

Master : STIC pour la Santé
Spécialité : « Bioinformatique, Connaissances, Données »

GMIN206 - Projet : TP XML

AGRET Clément & CHAKIACHVILI Marc

Présentation du projet

Le projet consiste à réaliser un programme en java qui récupère les fichiers XML d'un gène orthologue d'une base de données dans différents formats depuis le NCBI. Ce programme doit être capable de manipuler ces fichiers, effectuer quelques opérations sur leurs infos contenues et créer des informations qui ne sont pas accessibles, ou proposer par le site.

Notre programme a pour but de prendre en entrée le nom du fichier xml suivi d'un underscore ("_") et du format de celui-ci (BioSeq, GBCS ou Tiny(fasta)). Ex : séquence_BioSeq.

Notre application aligne les séquences qu'elle a récupérées grâce à Clustal W et elle renvoie la plus longue sous-séquence commune aux différentes séquences prises en entrée.

Réalisation

Traitement

Demande biologique : Retrouver et compter les exons présents dans les différentes séquences, retourner en sortie un fichier XML contenant uniquement les exons et leurs sous-séquences associées et un alignement global des séquences de ce gène.

DTD du fichier de sortie :

Listing 1 – Fichier de sortie

```
<output>
  <gene> </gene>
  <organisms>
    <organism>
      <GE></GE>
    </organism>
  </organisms>
  <alignment method="clustalW"> </alignment>
  <exons>
    <exon>
      <exon_location> </exon_location>
      <exon_sequence> </exon_sequence>
    </exon>
  </exons>
</output>
```

UML

