

Construction d'une ontologie par programmation basée sur l'exemple

Quentin Riché-Piotaix

Université de Poitiers - LIAS ENSMA - CHU de Poitiers

4 mai 2017

Introduction

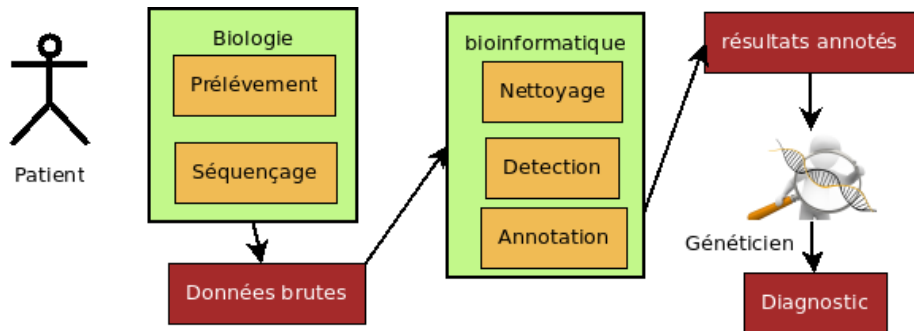
Encadrement

- Patrick Girard, LIAS - Université de Poitiers
- Ladjel Bellatreche, LIAS
- Frédéric Bilan, CHU de Poitiers - Université de Poitiers

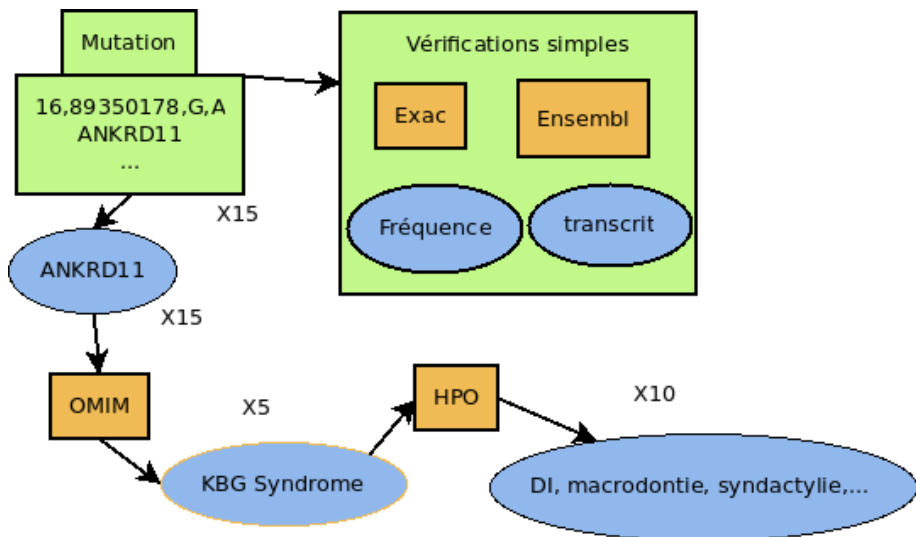
Sujet

Exploration à base ontologique de données issues de patients atteints de maladies génétiques rares.

Contexte : La Génétique



Scénario exemple



→ Besoin de rassembler les données

Contexte : La Génétique

Les bases de données génétiques

- sont nombreuses
- sont variées
- sont spécialisées
- se recoupent
- sont sémantiquement inconstantes
- sont disponibles en fichiers semi structurés

Module *ad hoc* de transfert ?

Ontologie ?

Contexte : La Génétique

Spécificité du domaine

Constante évolution, ajout et modification de savoir réguliers et fréquents

→ Ontologie globale impossible à maintenir

Spécificité du domaine

Les généticiens n'ont pas besoin d'une ontologie couvrant tout le domaine :
Poitiers s'intéresse particulièrement à la déficience intellectuelle

→ Créations de petites ontologies spécifiques

Spécificité du domaine

Les généticiens n'ont pas possibilité de travailler avec un ingénieur
spécialisé

→ Création de sa propre ontologie par le généticien

Contexte : les généticiens

Des experts...

- Formation médicale ou scientifique (Médecine ou doctorat en génétique)
- Différentes sous-spécialités : génétique moléculaire, génétique clinique, cytogénétique,...

...novices

- Pas de formation informatique
- Utilisation fréquente de base de données

Autres contraintes

- Pas de formation possible
- Confiance totale dans le résultat nécessaire

2 Contraintes

- Domaine très changeant
- Utilisateurs novices

2 Atouts

- Données semi structurées
- Utilisateurs experts

Sujet

Sujet

Exploration à base ontologique de données issues de patients atteints de maladies génétiques rares.

Notre approche

Peut on permettre à un utilisateur expert de créer sa propre ontologie à partir des données ?

Ontology learning

Processus de création ou d'augmentation d'une ontologie de façon automatique ou semi-automatique

Approches

- Regrouper des ontologies
- Spécialiser une ontologie de haut niveau
- Construire une ontologie depuis le départ

Bibliographie

- Principalement approches machine learning
- Quelques approches participatives à grande échelles

Limites des approches actuelles

- Interfaces faites pour des ingénieurs base de données
 - Réponse habituelle : proximité avec des standards du domaine (SPARQL)
 - Réponse non applicable pour nous
- Dépendance au format : malformation rendant complexe l'extraction d'information de fichiers semi-structurés
 - Aide de l'utilisateur et parsing intelligent
- Duplication d'entité lorsqu'elles sont présentes dans deux sources différentes
 - Nous allons être très concernés
- Évaluation de l'ontologie produite est compliquée
 - Solution classique : revue manuelle par un expert
 - Nous ne sommes pas concernés puisque l'ontologie est construite par l'expert

Hypothèses

Observations

- Les généticiens utilisent plusieurs bases de données
- Ils font des requêtes complexes en faisant plusieurs requêtes simples sur plusieurs bases

Première hypothèse :

- Les généticiens ont un modèle mental du domaine

Seconde hypothèse :

- Il est possible de trouver des correspondances entre leur modèle mental et le modèle réel

Test vocabulaire métier

- Valider première et seconde hypothèses
- Produire des données pour la phase suivante

Test prototype

- Construction d'un prototype à partir des résultats précédent
- Évaluation du prototype avec des utilisateurs

Détermination du vocabulaire métier

Objectifs

- Vérifier nos hypothèses
- Déterminer le vocabulaire utilisé, les structures grammaticales,...
- Vérifier la présence d'homonymie ou de synonymie

5 participants :

- un médecin généticien moléculaire
- un médecin clinicien
- une ingénieure en biologie moléculaire
- une médecin cytogénéticienne
- une biologiste moléculaire

Détermination du vocabulaire métier

3 bases de données :

- OMIM : gène, maladie
- HPO : maladie, phénotype, gène
- dbSNP : variation, gène

Aperçu des données :

```
#Format: diseaseId<tab>gene-symbol<tab>gene-id(entrez)<tab>HPO-ID<tab>HPO-term-name
OMIM:608980      FREM1    158326  HP:0000414    Bulbous nose
OMIM:608980      FREM1    158326  HP:0000077    Abnormality of the kidney
OMIM:608980      FREM1    158326  HP:0000322    Short philtrum
OMIM:608980      FREM1    158326  HP:0011803    Bifid nose
OMIM:608980      FREM1    158326  HP:0000143    Rectovaginal fistula
OMIM:608980      FREM1    158326  HP:0001545    Anteriorly placed anus
OMIM:608980      FREM1    158326  HP:0000007    Autosomal recessive inheritance
OMIM:611528      JUP      3728    HP:0004756    Ventricular tachycardia
OMIM:611528      JUP      3728    HP:0000006    Autosomal dominant inheritance
```

Résultats

- Les concepts sont assez facilement isolés
- Si les concepts sont bien isolés, les attributs sont bien répartis
- Le type des attributs n'est pas spontanément précisé
- Beaucoup de synonymes sont employés – jusqu'à 7 pour une personne
- Plusieurs homonymie, mais jamais chez un seul participant
- Les cardinalités peuvent être spontanément précisé via des modaux

Remarques

- Importance des exemples
- Attribut piège : identifiant OMIM
- Expertise réellement indispensable

Hypothèses

Conclusions du test

- Il existe un modèle mental pour les experts
- Il y a des correspondances avec le modèle réel

Prochaine étape

Peut on, à partir du modèle mental des experts, obtenir les informations permettant de bâtir l'ontologie ?

Mise en place d'une stratégie IHM

→ Développement d'un prototype pour validation

Développement du prototype

4 étapes :

- Phase de parsing des données : import de fichiers semi structurés, parsing et affichage
- Extraction des concept : à partir des données affichées, extraction des concepts et des attributs

Nom de la notion :

SNP

Attributs :

chromosome	Identifiant :	<input type="checkbox"/>
rsID	Identifiant :	<input checked="" type="checkbox"/>
position	Identifiant :	<input type="checkbox"/>

[Ajouter un champs](#)

valider

Passer à l'étape suivante

Développement du prototype

- Définition des relations : création des relations entre les concepts, avec leurs cardinalités

Concept 1	Type de relation	Nom de la relation	Nombre d'acteur 2	Acteur 2
<input type="text" value="chromosome"/>	<input type="text" value="peut"/>	<input type="text"/>	<input type="text" value="un seul"/>	<input type="text" value="chromosome"/>
genes	doit	presenter	un seul	chromosome
genes	peut	avoir	un ou plusieurs	transcrit

- Visionnage des résultats : récapitulatifs des concepts et relations saisis

Test du prototype

Objectif :

Évaluer la faisabilité de l'approche

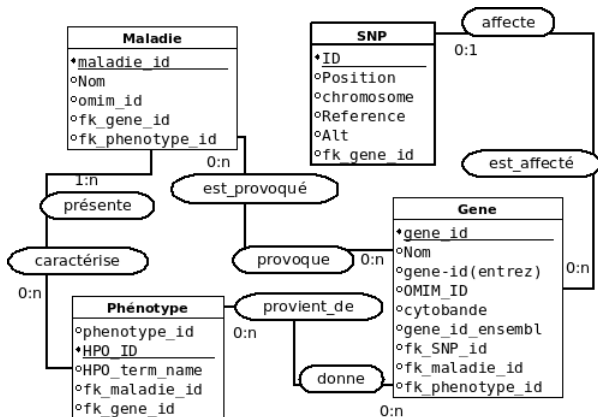
Scénario

Le participant doit rassembler trois bases de données biologiques

4 Participants :

- un médecin généticien moléculaire
- un médecin clinicien
- une ingénieure en biologie moléculaire
- nouveau participant : interne en génétique

Test du prototype



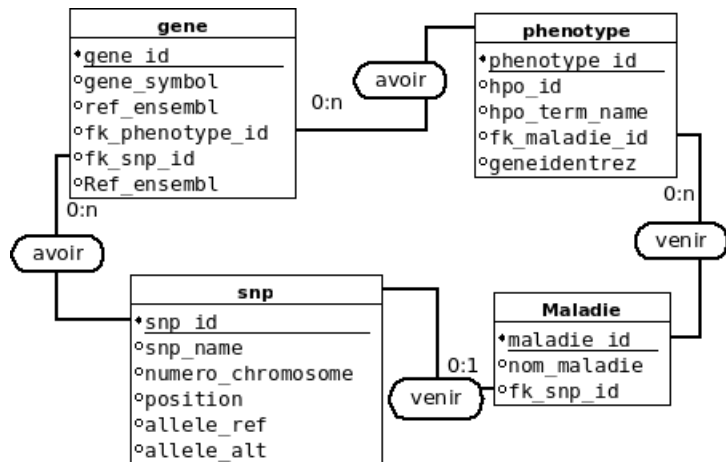
Évaluation

Comparaison des schémas entités associations obtenus avec l'objectif

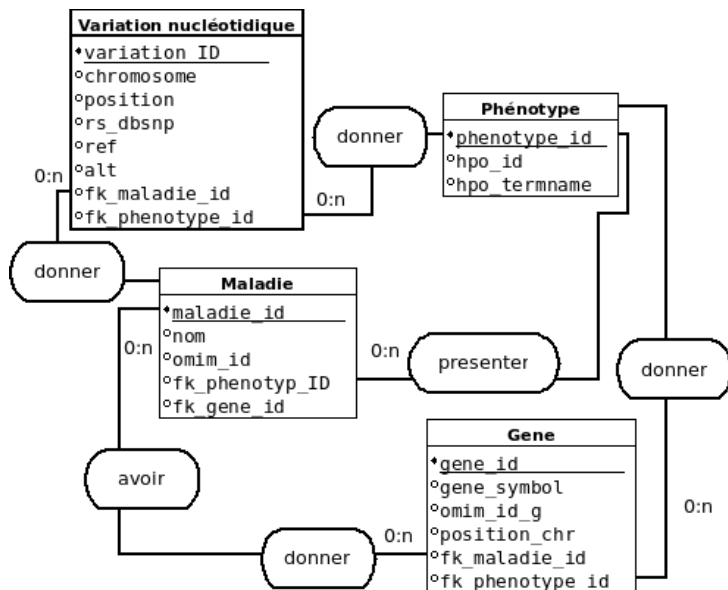
Observation des participants : points de blocages, réflexions, temps

Question aux participants

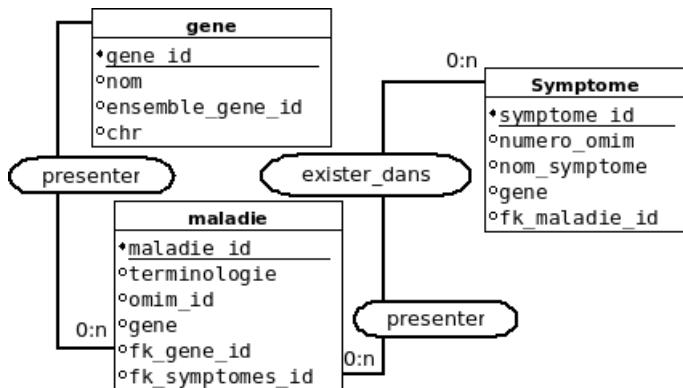
Résultats



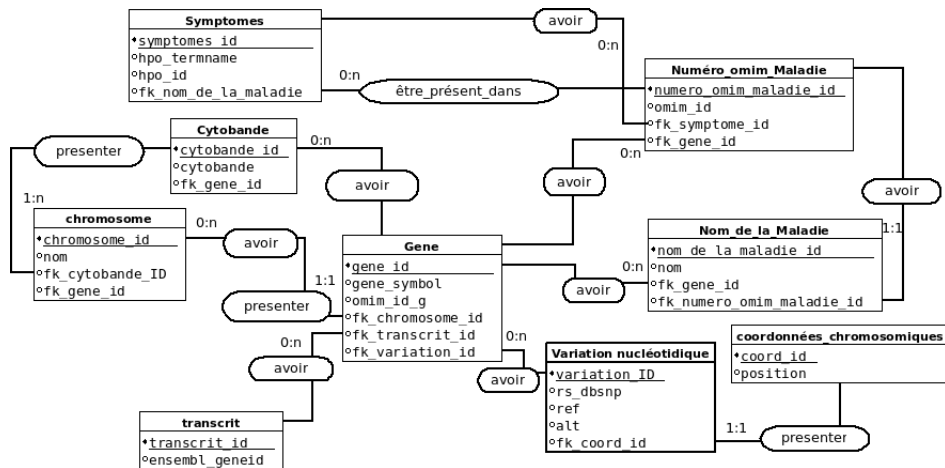
Résultats



Résultats



Résultats



Bilan résultats

Concepts

- 4 sur 4 pour 2 participants
- 1 manque
- une explosion

Relations

- des oublis...
- ...et des relations supplémentaires
- pas de concepts oubliés
- cardinalité parfois élargie

Bilan résultats

Remarques :

- Tous les systèmes sont implémentables
- Plusieurs synonymes pour les concepts et les attributs
- Cas du "faux" concept
- Réponses aux questions :

Question	F.	G.	S.	X.	Total sur 20
plausibilité du scénario	5	5	3	5	18
simplicité des tâches	3	3	4	3	13
clarté de l'interface	5	3	2	3	13

Situations particulières

- Relation 1 : 1
→ Détection et suggestion à l'utilisateur
- "Faux" concept
→ Impossible à détecter, mais résolution par l'utilisateur
- Oubli de concept
→ Vérifier l'attribution de chaque colonne, peut être volontaire
- Oubli de relation
→ Concept isolé, présence de plusieurs concepts dans un seul fichier
- Erreur de cardinalité
→ vérifiable pour certain cas (0 :1), pas pour d'autres (0 :n)

Conclusion générale

Conclusions

- Il existe un modèle mental du domaine chez les généticiens
- Ce modèle a des similitudes avec le modèle réel
- Il est possible de reconstruire une ontologie avec l'aide de l'expert

Prochaines étapes

- Implémentation des heuristiques d'aides
- Résolution des conflits sémantiques
- Système de requête adapté à l'utilisateur