo "none";
179         echo "</td>";
180         echo "</tr>";
181         $i++;
182     }
183 }
184 $i=0;
185 $result3 = mysql_query("SELECT XREF.REF FROM XREF WHERE XREF.ID ='".
186     $pat.''");
187 while(($row3 = mysql_fetch_array($result3, MYSQL_ASSOC)) != null){
188     foreach($row3 as $x3=>$x_value3)
189     {
190         if($i%2 == 0) echo "<tr class ='alt'>";
191         else echo "<tr>";
192         echo "<td class = 'one'>";
193         echo "Xref".$i;
194         echo "</td>";
195         echo "<td>";
196         if($x_value3 != "") echo $x_value3; else echo "none";
197         echo "</td>";
198         echo "</td>";
199         $i++;
200     }
201     echo "</table>";
202     $i++;
203 }
204 mysql_close($con);
205     ?
206     <div class="section">
207         <a href="#"></a>
208     </div>
209     </div>
210     </div>
211     </div>
212 </body>
213 </html>
214
215 </body>
216 </html>
```

## C.22 Result.php

```

1
2 <html>
3 <body>
4
5 <!DOCTYPE html>
6 <html>
```

```

7 <head>
8   <meta charset="UTF-8">
9   <title>UCbase - UC Data Mining</title>
10  <link rel="icon" href="http://localhost/appswesitetemplate/favicon.ico"
11    />
12  <link rel="stylesheet" href="css/style.css" type="text/css">
13 </head>
14 <body>
15   <div class="page">
16     <div class="sidebar">
17       <div id="logo">
18         <a href="index.html"></a>
19       </div>
20       <ul class="navigation">
21         <li>
22           <a href="index.html">Home</a>
23         </li>
24         <li class="selected">
25           <a href="uc_data_mining.php">UC Data Mining</a>
26         </li>
27         <li>
28           <a href="related_works.html">Related Works</a>
29         </li>
30         <li>
31           <a href="about_us.html">About Us</a>
32         </li>
33       </ul>
34       <p id="suggest"> Type an ID, a GENE NAME <br> or what you need and GO
35         !</p>
36       <form action="">
37         <input type="text" value="Rapid Search" onblur="this.value!=this.
38           value?'Rapid Search':this.value;" onfocus="this.select()"
39           onclick="this.value='';">
40         <input type="submit" value="" onclick="alert('Not yet implemented!')"
41           ;">
42       </form>
43
44       <p id="mail">Comments, questions? <br>
45         Email <a href="mailto:vinxlemons@gmail.com">vinxlemons@gmail.com</a>
46       </p>
47     </div>
48     <div class="content">
49       <div>
50         <h2>Search Result</h2>
51       </div>
52       <div class="data_mining">
53         <div class="result">
54           <h2> <?php
55             $uc = filter_input(INPUT_GET, "uc_id", FILTER_SANITIZE_STRING);
56             echo $uc;
57             echo "</h2>" ;
58             echo "</div>";
59             $con = mysql_connect("localhost", "client", "is1fs");
60             if (!$con){

```

```

55         die('Could not connect: ' . mysql_error());
56     }
57
58     mysql_select_db("UCbase", $con);
59     $result = mysql_query("SELECT * FROM UC WHERE UC_NAME=". $uc. "'");
60             );
61
62     $row = mysql_fetch_array($result); //solo una dovrebbe essere
63     echo "<table id='result'>";
64     if($row != null){
65         foreach($row as $key => $value) //elimino indici numerici
66         {
67             if(is_int($key)){
68                 unset($row[$key]);
69             }
70             $i=0;
71             foreach($row as $x=>$x_value)
72             {
73                 if($i%2 == 0) echo "<tr class ='alt'>";
74                 else echo "<tr>";
75                 echo "<td class = 'one'>";
76                 echo ucfirst(strtolower($x));
77                 echo "</td>";
78                 echo "<td>";
79                 if($x == "SEQUENCE") echo substr($x_value, 0, 20)."...";
80                 else {if($x_value != "") echo $x_value; else echo "none";}
81                 if($x == "SEQUENCE")
82                     echo "<a class='button' href='sequence.php?uc_id=".$uc.">".
83                         Get Sequence!</a>";
84                     echo "</td>";
85                     echo "</tr>";
86                     $i++;
87             }
88             echo "</table>";
89         }
90         else
91             echo "<p style='margin-left:22px;margin-top: 20px; font-size:24
92 px; color:black;'>Not found.</p>";
93         $gene = ""; //definisco per dopo.
94         if($row["ENSAMBL_GENE_ID"] != ""){
95             echo "<div class='result'><h2>". $row["ENSAMBL_GENE_ID"]. "</h2><
96             div>";
97
98             $gene = $row["ENSAMBL_GENE_ID"];
99
100            mysql_select_db("UCbase", $con);
101            $result = mysql_query("SELECT * FROM GENE WHERE ENSAMBL_GENE_ID
102 =".$gene. "'");
103
104            $row = mysql_fetch_array($result); //solo una dovrebbe essere
105            echo "<table id='result'>";
106            foreach($row as $key => $value) //elimino indici numerici
107            {

```

```

104         if(is_int($key)){
105             unset($row[$key]);
106         }
107     }
108     $i=0;
109     foreach($row as $x=>$x_value)
110     {
111         if($i%2 == 0) echo "<tr class ='alt'>";
112         else echo "<tr>";
113         echo "<td class = 'one'>";
114         echo ucfirst(strtolower($x));
115         echo "</td>";
116         echo "<td>";
117         echo $x_value;
118         echo "</td>";
119         echo "</tr>";
120         $i++;
121     }
122     echo "</table>";
123 }
124
125 $result = mysql_query("SELECT SPlicing.* FROM SPlicing,HAS WHERE HAS.
126     ENSAMBL_GENE_ID ='". $gene. "' AND HAS.EVENT_TYPE_COD = SPlicing.
127     EVENT_TYPE_COD ");
128     $i =0;
129     while(($row = mysql_fetch_array($result, MYSQL_ASSOC)) != null){
130         if($row != null && $i == 0) echo "<div class='result'><h2> Gene Splicing
131             </h2></div>";
132         echo "<table id='result'>";
133         $i=0;
134         foreach($row as $x=>$x_value)
135         {
136             if($i%2 == 0) echo "<tr class ='alt'>";
137             else echo "<tr>";
138             echo "<td class = 'one'>";
139             echo ucfirst(strtolower($x));
140             echo "</td>";
141             echo "<td>";
142             if($x_value != "") echo $x_value; else echo "none";
143             echo "</td>";
144             echo "</tr>";
145             $i++;
146         }
147     echo "</table>";
148     $i++;
149 }
150     $i =0;
151     while(($row = mysql_fetch_array($result, MYSQL_ASSOC)) != null){
152         if($row != null && $i == 0) echo "<div class='result'><h2> UC correled
153             SNP</h2></div>";

```

```

153     echo "<table id='result'>";
154     $i=0;
155     foreach($row as $x=>$x_value)
156     {
157         if($i%2 == 0) echo "<tr class ='alt'>";
158         else echo "<tr>";
159         echo "<td class = 'one'>";
160         echo ucfirst(strtolower($x));
161         echo "</td>";
162         echo "<td>";
163         if($x_value != "") echo $x_value; else echo "none";
164         echo "</td>";
165         echo "</tr>";
166         $i++;
167     }
168     echo "</table>";
169     $i++;
170 }
171
172 $result = mysql_query("SELECT PATHOLOGY.* FROM PATHOLOGY,RELATED_TO WHERE
173     RELATED_TO.ENSAMBL_GENE_ID ='".$gene."' AND RELATED_TO.ID = PATHOLOGY.
174     ID");
175
176 $i =0;
177 while(($row = mysql_fetch_array($result, MYSQL_ASSOC)) != null){
178     if($row != null && $i == 0) echo "<div class='result'><h2> Gene correled
179         pathology</h2></div>";
180
181     if(strstr($row['ID'] , "bis"))
182         continue;
183     echo "<table id='result'>";
184     foreach($row as $x=>$x_value)
185     {
186         if($i%2 == 0) echo "<tr class ='alt'>";
187         else echo "<tr>";
188         echo "<td class = 'one'>";
189         echo ucfirst(strtolower($x));
190         echo "</td>";
191         if($x == "ID"){ $pat = $x_value; }
192         if($x == "DEFINITION" && $x_value != ""){
193             $x_value = preg_replace('!^(http|ftp|scp)(s)?:\//([a-zA-Z0-9.?
194                 _=-]+)!','<a target=_blank href="'.$x_value.'>'.'</a>',
195                 $x_value);
196             echo $x_value;
197         }
198         else if($x_value != "") echo $x_value; else echo "none";
199         echo "</td>";
200         echo "</tr>";
201         $i++;
202     }
203     $i=0;
204     $result2 = mysql_query("SELECT SYNONYM.NAME FROM SYNONYM WHERE
205         SYNONYM.ID ='".$pat."'");

```

```

201     while(($row2 = mysql_fetch_array($result2, MYSQL_ASSOC)) != null){
202         foreach($row2 as $x2=>$x_value2)
203         {
204             if($i%2 == 0) echo "<tr class ='alt'>";
205             else echo "<tr>";
206             echo "<td class = 'one'>";
207             echo "Synonym".$i;
208             echo "</td>";
209             echo "<td>";
210             if($x_value2 != "") echo $x_value2; else echo "none";
211             echo "</td>";
212             echo "</tr>";
213             $i++;
214         }
215     }
216     $i=0;
217     $result3 = mysql_query("SELECT XREF.REF FROM XREF WHERE XREF.ID ='".
218     $pat."'");
219     while(($row3 = mysql_fetch_array($result3, MYSQL_ASSOC)) != null){
220         foreach($row3 as $x3=>$x_value3)
221         {
222             if($i%2 == 0) echo "<tr class ='alt'>";
223             else echo "<tr>";
224             echo "<td class = 'one'>";
225             echo "Xref".$i;
226             echo "</td>";
227             echo "<td>";
228             if($x_value3 != "") echo $x_value3; else echo "none";
229             echo "</td>";
230             $i++;
231         }
232     }
233     echo "</table>";
234     $i++;
235 }
236 mysql_close($con);
237     ?
238     <div class="section">
239         <a href="#"></a>
240         </div>
241     </div>
242     </div>
243     </div>
244 </body>
245 </html>
246 </body>
247 </html>

```

## C.23 Related\_works.html

```

2  <!DOCTYPE html>
3  <html>
4  <head>
5      <meta charset="UTF-8">
6      <title>UCbase - Related Works</title>
7      <link rel="icon" href="http://localhost/appswebsitetemplate/favicon.ico" /
8          >
9      <link rel="stylesheet" href="css/style.css" type="text/css">
10     </head>
11     <body>
12         <div class="page">
13             <div class="sidebar">
14                 <div id="logo">
15                     <a href="index.html"></a>
16                 </div>
17                 <ul class="navigation">
18                     <li>
19                         <a href="index.html">Home</a>
20                     </li>
21                     <li>
22                         <a href="uc_data_mining.php">UC Data Mining</a>
23                     </li>
24                     <li class="selected">
25                         <a href="related_works.html">Related Works</a>
26                     </li>
27                     <li>
28                         <a href="about_us.html">About Us</a>
29                     </li>
30                 </ul>
31                 <p id="suggest"> Type an ID, a GENE NAME <br> or what you need and GO
32                     !</p>
33                 <form action="">
34                     <input type="text" value="Rapid Search" onblur="this.value!=this.
35                         value?'Rapid Search':this.value;" onfocus="this.select()"
36                         onclick="this.value='';">
37                     <input type="submit" value="" onclick="alert('Not implemented yet
38                         !');">
39                 </form>
39
40             <div class="content">
41                 <div>
42                     <h2>Related works</h2>
43                     <p>In this section you will find all information relating to the UC
44                         sequences even outside of our research laboratory.
45                         Paper and news about it will be posted below.
46                     </p>
47                 </div>
48                 <div class="blog">
49                     <div class="date">

```

```

49      <span>12-03</span>
50      <span>2013</span>
51    </div>
52    <div>
53      <h2>Articles and Papers</h2>
54      <div>
55        border
56      </div>
57    </div>
58    <ul>
59      <li>
60        <div>
61          <div>
62            <a href="#"></a>
63            <a href="http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2686429
64              /">.pdf</a>
65          </div>
66          <div>
67            <h3>UCbase & miRfunc: a database of ultraconserved sequences
68              and microRNA function</h3>
69            <p>One hundred and eighty-one ultraconserved sequences (UCRs
70              ) longer than 200 bases were discovered in the genomes
71              of human, mouse and rat. These are DNA sequences showing
72              100% identity among the three species... </p>
73            <a onclick="alert('Not implemented yet!');">Read more>></a>
74          </div>
75        </div>
76      </li>
77      <li>
78        <div>
79          <div>
80            <a href="#"></a>
81            <a href="http://database.oxfordjournals.org/content/2011/
82              bar049.full">.pdf</a>
83          </div>
84          <div>
85            <h3>The BioMart project</h3>
86            <p>BioMart is a unique open source data federation
87              technology that provides unified access to distributed
88              databases storing a wide range of data. This DATABASE
89              issue recognizes BioMart's outstanding contributions to
90              bioinformatics and documents the achievements of the
91              BioMart community, which has grown impressively over the
92              last ten years to become what it is today...</p>
93            <a onclick="alert('Not implemented yet!');">Read more>></a>
94          </div>
95        </div>
96      </li>
97      <li>
98        <div>
99          <div>
100            <a href="#"></a>
101            <a href="http://www.mirbase.org/">.pdf</a>
102        </div>

```

```

91      <div>
92          <h3>miRBase: the microRNA database</h3>
93          <p>The miRBase database is a searchable database of
94              published miRNA sequences and annotation. Each entry in
95              the miRBase Sequence database represents a predicted
96              hairpin portion of a miRNA transcript (termed mir in the
97              database), with information on the location and
98              sequence of the mature miRNA sequence (termed miR). Both
99              hairpin and mature sequences are available for
100             searching and browsing, and entries can also...</p>
101             <a onclick="alert('Not implemented yet!');">Read more></a>
102         </div>
103     </div>
104 </li>
105 </ul>
106 <div class="section">
107     <a href="#"></a>
108 </div>
109 </div>
110 </body>
111 </html>

```

## C.24 Blastpat.php

```

1 <html>
2 <head>
3     <meta charset="UTF-8">
4     <title>UCbase - UC Data Mining</title>
5     <link rel="icon" href="http://localhost/appswebositetemplate/favicon.ico" />
6     <link rel="stylesheet" href="css/style.css" type="text/css">
7 </head>
8 <body>
9     <div class="page">
10         <div class="sidebar">
11             <div id="logo">
12                 <a href="index.html"></a>
13             </div>
14             <ul class="navigation">
15                 <li>
16                     <a href="index.html">Home</a>
17                 </li>
18                 <li class="selected">
19                     <a href="uc_data_mining.php">UC Data Mining</a>
20                 </li>
21             </ul>
22         </div>
23     </div>
24 
```

```

26      <li>
27          <a href="related_works.html">Related Works</a>
28      </li>
29      <li>
30          <a href="about_us.html">About Us</a>
31      </li>
32  </ul>
33  <p id="suggest"> Type an ID, a GENE NAME <br> or what you need and GO
34      !</p>
35  <form action="">
36      <input type="text" value="Rapid Search" onblur="this.value!=this.
37          value?'Rapid Search':this.value;" onfocus="this.select()"
38          onclick="this.value='';">
39      <input type="submit" value="" onclick="alert('Not yet implemented!')"
40          ;">
41  </form>
42
43  <p id="mail">Comments, questions? <br>
44      Email <a href="mailto:vinxlemons@gmail.com">vinxlemons@gmail.com</a>
45      </p>
46  </div>
47  <div class="content">
48      <div>
49          <h2>Search Result</h2>
50      </div>
51      <div class="data_mining">
52          <div class="result">
53              <h2> <?php
54                  $seq = filter_input(INPUT_GET, "sequence", FILTER_SANITIZE_STRING);
55                  $pat = filter_input(INPUT_GET, "pathology", FILTER_SANITIZE_STRING)
56                  ;
57                  if(strlen($seq)<=15) echo "Blast: ".$seq; else echo "Blast: ".
58                      substr($seq,0,15)."...";
59                  echo "</h2>";
60                  echo "</div>";
61                  echo "<div class='blast' >";
62                      $id = fopen("/var/www/UCbase/ncbi-blast-2.2.27+/bin/query.fasta"
63                          , 'w') or die("can't open file");
64                      fwrite($id, ">".$seq."\n".$seq);
65
66                      exec("cd /var/www/UCbase/ncbi-blast-2.2.27+/bin/; ./blastn -db
67                          UCdb -query query.fasta -task blastn-short 2>&1", $output);
68
69                      //mi connetto al database
70                      $con = mysql_connect("localhost", "client", "is1fs");
71                      if (!$con)
72                      {
73                          die('Could not connect: ' . mysql_error());
74                      }
75
76                      mysql_select_db("UCbase", $con);
77
78                      echo "<p>";

```

```

71     $final; //preparo output finale
72     $j=0;
73     $num_uc = 0;
74     $non_uc_count = 0;
75     $swap = 0;
76     for($i=10; $i<count($output); $i++){
77         $yn=false;
78         if(strpos($output[$i], "Lambda") === 0) { $final[$j] =
79             $output[0] ;break;}
80
81         if($i ==10){ $output[$i] = $output[$i]." (" .substr($output
82             [11],11).")";
83             $final[$j] = $output[$i];
84             $j++;
85         }
86         else if($i != 11 and $i!= 12 and $i!= 16){
87             if($i == 18) {
88                 $output[$i] = " Score E Score_pat Pathology";
89             }
90             if($i == 19) {
91                 $output[$i] = "Sequences alignments: (Bits) Value (Level)
92                     Name";
93             }
94
95             if((strpos($output[$i], "Query=") === 0)){
96                 $output[$i] = "<b>Click on the uc link to read specific
97                     Info.<br><br>".$output[$i]."</b>";
98             }
99
100            if((strpos($output[$i], "lcl") === 0)) {
101                $output[$i] = str_replace(" ", " ", $output[$i]);
102                $output[$i] = str_replace(" ", ".", $output[$i]);
103
104                preg_match("/uc\.[0-9]*/", $output[$i],$uc[$num_uc]);
105                $uc[$num_uc] = $uc[$num_uc][0];
106                preg_match("/[0-9]+\.[0-9]+/", $output[$i],$score[
107                    $num_uc]);
108                //echo "uc num ".$num_uc." : ".$score[$num_uc][0];
109
110                $query = "SELECT UC_NAME FROM
111                  (SELECT DISTINCT UC_NAME FROM UC INNER JOIN (SELECT
112                      ENSAMBL_GENE_ID,NAME FROM RELATED_TO INNER JOIN (
113                          SELECT ID,NAME FROM PATHOLOGY WHERE PREORDER >=
114                          (SELECT PREORDER FROM PATHOLOGY WHERE NAME='".$pat.')
115                          AND POSTORDER <=(SELECT POSTORDER FROM PATHOLOGY
116                          WHERE NAME='".$pat.')') ORDER BY LEVEL)AS B WHERE
117                          RELATED_TO.ID = B.ID)AS C WHERE UC.ENSAMBL_GENE_ID =
118                          C.ENSAMBL_GENE_ID) as X
119                          WHERE UC_NAME = '".$uc[$num_uc]."';
120                $result = mysql_query($query);
121                if(!$row = mysql_fetch_array($result, MYSQL_NUM))
122                    {$yn=true; $non_uc[$non_uc_count] = $uc[$num_uc] ;
123                     $non_uc_count++;}

124
125 //collego patologia

```

```

112 $query = "SELECT NAME FROM
113 (SELECT DISTINCT UC_NAME,NAME FROM UC INNER JOIN (SELECT
114     ENSAMBL_GENE_ID,NAME FROM RELATED_TO INNER JOIN (
115         SELECT ID,NAME FROM PATHOLOGY WHERE PREORDER >= (
116             SELECT PREORDER FROM PATHOLOGY WHERE NAME='".$pat."')
117             AND POSTORDER <=(SELECT POSTORDER FROM PATHOLOGY
118             WHERE NAME='".$pat.') ORDER BY LEVEL)AS B WHERE
119             RELATED_TO.ID = B.ID)AS C WHERE UC.ENSAMBL_GENE_ID =
120             C.ENSAMBL_GENE_ID) as X
121 WHERE UC_NAME = '".$uc[$num_uc]."';
122 $result = mysql_query($query);
123 if(($row = mysql_fetch_array($result, MYSQL_NUM)))
124 {
125     $output[$i] = $output[$i].".....<b>".$row[0]."</b>";
126     $pathology[$num_uc] = $row[0];
127 }
128 else { $pathology[$num_uc] = "null"; }
129 if(($row = mysql_fetch_array($result, MYSQL_NUM)))
130     $output[$i] = $output[$i]." and more..";
131
132 $index_row[$num_uc] = $i;
133
134 //echo ".$uc[$num_uc].": ".$pathology[$num_uc].".
135 //echo $index_row[$num_uc].";
136
137 $query = "select level from PATHOLOGY where name = '".
138     $pathology[$num_uc].'";
139 $result = mysql_query($query);
140 if($row = mysql_fetch_array($result, MYSQL_NUM))
141 {
142     $level[$num_uc] = $row[0];
143     list($pre, $post) = explode('<b>', $output[$i]);
144     $output[$i] = $pre.$level[$num_uc]."......<b>".$post;
145 }
146 else
147     $level[$num_uc] = 100;
148 //echo $level[$num_uc]."<br>";
149
150 //ordino anche per patologia rispetto a tutte le
151 //precedenti
152 if(!isset($non_uc))
153     $non_uc[] = "";
154 for($h=$num_uc-1; $h>=0; $h--){
155
156     if($score[$h] == $score[$num_uc] and $level[$h] > $level
157         [$num_uc] and !in_array($uc[$h], $non_uc)){
158         //echo "sono ".$uc[$num_uc]." scambio con ".$uc[$h].
159         //".<br>";
160         $swap = $index_row[$h];
161         $h_final = $h;
162     }
163 }
164 if($swap){

```

```

154         list($uc[$h_final],$uc[$num_uc]) = array($uc[$num_uc],
155                                         $uc[$h_final]);
156         list($level[$h_final],$level[$num_uc]) = array($level[
157                                         $num_uc],$level[$h_final]);
158     }
159     $num_uc++;
160 }
161
162
163     $output[$i] = str_replace(" ", " ", $output[$i]);
164     if(!(strpos($output[$i], "lcl") === 0))
165         $output[$i] = preg_replace('!uc\.[0-9]*!', "<a style='
166                                     color: blue;' target='_blank' href=\"result.php?uc.id
167                                     =\\0\\>\\0</a>", $output[$i]);
168     else
169         $output[$i] = preg_replace('!uc\.[0-9]*!', "<a style='
170                                     color: blue;' href=\"#\\0\\>\\0</a>", $output[$i]);
171     if($i != 10) $output[$i]= $output[$i]. "<br>";
172
173     if($swap){
174         list($final[$swap-$non_uc_count-13],$output[$i]) = array(
175             $output[$i],$final[$swap-$non_uc_count-13]);
176
177         $swap = 0;
178     }
179
180
181 }
182
183 }
184
185
186 //scrivo
187 echo $final[0];
188 echo "</p>";
189 echo "</div>";
190 echo "<div class='blast' id ='special'>";
191 echo "<p> <br>";
192 $comment = false;
193 for($i=3; $i<count($final); $i++){
194 /*
195     if((strpos($final[$i], "lcl") === 0)){
196         echo "sono numero: ".$i;
197     }
198 */
199     if((strpos($final[$i], "Query=") === 0)) echo "<b>Click on the
200         uc link to read specific Info.<br><br></b>";
//grassetto le lettere

```

```

201     $final[$i] = preg_replace('![AGCT]+[AGCT]+!', '<b>\0</b>',
202                               $final[$i]);
203     if((strpos($final[$i], ">lcl") === 0)){
204         preg_match("/uc\.[0-9]+/", $final[$i],$prova);
205         if(in_array($prova[0], $non_uc))
206             {$i +=1; continue;}
207         else{
208             echo "</p></div> ";
209             if($prova[0] == $uc[0])
210                 echo "<h2 style ='margin :10px 0px 10px 20px; width:400
211                               px; text-align: center;'> Unsorted Specific Info</h2>
212                               ";
213             echo "<div class='blast'> <p> <a name='".$prova[0]."'></a>
214                               <b> Click on the link for more info about the uc.</b>
215                               <br><br>";
216         }
217     }
218     else if (strpos($final[$i], "BLAS") === 0){
219         echo "</p></div> <div class='blast'> <p>";
220     }
221     /*
222     else if((strpos($final[$i], "<br>") === 0) and (strpos($final[$i
223                               +1], ">lcl")=== 0) )
224         echo "cazzoooooooooooooo";
225         if($comment)
226             {echo "-->; $comment=false;}
227         else
228             echo "</p> </div> <div class='blast'> <p>";
229     }
230
231     echo "</p>";
232     echo "</div>";
233     mysql_close($con);
234     ?>
235     <div class="section">
236         <a href="#"></a>
237         </div>
238     </div>
239     </div>
240     </div>
241     </body>
242 </html>
243 </body>
244 </html>
245
246

```

## C.25 Blastn.php

```
1  <html>
2  <body>
3
4
5  <!DOCTYPE html>
6  <html>
7  <head>
8      <meta charset="UTF-8">
9      <title>UCbase - UC Data Mining</title>
10     <link rel="icon" href="http://localhost/appswebsitetemplate/favicon.ico"
11         />
12     <link rel="stylesheet" href="css/style.css" type="text/css">
13 </head>
14 <body>
15     <div class="page">
16         <div class="sidebar">
17             <div id="logo">
18                 <a href="index.html"></a>
19             </div>
20             <ul class="navigation">
21                 <li>
22                     <a href="index.html">Home</a>
23                 </li>
24                 <li class="selected">
25                     <a href="uc_data_mining.php">UC Data Mining</a>
26                 </li>
27                 <a href="related_works.html">Related Works</a>
28             </li>
29             <li>
30                 <a href="about_us.html">About Us</a>
31             </li>
32         </ul>
33         <p id="suggest"> Type an ID, a GENE NAME <br> or what you need and GO
34             !</p>
35         <form action="">
36             <input type="text" value="Rapid Search" onblur="this.value!=this.
37                 value?'Rapid Search':this.value;" onfocus="this.select()"
38                 onclick="this.value='';">
39             <input type="submit" value="" onclick="alert('Not yet implemented!')"
40                 ;">
41         </form>
42
43         <p id="mail">Comments, questions? <br>
44             Email <a href="mailto:vinxlemons@gmail.com">vinxlemons@gmail.com</a>
45             </p>
46     </div>
47     <div class="content">
48         <div>
49             <h2>Search Result</h2>
50         </div>
```

```

46 <div class="data_mining">
47   <div class="result">
48     <h2> <?php
49     $seq = filter_input(INPUT_GET, "sequence", FILTER_SANITIZE_STRING);
50     if(strlen($seq)<=15) echo "Blast: ".$seq; else echo "Blast: ".
51       substr($seq,0,15)."...";
52     echo "</h2>";
53     echo "</div>";
54     echo "<div class='blast' >";
55     $id = fopen("/var/www/UCbase/ncbi-blast-2.2.27+/bin/query.fasta"
56       , 'w') or die("can't open file");
57     fwrite($id, ">".$seq."\n".$seq);
58
59     exec("cd /var/www/UCbase/ncbi-blast-2.2.27+/bin/; ./blastn -db
60           UCdb -query query.fasta -task blastn-short 2>&1", $output);
61     echo "<p>";
62
63     $final; //preparo output finale
64     $j=0;
65     $num_uc = 0;
66     for($i=10; $i<count($output); $i++){
67
68       if(strpos($output[$i], "Lambda") === 0) { $final[$j] =
69         $output[0] ;break;}
70
71       if($i ==10){ $output[$i] = $output[$i]." (" .substr($output
72         [11],11).");
73       $final[$j] = $output[$i];
74       $j++;
75     }
76     else if($i != 11 and $i!= 12 and $i!= 16){
77       if((strpos($output[$i], "lcl") === 0)){
78         preg_match("/uc\.[0-9]*/", $output[$i],$uc[$num_uc]);
79         $uc[$num_uc] = $uc[$num_uc][0];
80         $num_uc++;
81         $output[$i] = str_replace(" ", ".", $output[$i]);
82       }
83
84       if(strpos($output[$i], "Query=") === 0)){
85         $output[$i] = "<b>Click on the uc link to read specific
86           Info.<br><br>".$output[$i]."</b>";
87       }
88       if($i != 11)$output[$i] = str_replace(" ", " ", $output
89         [$i]);
89       if(!(strpos($output[$i], "lcl") === 0))
90         $output[$i] = preg_replace('!uc\.[0-9]*!', "<a style='
91           color: blue;' target='_blank' href=\"result.php?uc.id
92             =\\\"$>\\\"$</a>",$output[$i]);
93     else
94       $output[$i] = preg_replace('!uc\.[0-9]*!', "<a style='
95           color: blue;' href=\"#\\\"$>\\\"$</a>",$output[$i]);
96     if($i != 10) $final[$j]= $output[$i]."<br>";
97     else $final[$j] = $output[$i];
98     $j++;

```

```

90         }
91     }
92
93     echo $final[0];
94     echo "</p>";
95     echo "</div>";
96     echo "<div class='blast'>";
97     echo "<p> <br>";
98
99
100    for($i=3; $i<count($final); $i++){
101    /*
102        if((strpos($final[$i], "lcl") === 0)){
103            echo "sono numero: ".$i;
104        }
105    */
106        $final[$i] = preg_replace('![AGCT]+[AGCT]!', '<b>\0</b>',
107                                $final[$i]);
108
109        if((strpos($final[$i], ">lcl") === 0)){
110            preg_match("/uc\.[0-9]+/", $final[$i],$prova);
111            if($prova[0] == $uc[0]){
112                echo "<form style='width: 360px;' action='blastpat.php'
113                    method='GET'> <input type='text' name='sequence' value='
114                        '.$seq.'" style='display:none;'>";
115                echo "<span>&ampnbsp&ampnbsp&ampnbspFilter with pathology: </span> <
116                    select name='pathology'>";
117                $con = mysql_connect("localhost","client","is1fs");
118                if (!$con)
119                {
120                    die('Could not connect: ' . mysql_error());
121                }
122
123                mysql_select_db("UCbase", $con);
124                $query = "SELECT DISTINCT NAME FROM PATHOLOGY ";
125
126                $result = mysql_query($query);
127
128                while($row = mysql_fetch_array($result, MYSQL_NUM))
129                {
130                    if($row[0]!=""){
131                        if($row[0] == "disease")
132                            echo '<option value="'. $row[0]. '" selected="selected">
133                                .' $row[0]. '</option>';
134                        else
135                            echo '<option value="'. $row[0]. '">.' $row[0]. '</option>
136                                ';
137                    }
138                }
139                echo "<input class ='submit' type='submit' value=' '></select
140                    ></form>";
141            }
142            echo "</p></div> ";
143            if($prova[0] == $uc[0])

```

```

137         echo "<h2 style ='margin :10px 0px 10px 20px; width:400px;
138             text-align: center;'> Specific Info</h2>";
139         echo "<div class='blast'> <p> <a name='".$prova[0]."'></a><b>
140             Click on the link for more info about the uc.</b><br><br>"'
141             ;
142     }
143     /*
144     else if (strpos($final[$i], "BLAS") === 0){
145         echo "</p></div> <div class='blast'> <p>";
146     }
147     */
148     else if((strpos($final[$i], "<br>") === 0) and (strpos($final[$i
149         +1], ">lcl")=== 0) ) {
150         echo "cazzooooooooooooo";
151         if($comment)
152             {echo "-->; $comment=false;}
153         else
154             echo "</p> </div> <div class='blast'> <p>";
155     }
156     */
157     echo "</p>";
158     echo "</div>";
159     ?>
160     <div class="section">
161         <a href="#"></a>
162     </div>
163     </div>
164     </div>
165     </body>
166 </html>
167
168
169 </body>
170 </html>
```

## C.26 About\_us.html

```

1 <!DOCTYPE html>
2 <html>
3 <head>
4     <meta charset="UTF-8">
5     <title>UCbase - About Us</title>
6     <link rel="icon" href="http://localhost/appssite/template/favicon.ico" /
7         >
8     <link rel="stylesheet" href="css/style.css" type="text/css">
9 </head>
10 <body>
```

```

11 <div class="page">
12   <div class="sidebar">
13     <div id="logo">
14       <a href="index.html"></a>
15     </div>
16     <ul class="navigation">
17       <li>
18         <a href="index.html">Home</a>
19       </li>
20       <li>
21         <a href="uc_data_mining.php">UC Data Mining</a>
22       </li>
23       <li>
24         <a href="related_works.html">Related Works</a>
25       </li>
26       <li class="selected">
27         <a href="about_us.html">About Us</a>
28       </li>
29     </ul>
30     <p id="suggest"> Type an ID, a GENE NAME <br> or what you need and GO
31       !</p>
32     <form action="">
33       <input type="text" value="Rapid Search" onblur="this.value!=this.
34         value?'Rapid Search':this.value;" onfocus="this.select()"
35         onclick="this.value='';">
36       <input type="submit" value="" onclick="alert('Not yet implemented
37         !');">
38     </form>
39
40   <p id="mail">Comments, questions? <br>
41     Email <a href="mailto:vinxlemons@gmail.com">vinxlemons@gmail.com</a>
42   </p>
43 </div>
44 <div class="content">
45   <div>
46     <h2>About Us</h2>
47     <h3>Information Systems Group - University of Modena and Reggio
48       Emilia</h3>
49     <p>The work of the ISGroup, here at the Computer Engineering
50       Department (DII) of the University of Modena and Reggio Emilia,
51       mainly focuses on the design and development of new systems,
52       algorithms and data structures for the access and management of
53       Information.</p>
54     <h3>Join Our work!</h3>
55     <p>If you're experiencing issues and concerns about our work, please
56       email us or help us to create a better work working together
57       with our research group.</p>
58     <p>Design version: beta.<br>Code version: beta.<br><br>
59     <h3>Contact Us</h3>
60     <span><b>Telephone Nos. :</b> -----, -----</span>
61     <span><b>Email :</b> vinxelmons@gmail.com</span>
62     <span><b>Street Address :</b> ----- Modena, Italy </span>
63   </div>
64   <div class="links">

```

```

53      <div class="date">
54          <span>13-02</span>
55          <span>2013</span>
56      </div>
57      <div>
58          <h2>Useful links</h2>
59          <div>
60              border
61          </div>
62      </div>
63      <ul>
64          <li>
65              <a href="#"></a>
66          </li>
67          <li>
68              <a href="#"></a>
69          </li>
70          <li>
71              <a href="#"></a>
72          </li>
73          <li>
74              <a href="#"></a>
75          </li>
76          </ul>
77      </div>
78  </div>
81  </body>
82 </html>

```

## C.27 Style.css

```

1  /* Website template by freewebsitetemplates.com */
2  /*----- Layout styles -----*/
3  body{
4      margin: 0;
5      background: url(..../images/bg-body.gif) repeat-x #205eb0;
6  }
7  .page{
8      width: 950px;
9      padding: 0 5px;
10     margin: 0 auto;
11     position: relative;
12 }
13
14
15 /*----- Sidebar -----*/
16 .sidebar {
17     float:left;
18     padding:80px 0 0;
19     width:248px;

```

```

20  }
21 .sidebar div#logo {
22   background:url(..../images/interface.png) no-repeat 0 -142px;
23   height:200px;
24   left:-85px;
25   padding:60px 0 0;
26   position:absolute;
27   top:0;
28   width:383px;
29 }
30 .sidebar div#logo a {
31   display:block;
32   height:57px;
33   margin:0 0 0 90px;
34   outline:none;
35   width:230px;
36 }
37 .sidebar div#logo a img {
38   border:0;
39 }
40 .sidebar ul.navigation {
41   background:url(..../images/border-dashed.gif) repeat-x left bottom;
42   list-style:none;
43   margin:80px 0 0 20px;
44   padding:0 0 10px;
45   text-align:center;
46   width:192px;
47   position:relative;
48 }
49 .sidebar ul.navigation li {
50   height:36px;
51   margin:0 0 10px;
52 }
53 .sidebar ul.navigation li.selected {
54   background:url(..../images/interface.gif) no-repeat;
55 }
56 .sidebar ul.navigation li.selected a,.sidebar ul.navigation li a:hover {
57   color:#d3d4de;
58 }
59 .sidebar ul.navigation li a {
60   color:#ffffff;
61   font-family:Times New Roman;
62   font-size:20px;
63   font-weight:700;
64   line-height:36px;
65   outline:none;
66   text-decoration:none;
67   text-shadow:0 2px 0 #00007c;
68 }
69 .sidebar form {
70   background:url(..../images/border-dashed.gif) repeat-x left bottom;
71   margin:15px 0 0 20px;
72   overflow:hidden;
73   padding:0 0 20px;

```

```

74     width:192px;
75 }
76 .sidebar input:first-child {
77   background:none;
78   background-color:#dadaed;
79   border:0;
80   color:#1325f0;
81   cursor:auto;
82   float:left;
83   font-family:helvetica;
84   font-size:12px;
85   font-style:italic;
86   font-weight:700;
87   height:auto;
88   margin:0;
89   padding:5px 9px;
90   width:132px;
91 }
92 .sidebar input {
93   background:url(..../images/interface.gif) no-repeat 0 -56px;
94   border:0;
95   cursor:pointer;
96   float:right;
97   height:25px;
98   width:30px;
99 }
100
101 .sidebar #suggest {
102   color:#ffffff;
103   font-family:Times New Roman;
104   font-size:12px;
105   margin-left:45px;
106   width:160px;
107 }
108
109 .sidebar #mail {
110   color:#ffffff;
111   font-family:Times New Roman;
112   font-size:12px;
113   margin:135px 0 0 20px;
114   text-align:center;
115   width:200px;
116 }
117
118 /*----- Content -----*/
119 .content{
120   float:left;
121   padding:55px 0 24px;
122   width:702px;
123 }
124 .content .article{
125   margin:0 0 0 24px;
126 }
127 .content .article h2{

```

```

128 /* color:#FEDD14;*/
129 color:#FFFFFF;
130 font-family:Times New Roman;
131 font-size:35px;
132 margin:0;
133 text-shadow:0 2px 10px #00007c;
134 }
135 .content .article p{
136 color:#fff;
137 font-family:Times New Roman;
138 font-size:16px;
139 font-style:normal;
140 line-height:21px;
141 margin:10px 0 15px;
142 text-align:justify;
143 }
144
145 .content .blog{
146 background-color:#dce3e8;
147 margin:24px 0 0 24px;
148 padding:18px 0 0 ;
149 }
150 .content .blog .date,.content .links .date{
151 background:url(..../images/interface.png) no-repeat;
152 float:left;
153 height:54px;
154 margin:0 0 0 -24px;
155 width:90px;
156 }
157 .content .blog .date span:first-child,.content .links .date span:first-child{
158 font-size:14px;
159 margin:0 0 0 35px;
160 padding:2px 0 0 ;
161 }
162 .content .blog .date span,.content .links .date span{
163 color:#ffecc2;
164 display:block;
165 font-family:Times New Roman;
166 font-size:12px;
167 margin:0 0 0 45px;
168 }
169 .content .blog div,.content .links div{
170 float:left;
171 margin:0;
172 }
173 .content .blog div h2,.content .links div h2{
174 color:#205eb0;
175 float:left;
176 font-family:Times New Roman;
177 font-size:34px;
178 margin:4px 0 0 18px;
179 text-shadow:0 1px 0 #fff7e4;
180 }
181 .content .blog div div,.content .links div div{

```

```

182    background:url(..../images/border-dashed2.gif) repeat-x center center;
183    margin:15px 0 0 10px;
184    text-indent:-9999px;
185    width:210px;
186  }
187 .content .blog ul{
188   background:url(..../images/border-dashed2.gif) repeat-x center bottom;
189   clear:both;
190   list-style:none;
191   margin:0 0 85px;
192   overflow:hidden;
193   padding:15px 0 25px;
194   width:527px;
195  }
196 .content .blog ul li{
197   margin:50px 0 0;
198   overflow:hidden;
199  }
200 .content .blog ul li div{
201   background:none;
202   float:none;
203   margin:0;
204   overflow:hidden;
205   text-indent:0;
206  }
207 .content .blog ul li div div:first-child{
208   background:none;
209   margin:0;
210   width:84px;
211  }
212 /*.content .blog ul li div div:first-child a:first-child{
213   background:url(..../images/shadow.jpg) no-repeat center bottom;
214   height:100px;
215   width:84px;
216 }*/
217 .content .blog ul li div div:first-child a{
218   background-color:#989eb3;
219   color:#000000;
220   display:block;
221   font-family:Times New Roman;
222   font-size:16px;
223   font-style:normal;
224   height:83px;
225   line-height:83px;
226   outline:none;
227   text-align:center;
228   text-decoration:none;
229   text-indent:0;
230   text-shadow:none;
231   width:84px;
232  }
233 .content .blog ul li div div{
234   float:left;
235   margin:0 0 0 18px;

```

```

236     width:425px;
237 }
238 .content .blog ul li div div h3{
239     color:#0a578a;
240     font-family:Times New Roman;
241     font-size:22px;
242     font-style:normal;
243     font-weight:500;
244     margin:0;
245     text-transform:Capitalize;
246     text-shadow: none;
247 }
248 .content .blog ul li div div p{
249     color:#000000;
250     font-family:Times New Roman;
251     font-size:16px;
252     font-style:normal;
253     margin:12px 0 35px;
254     text-align:justify;
255 }
256 .content .blog ul li div div a{
257     color:#000000;
258     float:right;
259     font-family:Times New Roman;
260     font-size:14px;
261     font-style:italic;
262     outline:none;
263     text-decoration:none;
264     text-shadow:0 1px 0 #f6f8f5;
265 }
266 .content .blog ul li div div a:hover{
267     color:#0a578a;
268     cursor: hand;
269     cursor: pointer;
270 }
271 .content .blog .section{
272     float:none;
273     height:60px;
274     margin:0 0 0 85px;
275     width:527px;
276 }
277 .content .blog .section a:first-child{
278     height:48px;
279     outline:none;
280     background:url(../images/interface.png) no-repeat 0 -74px;
281     display: block;
282     width:48px;
283 }
284 .content .blog .section a{
285     color:#e8cb93;
286     margin-left: 245px;
287     margin-top: 10px;
288     font-family:Times New Roman;
289     font-style:italic;

```

```
290    line-height:84px;
291    outline:none;
292    text-decoration:none;
293 }
294 .content .blog .section a:hover{
295   color:#F89F1F;
296 }
297 .content .blog .section .paging{
298   background:none;
299   margin:39px 0 0 225px;
300   width:144px;
301 }
302 .content .blog .section .paging a{
303   background-color:#FEC97B;
304   color:#FFFBF5;
305   font-style:normal;
306   font-weight:700;
307   height:18px;
308   line-height:17px;
309   margin:0 0 0 12px;
310   outline:none;
311   text-align:center;
312   text-indent:0;
313   width:66px;
314 }
315 .content .blog .section .paging a:hover{
316   background-color:#e6cc93;
317   color:#ffeabf;
318 }
319 .content div{
320   margin:0 0 0 24px;
321 }
322 .content div h2{
323   color:#ffffff;
324   font-family:'Times New Roman';
325   font-size:35px;
326   margin:0;
327   text-shadow:#00007C 0 2px 10px;
328 }
329 .content div h3{
330   color:#ffffff;
331   font-family:Times New Roman;
332   font-size:18px;
333   font-style:italic;
334   font-weight:600;
335   margin:10px 0 0;
336   text-transform:uppercase;
337   text-shadow:#00007C 0 2px 0;
338 }
339 .content div p{
340   color:#FFFFFF;
341   font-family:'Times New Roman';
342   font-size:16px;
343   font-style:normal;
```

```
344 line-height:21px;
345 margin:10px 0 15px;
346 text-align:justify;
347 }
348 .content div p a:hover{
349   color:#FEDD14;
350 }
351 .content div p a{
352   color:#fff;
353   outline:none;
354 }
355 .content .links ul{
356   background:none;
357   clear:both;
358   list-style:none;
359   margin:0;
360   overflow:hidden;
361   padding:0px 0 50px;
362 }
363 .content .links ul li{
364   float:left;
365   margin:0 0 0 63px;
366   width:84px;
367 }
368
369 .content .links ul li a{
370   color:#ffebc4;
371   display:block;
372   font-family:Times New Roman;
373   font-size:16px;
374   font-style:normal;
375   height:24px;
376   line-height:25px;
377   outline:none;
378   text-align:center;
379   text-decoration:none;
380   width:84px;
381   margin-bottom: 20px;
382 }
383 .content .data_mining h2{
384   color:#ffffff;
385   float:left;
386   font-family:Times New Roman;
387   font-size:34px;
388   margin:4px 0 84px;
389   padding:5px 120px;
390   background-color: #205EB0;
391   text-shadow:#00007C 0 2px 10px;
392 }
393 }
394 .content .data_mining select#tre{
395 width: 250px;
396 }
397 .content .data_mining select{
```

```

398 height:28px;
399 }
400 .content .data_mining input,select{
401 background:none;
402 float: left;
403 background-color:#ffffff;
404 border: 1px solid #5b76e3;
405 color:#205EB0;
406 cursor:auto;
407 font-family:helvetica;
408 font-size:14px;
409 font-style:italic;
410 font-weight:700;
411 height:20px;
412 margin:6px;
413 padding:2px 0px 5px 9px;
414 width:132px;
415 }
416
417 .content .data_mining input#second{
418 width:560px;
419 }
420
421 .content .data_mining input.submit{
422 background:url(..../images/interface.gif) no-repeat 0 -56px;
423 border:0;
424 cursor:pointer;
425 float:right;
426 margin-top: 8px;
427 height:25px;
428 width:30px;
429 }
430
431 .content .data_mining span{
432 color:#000000;
433 display: block;
434 font-family:"Times New Roman",Georgia,Serif;
435 font-size:17px;
436 # font-style:bold;
437 font-weight: bold;
438 margin:10px 0 10px;
439 float: left;
440 }
441
442 .content .data_mining h2#second{
443 margin:20px 0 10px 84px;
444 padding:5px 90px;
445 }
446
447 .content .data_mining div{
448 margin:0;
449 overflow:hidden;
450 }
451 }

```

```

452 .content .data_mining div div{
453   background:url(..../images/border-dashed2.gif) repeat-x scroll center center
454   transparent;
455   float:left;
456   margin:15px 0 0 10px;
457   text-indent:-99999px;
458   width:325px;
459 }
460 .content div.result{
461   background-color: #205eb0;
462   width: 640px;
463   margin-left: 20px;
464 }
465 .content div.blast#special{
466   overflow:scroll;
467 }
468 .content div.blast{
469   background-color: #a1c5df;
470   width: 640px;
471   margin-left: 20px;
472   margin-top: 10px;
473 }
474
475 .content div.result h2{
476   margin-left: 10px;
477   padding: 0;
478 }
479 .content .data_mining ul{
480   background-color: #c0c8ed;
481   list-style:none;
482   margin: 8px 20px 0px 20px;
483   overflow:hidden;
484   width:630px;
485   height: 41px;
486   padding-left: 10px;
487 }
488 .content .data_mining .section#result{
489   background:none;
490   float:none;
491   margin:10px 0 0 20px;
492   padding:0 0 15px;
493   overflow-x:scroll;
494   overflow-y:scroll;
495   height: 500px;
496 }
497
498 .content .data_mining .section{
499   background:none;
500   float:none;
501   margin:0 0 0 87px;
502   padding:0 0 15px;
503
504 }

```

```

505 .content .data_mining .section a#end{
506   height:48px;
507   margin:15px 0 0 225px;
508   background:url(..../images/interface.png) no-repeat 0 -74px;
509   display: block;
510   width:48px;
511 }
512 .content .appspage{
513   background-color:#FFECC1;
514   margin:24px 0 324px 45px;
515   padding:18px 0 33px;
516 }
517 .content .appspage div h2{
518   color:#F27109;
519   float:left;
520   font-family:Times New Roman;
521   font-size:34px;
522   margin:4px 0 0 18px;
523   text-shadow:0 1px 0 #FFF7E4;
524 }
525 .content .appspage div div{
526   background:url(..../images/border-dashed2.gif) repeat-x scroll center center
      transparent;
527   float:left;
528   margin:15px 0 0 10px;
529   text-indent:-99999px;
530   width:350px;
531 }
532 .content .appspage .previous{
533   background:url(..../images/interface.png) no-repeat;
534   float:left;
535   height:54px;
536   margin:0 0 0 -24px;
537   text-align:center;
538   width:90px;
539 }
540 .content .appspage .previous a{
541   color:#FFECC2;
542   font-family:Times New Roman;
543   font-size:14px;
544   line-height:35px;
545   outline:none;
546   text-align:center;
547   text-decoration:none;
548 }
549 .content .appspage .section{
550   background:url(..../images/border-dashed2.gif) repeat-x scroll center bottom
      transparent;
551   margin:0 0 0 63px;
552   overflow:hidden;
553   padding:0 0 99px;
554   width:535px;
555 }
556 .content .appspage .section div:first-child{

```

```

557     float:left;
558     width:84px;
559 }
560 /*.content .appspage .section div:first-child a:first-child{
561     background:url(..../images/shadow.jpg) no-repeat scroll center bottom
562         transparent;
563     height:100px;
564     margin:0;
565     width:84px;*/
566 }
566 .content .appspage .section div:first-child a{
567     background-color:#F89F1F;
568     color:#FFECC1;
569     display:block;
570     font-family:Times New Roman;
571     font-size:16px;
572     font-style:normal;
573     height:24px;
574     line-height:25px;
575     margin:5px 0 0;
576     outline:none;
577     text-align:center;
578     text-decoration:none;
579     text-indent:0;
580     text-shadow:none;
581     width:84px;
582 }
583 .content .appspage .section div{
584     background:none;
585     margin:0 0 0 20px;
586     text-indent:0;
587     width:410px;
588 }
589 .content .appspage .section div h3{
590     color:#D13D01;
591     font-family:Times New Roman;
592     font-size:24px;
593     font-style:normal;
594     font-weight:400;
595     margin:0;
596     text-transform:none;
597 }
598 .content .appspage .section div p{
599     color:#F69F1C;
600     font-family:Times New Roman;
601     font-size:12px;
602     font-style:normal;
603     line-height:21px;
604     margin:12px 0 15px;
605     text-align:justify;
606 }
607 .content .appspage .section div ul{
608     list-style:none;
609     margin:0;

```

```

610     padding:0;
611 }
612 .content .appspage .section div ul li{
613     float:left;
614     margin:0 0 0 7px;
615 }
616 .content .appspage .section div ul li a{
617     outline:none;
618 }
619 .content div span{
620     color:#fff;
621     display:block;
622     font-family:Times New Roman;
623     font-size:16px;
624     font-style:italic;
625     margin:0 0 10px;
626 }
627 .content div span b{
628     color:#ffd215;
629     font-family:Times New Roman;
630     font-size:16px;
631     font-style:italic;
632 }
633 .content .article a:hover,.content .blog ul li div div:first-child a:hover,.
    content .links ul li a:hover,.content .data_mining ul li a:hover,.content
    .appspage .previous a:hover{
634     color:#fff;
635 }
636 .content .blog ul li:first-child,.content .blog .section .paging a:first-child
    {
637     margin:0;
638     background:none;
639 }
640 .content .data_mining ul li:first-child,.content .appspage .section div ul li:
    first-child{
641     margin: 0;
642 }
643 .content .blog ul li div div:first-child a:first-child img,.content .blog .
    section a:first-child img,.content .links ul li a:first-child img,.
    content .data_mining ul li a:first-child img,.content .data_mining .
    section a img,.content .appspage .section div:first-child a img,.content
    .appspage .section div ul li a img{
644     border:0;
645 }
646 .content .blog ul li div div:first-child a:first-child img:hover,.content .
    links ul li a:first-child img:hover,.content .data_mining ul li a:first-
    child img:hover,.content .appspage .section div:first-child a img:hover{
647     opacity:0.8;
648 }
649 .content .links,.content .data_mining{
650     background-color:#DCE3E8;
651     margin:24px 0 0 24px;
652     padding:18px 0 0;
653 }

```

```
654
655 table#result
656 {
657 margin-top: 20px;
658 font-family:"Trebuchet MS", Arial, Helvetica, sans-serif;
659 width:640px;
660 margin-left: 20px;
661 margin-bottom: 20px;
662 border-collapse:collapse;
663 table-layout:fixed;
664 }
665 #result td, #result th
666 {
667 font-size:1em;
668 border:1px solid #ffffff;
669 padding:3px 7px 2px 7px;
670 }
671
672 #result tr.alt{
673 background-color:#cbd5ff;
674 }
675
676 #result td.one
677 {
678 width: 170px;
679 font-size:18px;
680 font-weight:bold;
681 text-align:right;
682 padding-top:5px;
683 padding-bottom:4px;
684 background-color:#205EB0;
685 color:#ffffff;
686 }
687
688 a.button{
689 width: 115px;
690 float: right;
691 text-align: center;
692 background-color:#205EB0;
693 color: #ffffff;
694 text-decoration: none;
695 margin-top:1px;
696 }
697
698 a.button:hover{
699 background-color:#3768f9;
700 }
701
702 table#big
703 {
704 font-family:Lucidatypewriter, monospace;
705 width:100%;
706 border-collapse:collapse;
707 }
```

```
708 #big td, #big th
709 {
710 font-size:12px;
711 border:1px solid #205EB0;
712 padding:3px 7px 2px 7px;
713 }
714 #big th
715 {
716 font-family:"Trebuchet MS", Arial, Helvetica, sans-serif;
717 font-size:14px;
718 text-align:left;
719 padding-top:5px;
720 padding-bottom:4px;
721 background-color:#205EB0;
722 color:#ffffff;
723 }
724 #big tr.alt td
725 {
726 color:#000000;
727 background-color:#e7eaf9;
728 }
729
730 div.blast p
731 {
732 color:black;
733 font-family:monospace;
734 font-size: 12px;
735 margin-left:10px;
736 }
```

# Bibliografia

- [1] Isgroup. Proposed Thesis. URL:  
<http://www.isgroup.unimo.it/thesis.asp>.
- [2] C. J. Date and Hugh Darwen. *A Guide to the SQL Standard*. Addison-Wesley, 4th edition, 2002.
- [3] Philipp K. Janert. *Data Analysis with Open Source Tools*. O'Reilly Media.
- [4] GNU's Not Unix. The R Project for Statistical Computing. URL:  
<http://www.r-project.org/>.
- [5] Microsoft. Microsoft Excel. URL:  
<http://office.microsoft.com/it-it/excel/>.
- [6] D. Beneventano, S. Bergamaschi, F. Guerra, and M. Vincini. *Progetto di basi di dati relazionali*. Pitagora Editrice s.r.l., 4th edition, 2007.
- [7] McKusick-Nathans Institute. Online Mendelian Inheritance in Manl. URL: <http://www.omim.org/>.
- [8] EBI and EMBL. Ensembl. URL: <http://www.ensembl.org/>.
- [9] Bioconductor. BiomaRt. URL:  
<http://www.bioconductor.org/packages/release/bioc/html/biomaRt.html>.
- [10] Arek Kasprzyk. BioMart: driving a paradigm change in biological data management. *Oxford University Press*, 2011.
- [11] Apache Software Foundation. openoffice. URL:  
<http://www.openoffice.org/>.
- [12] Oracle Corporation. MySQL. URL: <http://www.mysql.it/>.
- [13] The PHP Group. PHP: Hypertext Preprocessor. URL:  
<http://www.php.net/>.
- [14] IHTSDO. SNOMED CT. URL: <http://www.ihtsdo.org/snomed-ct/>.

- [15] The Open Biological and Biomedical Ontologies. OBO Foundry. URL: <http://www.obofoundry.org/>.
- [16] G. Cabri and F. Zambonelli. *Programmazione a oggetti in Java: dai fondamenti a Internet*. Pitagora Editrice s.r.l., 2003.
- [17] James Gosling and Sun Microsystems. Java. URL: <http://www.java.com/>.
- [18] Eclipse Foundation. Eclipse. URL: <http://www.eclipse.org/>.
- [19] Nomi Harris, John Day-Richter, Chris Mungall, Amina Abdulla, Nomi Harris, and Jennifer Deegan. Obo-Edit. URL: <http://oboedit.org/>.
- [20] Pavel Zezula, Federica Mandreoli, and Riccardo Martoglia. Tree Signatures and Unordered XML Pattern Matching. 2004.
- [21] Taccioli C, Fabbri E, Visone R, Volinia S, Calin GA, Fong LY, Gambari R, Bottoni A, Acunzo M, Hagan J, Iorio MV, Piovan C, Romano G, and Croce CM. UCbase & miRfunc: a database of ultraconserved sequences and microRNA function. *Database issue journal*, 2008.
- [22] Gill Bejerano, Michael Pheasant, Igor Makunin, Stuart Stephen, W. James Kent, John S. Mattick, and David Haussler. Ultraconserved elements in the human genome. *Science Magazine*, 2004.
- [23] Wikimedia Foundation. Wikipedia. URL: <http://www.wikipedia.org>.
- [24] Istituto dell'Enciclopedia Italiana. Treccani. URL: <http://www.treccani.it/>.
- [25] Sam Griffiths-Jones. The microRNA Registry. *Database issue journal*, 2003.
- [26] R01RR025342. Human disease ontology. URL: <http://diseaseontology.sourceforge.net/>.
- [27] National Center for Biotechnology Information. NCBI. URL: <http://www.ncbi.nlm.nih.gov/>.
- [28] National Center for Biotechnology Information. BLAST. URL: <http://blast.ncbi.nlm.nih.gov/Blast.cgi>.