



TP : Pipeline d'annotation des variants génétiques sous Galaxy



Ecole de Bioinformatique Roscoff 2013 – Atelier Annotation des SNPs Sabrina Rodriguez et Maria Bernard



Objectif:

Cette formation a pour but de vous apprendre à annoter fonctionnellement des SNPs, annotations qui sont des critères de choix pour ensuite approfondir vos analyses. Vous avez vu dans les ateliers précédents « Détections des SNPs » comment identifier des SNPs. Nous verrons ici comment sélectionner des SNPs sur différents critères, comment prédire leurs conséquences, et enfin comment les annoter fonctionnellement.

Indications Pratiques :

Pour vous connecter à Galaxy :

- instance de Roscoff : <u>http://galaxy.sb-roscoff.fr/</u>
- instance de Toulouse : <u>http://sigenae-workbench.toulouse.inra.fr/</u>

Identifiant de formation pour l'instance de Toulouse :

- <u>login :</u>
 - anemone
 - aster
 - bleuet
 - iris
 - muguet
 - narcisse
 - pensee
 - rose
 - tulipe
 - violette
 - lilas
 - pervenche
 - laurier
 - lavande
 - lis
 - capucine
 - coquelicot
 - geranium
 - liseron
 - arome
 - chardon

- Password

• f1o2r3!

Tous les outils concernant l'annotation sont classés dans la section « SNP annotation», dans le menu de gauche.

D'autres outils généraux vous seront très souvent utiles pour manipuler vos fichiers, n'hésitez pas à naviguer dans ces catégories.



ETAPE 1 : Initialisation de votre projet d'analyse

1) Les données

Pour la formation nous mettons à votre disposition un jeu de données test « Fichier_test_annot.vcf ». Le fichier est téléchargeable à cette adresse : <u>http://snp.toulouse.inra.fr/~sigenae/Galaxy_Formation/Annotation_SNP/Fichier_test_annot_.vcf</u>

Ces SNPs ont été détectés chez le bovin, c'est donc sur cet organisme que nous allons travailler.

Version du génome : Cow December 2009 (UMD 3.1 assembly)

2) Charger les données dans Galaxy : outil « Get Data »

Dans Galaxy pour accéder à vos fichiers vous avez 2 modes :

• Upload from your computer :

Cet outil vous permet de télécharger un fichier de votre ordinateur ou bien disponible via une URL sur internet (comme c'est le cas ici).

Une fois téléchargé vous pouvez renommer le nom du fichier grâce à l'outil « pencil » à côté de votre fichier. (exemple : « Fichier_test_annot.vcf »).

Attention : ce mode copie votre fichier sur le serveur Galaxy, il consomme donc une partie de votre quotat.

• Upload local file from filesystem path

Cet outil permet de communiquer le chemin de votre fichier sur Genotoul/Roscoff au serveur Galaxy sans pour autant le copier. Ce mode vous permet d'économiser votre quotat sur le serveur Galaxy.

Pour Genotoul (comme c'est le cas ici), le fichier est déjà déposé sur vos comptes de formation il vous suffit d'indiquer les informations suivantes :

File Name : Fichier_test_annot Filte Type: VCF Path to file: /work/USERNAME/galaxy/Fichier_test_annot.vcf

Maintenant que notre fichier est chargé, notre projet d'analyse peut commencer, renommons donc notre historique dans le menu de droite,exemple « **TP_Annotation** ». Un bilan des fichier enregistrés sur votre compte est disponible dans « Saved Datasets ».

	Analyze Data	Workflow	Shared Data	Visualization	Help	User Welcome mbernard
Saved Datasets	;					Logged in as mbernard@toulouse.inra.fr
search Advanced Search						Logout Saved Histories
□ <u>Name</u>			<u>History</u>	т	ags	Saved Datasets
Eichier test annot.	<u>vcf</u> 🔻		TP Annotation	<u>0</u>	<u>Tags</u>	Public Name
For 0 selected datas	ets: Copy to c	urrent histo	ry			



3) Explorer et sélectionner des SNPs

Vous pouvez visualiser un fichier en utilisant l'outil « eye » de ce fichier.

Exercice :

Supprimez les lignes d'en-tête qui vous gêneront lors de l'import du fichier dans Excel et téléchargez le fichier VCF sur votre ordinateur.

Créez un fichier qui ne contient pas le chromosome 8.



 Pour supprimer l'entête il suffit de supprimer les 26 premières lignes de notre fichier VCF grâce à l'outil « *Remove beginning* » de la catégorie « *Text Manipulation* »



Renommez le nouveau fichier créé puis utilisez l'outil « *download* » sur le nouveau fichier généré.

	History	/			Optio	ns 🔻
						12 📑
	TP_An	notation			52	.1 Kb
	- 1				-	0.00
	5:				۲	0 23
	Fichie	r_test_a	not	t_no	_head	er.vc
	<u>1</u>					
	71 line	S				
	format	: vcf, dat	abas	e: <u>?</u>		
I	Info	pilog : jo	b fir	nishe	d at m	er.
	Downloa	5:15:5	7 CE	T 20	13	
		0 🖄				47 📄
	1 Chro	m 2.Pos	3.10) 4.R	ef 5.A	lt 6.Q
	Chr1	18073037		С	Α	20.2
	Chr1	42531330	•	С	т	22
	Chr1	53156188		С	Т	65.1
	Chr1	58593687	•	G	Α	4.13
	Chr1	58792664	•	G	Α	110
	Chr1	68929240	•	C	Α	22.1
	((<u></u>)))))
	1: Fich	nier_test	anı	not.	<u>/cf</u> @	0 %

• Pour sélectionner uniquement certaines lignes, utilisez l'outil « *Filter* » de la catégorie « *Filter and Sort* ».

Cet outil permet de sélectionner des lignes en fonction des valeurs qu'elles contiennent dans 1 ou plusieurs colonnes. Votre fichier doit nécessairement être un fichier de type tabulé.

Dans le champ « With following condition: » écrire c1!= « Chr8».



Note :

- Si la valeur recherchée est une chaine de caractères, on la met entre simple côte
- S'il s'agit d'un nombre, il ne faut pas de côtes

N'oubliez pas de renommer votre nouveau fichier, Fichier_test_annot_except_chr8.vcf

Pour se faire, vous pouvez cliquez sur le « crayon » de la boite verte en haut à droite contenant le résultat de votre sélection de chromosomes, vous permettant d'éditer les attributs :

	Edit Attributes	-	History	Options 👻	
<	Name: Fichier_test_annot_except_chr8.vcf Fichier_test_annot_except_chr8.vcf Filtering with cl1=chr8. Filtering with cl1=chr8. Annotation / Notes: None		 ieudi_100113_pa es B: Fichier test a header_noCHR8 7: Fichier test a header.vcf 	r_modul 57.3 Kb cfit attributes nnet_no_@ 0 %	>
4	Add an annotation or notes to a dataset; annotations are available when a history is viewed. Database/Build: [Click to Search or Select Number of comment lines:		1: Fichier test a	nnot.vcf 👁 0 🔀	
	Sore column for visualization:				



1) <u>Génération des conséquences</u>

<u>Rappel</u>

La première annotation que l'on peut générer est ce que l'on appelle la conséquence de chaque SNP. Cette annotation se fait grâce à l'outil : « *Variant_effect_Predictor* ».

Ce programme va parcourir un fichier VCF, et se connecter à Ensembl selon l'espèce du génome qu'on lui précisera. Il va ensuite produire des annotations de bases ou conséquences de chaque SNP en fonction de la localisation du SNP sur le génome dans Ensembl.

Lancement

Γ.	
	* Variant effect predictor (version 1.0.0)
	Your input file: 10: Fichier_test_annopt_chr8.vcf 🔍
	File type: vcf 💌
	Specie name:
	Single project name - ATTENTION : Please, repeat this same name for each Galaxy module used: mbernard_TP_Annotation
	Execute

Attention de bien préciser un nom de projet. Celui-ci permettra à Galaxy de retrouver tous les fichiers de résultats lorsque vous voudrez les fusionner. Exemple : « USERNAME_TP_Annotation ».

Vous devez garder ce nom de projet tout au long de l'analyse de votre échantillon (ou fichier .vcf), mais, ce nom ne pourra être utilisé que pour 1 analyse (pour une 2eme analyse, il faudra changer de nom de projet).

Ne pas oublier de préciser l'espèce : « cow ».



<u>Résultat</u>

Variant_effect_predictor retourne différentes annotations :



Others: Within non-coding gene, Within mature miRNA, NMD transcript

Il retourne également des informations concernant le gène impacté, et l'élément du gène impacté ainsi que des informations sur la position du SNP dans les différents éléments du gène et la conséquence au niveau de la séquence protéique.



Renommez le fichier, par exemple : Fichier_test_annot_except_chr8.vep

Ecole de Bioinformatique Roscoff 2013 – Atelier Annotation des SNPs Sabrina Rodriguez et Maria Bernard



Exercice :

Sélectionnez les SNPs qui sont catégorisés comme « NON_SYNONYMOUS_CODING » et « INTRONIC ».

Sachant que nous avons 1 SNP par ligne, comment savoir le nombre de SNPs sélectionnés ?



• Pour sélectionner les SNPs, vous pouvez utiliser « Select ».

L'outil « Select » est plus généraliste que l'outil « Filter » qui ne fonctionne que sur des fichiers tabulés, mais il ne vous permet pas de sélectionner la colonne sur laquelle vous voulez faire votre recherche.

Dans le champ « the pattern: » écrire : NON_SYNONYMOUS_CODING|INTRONIC

Attention ne pas mettre d'espace, si tel était le cas l'outil chercherait également ce caractère, or nous avons une tabulation entre chaque colonne !!!

Renommez votre fichier de sortie : exemple : « Fichier_test_annot_except_chr8_NCS_INTRONIC.vep »

- Pour connaître le nombre de lignes d'un fichier, donc le nombre de SNPs,
 - o Vous pouvez lancer l'outil « Line/Word/Character count » en cochant « line »,
 - Ou plus simplement cliquer sur le fichier d'intérêt et vous avez en dessous une brève description du fichier qui contient parfois le nombre de lignes.



2) Formatage et génération des séquences protéiques pour la recherche des conséquences des SNPs

<u>Rappel</u>

Le formatage des données consistes à récupérer les séquences protéiques de chaque couple [allèle référence / allèle alternatif] de chaque SNP.

Dans le cas ou le SNP est multiallèlique (> 2 allèles), il peut y avoir plus de 2 couples d'allèles par SNP.

Si un SNP affecte plus de 2 transcrits, il y aura aussi plus de 2 couples, c'est à dire au moins 4 protéines générées (référence / alternatif 1 et référence / alternatif 2).

Les noms initiaux des séquences sont formés à partir de leur localisation sur le génome, du nom du transcrit auquel il appartient et des allèles référence et alternatif qu'il contient suivi de « alléle1 » pour la référence et de « allèle 2 » pour l'alternatif. Les séquences sont ensuite renommées avec un encodage de noms pour faciliter l'utilisation de certains outils d'annotation en utilisant la syntaxe « _seq1, seq_2... ».

L'outil qui permet de faire ce formatage est l'outil « Format SNP Effect - Protseq ».

Lancemer	٦t
	-

** Format SNP effect - Protseq (version 1.0.0)	
Your input file Varient Effect Predictor (format .vep 13: Fichier_test_annoTRONIC.vep)):
Specie name: Cow	
value between 0 and 1. Discriminant for the score 0.5	comparison. (0.5 per default):
Your unique project name:	
mbernard_TP_Annotation	

Attention utilisez le même nom de projet que pour « *Variant effect Predictor* » (cf : « USERNAME_TP_Annotation »).

Astuce double cliquez dans la zone de texte pour faire apparaître les noms de projets tapés précédemment ou cliquer sur votre fichier VEP initial, cliquer sur l'outil informations, vous trouverez les paramètres que vous avez utilisés pour générer ce fichier.





<u>Résultats</u>

[★] 20: ** Format SNP
[*] <u>19: ** Format SNP</u> ● Ø × effect - Protseq on data 15
18: ** Format SNP

4 boites de dialogues sont lancées en même temps pour ce module.

En effet, Protseq va générer différents fichiers que vous allez devoir ensuite utiliser avec les modules de prédiction des conséquences.

Renommez vos fichiers en fonction de ce qu'ils contiennent :

Le premier fichier généré est un fichier html bilan des différentes sorties du programme. Renommez le « Protseq_bilan.html »

Voici l'ensemble des fichiers generes par le script PROTSEQ : Veuillez cliquer sur les liens pour ouvrir vos fichiers resultats.

Link to specie form file. This file will contain the protein sequences FASTA FORMAT, 1 per allele from the input file

Link to specie.fasta file. This file with protein sequences, 1 per allele, with short encoded/shorter name

Le second fichier est un fichier de séquences protéiques au format fasta. Renommez le « sequences_prot.fasta »

> seq1
$M \overline{w} plvvvlllgsvrcgsaqlifnaiksveytlcnqtvvipcfvnnvetknitelyvrwkfkgenififdgsqrmskpssnfssaeiapsellrgiaslkmaksdavlgnytcevtelsregetiielkyrvvswfspnenil$
>_seq2
$\tt M wplvvvlllgsvrcgsaqlifnaiksveytlcnqtvvipcfvnnvetknitelyvrwkfkgenififdgnqrmskpssnfssaeiapsellrgiaslkmaksdavlgnytcevtelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelkyrvvswfspnenilitelkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelk$
>_seq3
$\tt Mwpluvullgsvrcgsaqlifnaiksveytlcnqtvupcfvnnvetknitelyvrwkfkgenififdgsqrmskpssnfssaelapsellrgiaslkmaksdavlgnytcevtelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelsregetiielkyrvvswfspnenilitelkyrvvswfspnenilitelkyrvvswfspnenilitelkyrvvswfspnenilitelkyrvvswfspnenilitelkyrvvswfspnenilitelkyrvvswfspnenilitelkyrvvswfspnenilitelkyrvvswfspnenilitelkyrvvs$
>_seq4
$\tt Mwplvvvlllgsvrcgsaqlifnaiksveytlcnqtvvipcfvnnvetknitelvvrwkfkgenififdgnqrmskpssnfssaeiapsellrgiaslkmaksdavlgnytcevtelsregetiielkyrvvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiielkyrvswfspnenilitelserecturgiaslkmaksdavlgnytcevtelsregetiitelserecturgiaslkmaksdavlgnytcevtelsregetiitelserecturgiaslkmaksdavlgnytcevtelsregetiitelserecturgiaslkmaksdavlgnytcevtelsregetiitelserecturgiaslkmaksdavlgnytcevtelsregetiitelserecturgiaslkmaksdavlgnytcevtelsregetiitelserecturgiaslkmaksdavlgnytcevtelsregetiitelserecturgiaslkmaksdavlgnytcevtelsregetiitelserecturgiaslkmaksdavlgnytcevtelsregetiitelserecturgiaslkmaksdavlgnytcevtelsregetiitelserecturgiaslkmaksdavlgnytcevtelsregetiitelserecturgiaslkmaksdavlgnytcevtelsregetiitelserecturgiaslkmaksdavlgnytcevtelsregetiitelserecturgiaslkmaksdavlgnytcevtelsregetiitelserecturgiaslkmaksdavlgnytcevtelserecturgiaslkmaksdavlgnytcevtelserecturgiaslkmaksd$
>_seq5
$\tt MSFVRVNRYGPRGGGRKTLKVKKKTSVKQEUDNTVTDLTVHRATPEDLIRRHEIHKSKNRALVHUELQEKALKRRUKKQKPEISNLEKRRLSIMKEILSDQYQLQDVLEKSDHLMATAKGLFVDFPRRRTGFPNVTMAPESSTATISTICKTSVKQEUDNTVTDLTVHRATPEDLIRRHEIHKSKNRALVHUELQEKALKRRUKKQKPEISNLEKRRLSIMKEILSDQYQLQDVLEKSDHLMATAKGLFVDFPRRRTGFPNVTMAPESSTATISTICKTSVKQEUDNTVTDLTVHRATPEDLIRRHEIHKSKNRALVHUELQEKALKRRUKKQKPEISNLEKRRLSIMKEILSDQYQLQDVLEKSDHLMATAKGLFVDFPRRRTGFPNVTMAPESSTATISTICKTSVKQEUDNTVTDLTVHRATPEDLIRRHEIHKSKNRALVHUELQEKALKRRUKKQKPEISNLEKRRLSIMKEILSDQYQLQDVLEKSDHLMATAKGLFVDFPRRRTGFPNVTMAPESSTATISTICKTGFPNVTMAPESSTATISTICKTGFPNVTMAPESSTATISTICKTGFPNVTMAPESSTATISTICKTGFPNVTMAPESSTATISTICKTGFPNVTMAPESSTATISTICKTGFPNVTMAPESSTATISTICKTGFPNVTMAPESSTATISTICKTGFPNVTMAPESSTATISTICKTGFPNVTMAPESSTATISTICKTGFPNVTMAPESSTATISTICKTGFPNVTMAPESSTATISTICKTGFPNVTMAPESSTATISTICKTGFPNVTMAPESSTATISTICKTGFPNVTMAPESSTATISTICKTGFPNVTMAPESSTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTICKTGFPNVTMAPESTATISTITATISTICKTGFPNVTMAPESTATISTITATISTICKTGFPNVTMAPESTATISTICKTGFPNT$

 Le troisième fichier permet la correspondance entre le nom des séquences protéiques et les différents allèles des SNPs « NON_SYNONYMOUS_CODING ».
 Renommez le « seq_prot_snp.name. »

>1 52156199 C/T ENGETATOOOOOO0466419/N 711 ollolo1	aca1	
>1_3136166_C/1M3B1X10000004664[5/N-71]_a11e1e1	_seqr	Sub
>1_53156188_C/TENSBTAT00000004664[S/N-71]_allele2	_seq2	snp
>1_53156188_C/TENSBTAT00000055450[S/N-71]_allele1	_seq3	snp
>1_53156188_C/TENSBTAT00000055450[S/N-71]_allele2	_seq4	snp
>1_58593687_G/AENSBTAT00000057527[S/F-633]_allele1	_seq5	snp

 Le quatrième fichier est un fichier VEP (fichier de sortie de Variant Effect Predictor) qui ne contient que les SNPs « NON_SYNONYMOUS_CODING ». Renommons le « Fichier_test_annot_except_chr8_NCS.vep »

	1_53156188_C/T	1:53156188	Т	ENSBTAG0000003585	ENSBTAT0000004664	Transcript	NON_SYNONYMOUS_CODING	330	212
	1_53156188_C/T	1:53156188	Т	ENSBTAG0000003585	ENSBTAT00000055450	Transcript	NON_SYNONYMOUS_CODING	330	212
	1_58593687_G/A	1:58593687	A	ENSBTAG0000013937	ENSBTAT00000057527	Transcript	NON SYNONYMOUS CODING	1977	1898
	1_58792664_G/A	1:58792664	A	ENSBTAG0000021000	ENSBTAT0000027966	Transcript	NON_SYNONYMOUS_CODING	4856	4856
	1_68929240_C/A	1:68929240	A	ENSBTAG0000000566	ENSBTAT0000000740	Transcript	NON SYNONYMOUS CODING	239	239
	1_68929240_C/A	1:68929240	A	ENSBTAG0000000566	ENSBTAT00000054897	Transcript	NON_SYNONYMOUS_CODING	198	95
	1 95205755 G/A	1:95205755	A	ENSBTAG0000037965	ENSBTAT00000052461	Transcript	NON SYNONYMOUS CODING	512	512
	1_145178567_G/A	1:145178567	A	ENSBTAG00000017010	ENSBTAT0000022620	Transcript	NON_SYNONYMOUS_CODING	710	643
l	1 154285564 C/A	1:154285564	A	ENSBTAG0000040193	ENSBTAT0000030369	Transcript	NON SYNONYMOUS CODING	158	86



Bilan des dataSets enregistrés :

🗌 Fichier test annot except chr8 NCS.vep 🔻	TP Annotation	<u>O Taqs</u>	2 days ago
🗖 seq prot snp.fasta 🔻	TP Annotation	<u>O Taqs</u>	2 days ago
🗇 <u>sequences prot.fasta</u> 👻	TP Annotation	<u>O Taqs</u>	2 days ago
Protseq_bilan.html +	TP Annotation	<u>O Taqs</u>	2 days ago
Fichier test annot except chr8 NCS INTRONIC.vep 💌	TP Annotation	<u>O Taqs</u>	2 days ago
🗖 Fichier test annot except chr8.vep 🔻	TP Annotation	<u>O Taqs</u>	2 days ago
🗆 Fichier test annot except chr8.vcf 🔻	TP Annotation	<u>O Taqs</u>	2 days ago
🗖 Fichier test annot no header.vcf 🔻	TP Annotation	<u>O Taqs</u>	2 days ago
□ Fichier test annot.vcf 👻	TP Annotation	<u>O Taqs</u>	2 days ago
For 0 selected datasets: Copy to current history			

Exercice

Combien y a-t-il de séquences dans le fichier fasta ?





• Pour connaître le nombre de séquence d'un fichier fasta

Il faut simplement cliquer sur le nom de ce fichier, et dans la brève description nous n'avons pas le nombre de lignes mais le nombre de séquences du fichier.



ETAPE 3 : Annotation fonctionnelle des SNPs « non-synonymous coding ».

Nous allons maintenant lancer les différents modules de prédictions des conséquences des SNPs « NON_SYNONYMOUS_CODING ».

Les fichiers de sortie indiqueront au niveau de la colonne « case » : « loss », « gain » ou vide. Les colonnes suivantes indiqueront les scores pour l'allèle de référence et l'allèle alternatif (sauf pour le module conservation). La perte « loss » d'une modification post-traductionnelle correspond à l'allèle 1 avec signal prédit et allèle 2 sans signal prédit pour le module d'annotation fonctionnelle concerné.

Le gain « gain » d'une modification post-traductionnelle correspond à l'allèle 2 avec signal et allèle 1 sans signal. A cette comparaison s'ajoute une comparaison des scores. Dans le cas ou le programme génère un score de prédiction mais, sans fournir de signal fort; si les 2 allèles ont fourni un tel score et que la différence de score est supérieure ou égale au seuil minimum fixée en paramètre (0.5), alors, une sortie avec « perte? » ou « gain? » est indiquée ainsi que les scores des allèles 1 et 2 ayant permit de générer l'hypothèse de prédiction.

Pour chaque module, sont générés :

- un fichier contenant les résultats « bruts » du programme
- un fichier tabulaire contenant les résultats interprétés que l'on nommera
- « [nom_de_module]_results.txt »
- un fichier contenant les titres des colonnes des résultats que l'on nommera « [nom de module] title.txt ».

Les fichiers « _results.txt » contiennent la liste des SNPs avec les gains ou pertes de signaux. Le fichier « title » contient les noms des colonnes des informations trouvées dans le fichier « results».

1) <u>Mitoprot</u>

Pour étudier les signaux d'adressage à la mitochondrie, sélectionnez l'outil « *run_Mitoprot »*.

** Run Mitoprot (version 1.0.0)
Protein sequence query file - !!!the sequence name should be the short name!!! - Fasta file: 15: sequences_prot.fasta 💌
file containing the encoded protein names (long name 'allele1' - short name '_seq') - txt format:
16: seq_prot_snp.name 🔍 🔻
value between 0 and 1. Discriminant for the score comparison. (0.5 per default):
Your single project name:
mbernard_TP_Annotation
Execute



Les fichiers de sorties mitoprot

Renommez les fichiers de sorties en fonction des exemples ci dessous

- Extrait du fichier contenant les titres des colonnes de résultats « mitoprot_title.txt »
 - Sequence Name Case score allele1 score allele2
- Extrait du fichier contenant les résultats formatés « mitoprot_results.txt » :

1_68929240_C/AENSBTAT00000054897[W/L-32]	loss	0.3961	0.4230
1_154285564_C/AENSBTAT00000030369[S/I-29]	loss	0.3758	0.4093
1_42531330_C/TENSBTAT00000005398[E/K-29]	loss	0.0900	0.6499
2_126676995_C/TENSBTAT00000055788[P/S-17]	gain	0.3499	0.3855
3_25969298_G/TENSBTAT00000039902[P/T-21]	loss	0.2367	0.3668
4_94912500_C/AENSBTAT00000017924[R/S-11]	loss	0.7161	0.3288
4_106869716_G/AENSBTAT00000064277[E/K-4]	gain	0.1166	0.4817
4_116237864_G/CENSBTAT00000011851[W/S-9]	loss	0.6446	0.4435
4_12724969_A/GENSBTAT00000065396[S/G-27]	loss	0.6108	0.5940
13_61584514_T/AENSBTAT00000027390[L/Q-11]	gain	0.1882	0.2982

• Extrait du fichier contenant les résultats « mitoprot_bruts.txt » :

>1_53156188_C/TENSETAT00000004664[S/N-71]_allele1 0.2119	not imported
>1_53156188_C/TENSBT&T00000004664[S/N-71]_allele2 0.2119	not imported
>1_53156188_C/TENSBT&T00000055450[S/N-71]_allele1 0.2119	not imported
>1_53156188_C/TENSBTAT00000055450[S/N-71]_allele2 0.2119	not imported
>1_58593687_G/AENSBTAT00000057527[S/F-633]_allele1 0.9538	imported
>1_58593687_G/AENSBTAT00000057527[S/F-633]_allele2 0.9538	imported

2) <u>Netphos</u>

Pour prédire s'il y a acquisition ou perte de sites de phosphorylation, sélectionnez l'outil « *run_Netphos ».*

** Run netphos (version 1.0.0)
Protein sequence query file - !!!the sequence name should be the short name!!! - Fasta file: 15: sequences_prot.fasta 💌
file containing the encoded protein names (long name 'allele1' - short name '_seq') - txt format: 16: seq_prot_snp.name
value between 0 and 1. Discriminant for the score comparison. (0.5 per default): 0.5
Your unique project name: mbernard_TP_Annotation
Execute



Les fichiers de sorties netphos

Renommez les fichiers de sorties en fonction des exemples ci-dessous.

• Extrait du fichier contenant les titres des colonnes de résultats « netphos_title.txt »

Name	ATM result	ATM all	ele1	score	ATM	allele2	score	CKI result
• E	xtrait du fichier c	ontenant	les ré	sultats	formatés	« netph	os_results.t>	kt » :
1_5315618	88_C/TENSBTAT000000554	450[S/N-71]						
1_5859368	37_G/AENSBTAT000000575	527[S/F-633]	loss	0.553				
1_5879266	54_G/AENSETATOOOOOO279 10 C/A - ENSETATOOOOOO007	966[S/F-1619] 740[W/L-80]	loss	0.529				
1_6892924	O_C/AENSBTAT000000548	897[W/L-32]						
1_9520575	55_G/AENSBTAT000000524	ł61[S/N−171]	loss	0.546				
2_1266769	95_C/TENSBTAT00000055	5788[P/S-17]						
2_1333210)47_C/TENSBTAT00000017	7333[S/F-83]	loss	0.548				
2_1791666	50_G/AENSBTAT000000401	L94[G/S-342]	gain		0.546		gain	0.536
2_8525906	59_G/AENSBTAT000000383	807[P/S-441]						
3_1004985	539_A/CENSBTAT00000017	7715[L/R-240]						

• Extrait du fichier contenant les résultats « netphos_bruts.txt » :

This a <u>Show</u>	dataset is <u>/ all</u> <u>Sav</u>	s large e	and only the first megabyte is shown below.
 12 12 12 12 12 12 12 12 12 12 12 12 12 1	PKC cdc2 CaM-II GSK3 CKI P38MAPK DNAPK DNAPK ATM RSK CKII PKG PKA cdb5	0.712 0.514 0.455 0.350 0.340 0.343 0.301 0.245 0.228 0.228 0.228 0.228 0.228 0.228	YE 3 YE 3

3) <u>Gpi</u>

Pour analyser les changements d'ancrage GPI, selectionnez l'outil « Run_gpi ».

** Run gpi (version 1.0.0)
Protein sequence query file - !!!the sequence name should be the short name!!! - Fasta file: 15: sequences_prot.fasta
file containing the encoded protein names (long name 'allele1' - short name '_seq') - txt format:
16: seq_prot_snp.name 🔍
value between 0 and 1. Discriminant for the score comparison. (0.5 per default): 0.5
Your single project name:
mbernard_TP_Annotation
Execute

GPI produit 2 fichiers de sorties, un fichier title et un fichier « résults » comme précédemment.



4) <u>Conservation</u>

Pour analyser la conservation protéique entre référence et alternatif, utilisez l'outil « extract_conservation_results.pl »

** Extract conservation results (version 1.0.0)
Your vep query file:
17: Fichier_test_annohr8_NCS.vep 🔍
Specie name: Cow
Your single project name:
mbernard_TP_Annotation
Execute

Attention cette fois ci nous utilisons le fichier VEP produit par l'outil de formatage Protseq. Ce VEP ne contient que les NON_SYNONYMOUS_CODING SNPs : « Fichier_test_annot_except_chr8_NCS.vep ».

Conservation produit 2 fichiers de sorties, un fichier « title » et un fichier « résults » comme précédemment, renommez les.

Pour cet outil, les résultats sont constitués d'un tableau à 4 colonnes:

- 1 colonne pour le résultat généré à partir de la matrice Grantham de conservation des propriétés physico chimique au sein de la protéine entre référence et alternatif :
 - o conservées (0-50)
 - o modérément conservées (51-100),
 - o modérément radicales (101-150)
 - o radicales (≥ 151)
- et 3 colonnes concernant les résultats générés à partir de la base Ensembl (score attendu, score observé et différence de score) qui indique un taux de conservation de la mutation par rapport aux protéines orthologues.
 - Plus le score de différence est élevé, plus la position est conservée et donc plus une mutation à cette position a de fortes chances d'être délétère.

Exercice

Fusionnez tous les fichiers « title » avec leur correspondant « results ». Renommez les par exemple : **netphos_final, mitoprot_final, conservation_final,** et **gpi_final.**



• La concaténation de plusieurs fichier se fait grâce à l'outil « Concatenate datasets »

Selectionnez un fichier « title » suivi du fichier « result » correspondant, exécutez l'outil et renommez le résultat.

Concatenate datasets (version 1.0.0)
Concatenate Dataset: 19: netphos_title.txt
Datasets
Dataset 1
Select:
18: netphos_results.txt
Remove Dataset 1
Add new Dataset
Execute



Etape 4 : fusion et export des résultats.

<u>Join</u>

Le choix du fichier d'entrée sert uniquement à retrouver le chemin du répertoire de projet « USERNAME_TP_Annotation », vous pouvez par exemple utiliser le fichier fasta des séquences protéiques « sequences_prot.fasta ».

** Join annotation results (version 1.0.0)
Protein sequence query file - !!!the sequence name should be the short name!!! - Fasta file: 15: sequences_prot.fasta 💌
Your unique project name:
mbernard_TP_Annotation
Execute

<u>Remarque</u>

Cette fois nous avons 48 lignes dans notre fichier de sortie. Tous les SNPs contenus dans notre fichier VEP sont retournés, mais seulement les SNPs « NON_SYNONYMOUS_CODING » ont des résultats de score.

Le fichier de résultats est un fichier de type texte. Pour pouvoir faire des sélections sur les colonnes il est nécessaire de changer le type du fichier en « tabular ». Utilisez l'outil « pencil » et sélectionnez « tabular » dans « New Type ».

Exercice:

Combien de SNPs provoquent une perte de signalisation à la mitochondrie ? Combien de SNPs ont un score de conservation de Grantham > 100 ?



• Combien de SNPs provoquent une perte de signalisation à la mitochondrie ?

La sélection de la colonne correspondant aux résultats de perte de signalisation à la mitochondrie doit se faire grâce à une recherche du mot « loss » sur la 8^e colonne, l'outil « Filter » est dédié à ce type de recherche.

	<u>35: mito_signal_loss</u> ● Ø ※
	i viines Normat: tabular, database: <u>?</u>
Filter (version 1.1.0)	Info: Filtering with c8=='loss',
Filter: 34: Joint_annotation_results Dataset missing? See TIP below.	total lines). Epilog : job finished at ven. janv. 11 17:38:16 CET 2013
With following condition:	→
c8=='loss' Double equal signs, ==, must be used as shown above. T Execute	1 1_154285564_C/AENSETAT000000303 1_42531330_C/TENSETAT0000000539 1_68929240_C/AENSETAT0000005489 3_2596298_G/TENSETAT0000003990 4_116237864_G/CENSETAT000000118 4_12724969_A/GENSETAT0000006539

• Combien de SNPs ont un score de conservation de Grantham > 100 ?

Cette fois ci il s'agit de faire une comparaison numérique sur la 2^e colonne. L'outil « Filter » est toujours l'outil à utiliser.

	36: • • • ×
	grantham sup 100
	17 lines
	Info: Filtering with $c2 > 100$,
Filter (version 1.1.0)	kept 35.42% of 48 valid lines (48
Filter (Version 1.1.0)	total lines). Skinned 14 invalid line(s)
Filter	starting at line #1: "Name
	conservation Grantham_score
134: Joint_annotation_results	conservation Observed_score
Dataset missing? See TIP below.	conservation Expected_score
With following condition:	conservation Difference_score
	gpi case gp
c2>=100	
Double equal signs, ==, must be used as shown above. To filter for	
	12 61584514 T/b - FNSETAT000000272
Execute	1 154285554 C/A - EN3BTAT000000303
	1 154285554 C/A - EN3BTAT000000513
	1 18073037 C/A - ENSBTAT0000000078
	1_58593587_6/AEN3BTAT0000005752
	1_58792664_6/AEN3BTAT0000002796

Attention de ne pas mettre de côte lorsqu'il s'agit d'une comparaison numérique.



ETAPE 5 : Automatisation de l'analyse

L'utilisation de workflow dans Galaxy (démo)

Un workflow est déjà disponible pour faire toutes ces analyses.

Dans l'onglet workflow, en haut à droit cliquez sur « upload or import workflow », et entrez l'URL suivante :

http://snp.toulouse.inra.fr/~sigenae/Galaxy_Formation/Annotation_SNP/Galaxy-Workflow-Pipeline_annotation.ga

Le pipeline automatisé apparaît dans votre onglet workflow.

Cliquez sur Edit (flèche à côté du nom du workflow), pour voir les différents éléments qui le composent et l'enchainement des traitements.

output Q-Q	Your input file	** Format SNP effect - Protseq 🕱
	output_vep (txt)	Your input file Varient Effect Predictor (format .vep) Filter sequences by length &
		output_html (html) output_fasta (fasta) output_txt (txt) output_txt (txt)
* Run netphos 🛛 🕱		
<pre>>quence name should be the short ame!!! - Fasta file e containing the encoded protein ames (long name 'allele1' -</pre>	** Run gpi 🛛 😵	** Run Mitoprot
ort name '_seq') - txt format utput_results (txt) utput_title (txt) of atphos_output (txt)	Protein sequence query file - IIIthe sequence name should be the short nameIII - Fasta file file containing the encoded protein names (long name'allele1' - short name'.equ) - byt format	Protein sequence query file - IIIthe sequence name should be the short nameIII - Fasta file file containing the encoded protein n ames [long nameallele1' -
	output title (txt)	short name _seq) - txt format

Tous les programmes d'analyses que nous avons lancé précédemment se retrouvent dans le workflow, sauf le programme Join qu'il faut ensuite lancer manuellement à la fin..

Pour lancer le workflow, créer un nouvel historique et réimporter notre fichier d'entrée VCF. Sur l'onglet workflow, sélectionner « run » du pipeline d'annotation. Sélectionné votre VCF ainsi qu'un nouveau nom de projet, enfin cliquez sur « Run workflow ».



ETAPE 6 : Conclusion

• Tous les résultats intermédiaires sont sur le site : http://snp.toulouse.inra.fr/~sigenae/Galaxy_Formation/Annotation_SNP/

- Source vers les sites des modules
 - o Mitoprot

http://ihg2.helmholtz-muenchen.de/ihg/mitoprot.html MITOPROT: Prediction of mitochondrial targeting sequence

o Netphos

http://www.cbs.dtu.dk/services/NetPhos/

Sequence and structure based prediction of eukaryotic protein phosphorylation sites.

o GPI

http://gpi.unibe.ch/

GPI-SOM: Identification of GPI-anchor signals by a Kohonen Self Organizing Map

- o Conservation
 - Matrice Grantham

Amino acid difference formula to help explain protein evolution.

API Ensembl Compara

http://www.ensembl.org/info/docs/api/compara/index.html

• Installation galaxy et workflow SNP Annotation

Si votre laboratoire installe une instance de Galaxy et que vous désirez installer ces outils et ce workflow vous trouverez toutes les instructions en tapant cet commande : hg clone <u>http://inra-sigenae-sarah-maman@toolshed.g2.bx.psu.edu/repos/inra-sigenae-sarah-maman/snp_annotation</u>

