

Génotypage par séquençage de SNP de Pois

développement d'une approche par
séquençage HiSeq2000 d'ADN génomique



Groupe Légumineuses UMR IGEPP

Résistance Partielle Quantitative (RPQ)

1. Combinaisons de facteurs génétiques pour une RPQ efficace/durable

1.1- Sources

1.2- Diversité, efficacité, stabilité

1.3- Structure moléculaire QTL/gènes

1.4- Durabilité

2. Fonctions associées à la RPQ

2.1- Composantes biologiques

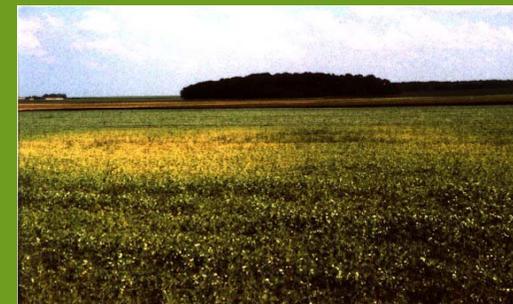
2.2- Gènes candidats régulés

2.3- Voies métaboliques régulées

Ascochytose (*Mycosphaerella pinodes*)



Aphanomyces (*Aphanomyces euteiches*)



Projet INRA-BIOGEMMA-SOFIPROTEOL PEAPOL 2011-2014

Demande de Sofiprotéol: développer rapidement un grand nombre de marqueurs SNP chez le pois, utiles à la sélection et accessibles à la communauté scientifique.

Co-coordonateurs : G. Boutet (INRA IGEPP) et N. Rivière (Biogemma)

4 Partenaires :

P1 : Biogemma

P2 : INRA UMR IGEPP Rennes

P3 : INRA Genotoul

P4 : INRIA-GenScale-Genouest



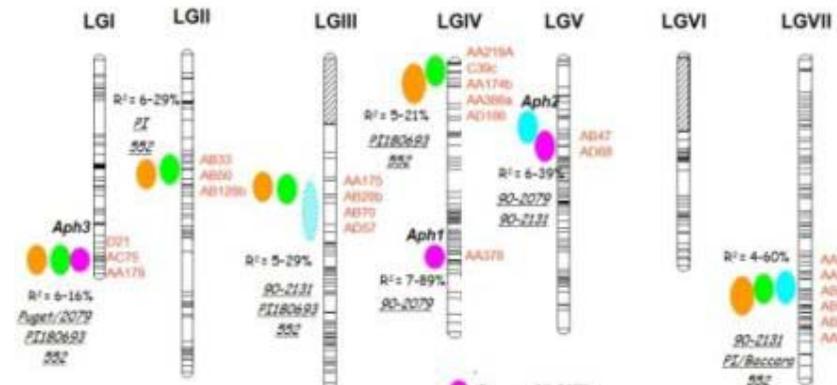
Délivrables attendus pour le projet :

- Approche séquençage cDNA normalisés (Biogemma) : 30 k SNP issus de cDNA full length
- Approche séquençage gDNA enrichi en séquences géniques (INRA) : 15 k SNP ordonnés sur le génome du pois
- Génotypage 10 OPA vera-code : 3k SNP cartographiés

Objectif Scientifiques

- ❖ Saturer la carte génétique du pois en marqueurs SNP
- ❖ Réduire la taille des intervalles de confiance des QTL
- ❖ Établir des ponts entre les cartes issues de différentes populations de RILs
- ❖ Suivre efficacement les allèles favorables aux QTL dans des populations de sélection

assistée par marqueurs (SAM)



- ❖ Cloner certains QTL ciblés d'intérêt.
- ❖ Développer des projets de génétique d'association pour l'identification de nouveaux QTL
- ❖ Préciser la synténie entre le pois et l'espèce légumineuse modèle *M. truncatula* dans les

régions génomiques comportant des QTL

Approche cDNA normalisés (Biogemma)

Délivrable attendu : 30 k SNPs issus de cDNA full length

Développement rapide axé sur des technologies et méthodes validées chez BGA

- ❖ Développement rapide axé sur des technologies et méthodes validées chez BGA
- ❖ Matériel:
 - ❖ 8 géotypes représentant la diversité pois
 - ❖ 2 parents de pop de cartographie (Champagne, Tèreèse)
- ❖ Décomplexification du génome:
 - ❖ séquençage cDNA: séquences exoniques exprimées
 - ❖ Normalisation des banques: lissage des variations de niveau d'expression
- ❖ Séquençage « long read » 454
- ❖ Assemblage et détection des SNP avec le pipeline de bioinformatique BGA

Bilan Approche cDNA normalisés (Biogemma)

- ❖ **68850 contigs** ont été assemblés, dont **50636** ont un **hit** sur le génome de **Medicago truncatula**
- ❖ **74747 SNP** ont été découverts, dont **35455 SNP robustes** caractérisés comme suit :
 - Répartis **sur 10369 contigs pois**, dont 10111 présentant un hit sur le génome de *Medicago truncatula*
 - **11803 SNP** polymorphes entre les lignées Champagne et Térése
 - **26100 SNP** pour lesquels au moins **5 génotypes sur 8 séquencés**
 - **20504 SNP** apportés par un seul génotype (dont **8900 SNP** apportés par Champagne)
- ❖ **1920 SNP** (5 OPA Vera-code) choisis parmi les 35k développés ont été génotypés, et **1360 cartographiés.**

Duarte *et al.* Submitted

Approche gDNA enrichi en séquences géniques (INRA)

Délivrable attendu : 15 k SNP « génomiques » ordonnés sur le génome du pois

VERROUS :

POIS = **gros génome** (4.5 Gb) pour lequel il n'existe **pas de séquence de référence** et composé majoritairement de **séquences répétées**.

DEVELOPPEMENTS PREVUS POUR LE PROJET :

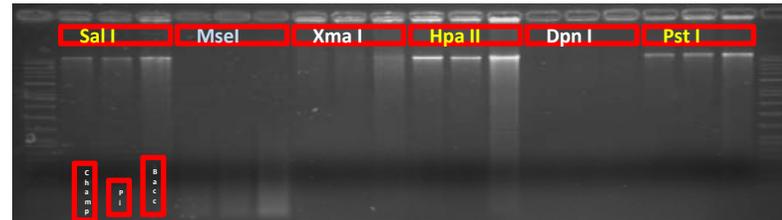
1. Approche de réduction de génome par methyl-filtration (INRA IGEPP / GeT-PlaGe)
2. Read2SNP (INRIA équipe Genscale Pierre Peterlongo / Raluca Uricaru)
3. Cartographie par séquençage d'une pop de 48 RILs de Pois

Stratégie Réduction de génome par methyl-filtration

Développement méthyl-Filtration sur ADN genomique de Pois

1. Réduction de génome par enzymes de restriction sensibles à la méthylation et sizing des fragments d'intérêt

Choix de 2 enzymes parmi 12 testées :
HpaII (C/CGG) et PstI (CTGCA/G)



2. Test de la stratégie méthyl filtration par séquençage sur 1 ligne de 454 GS FLX ; d'1 banque (1 des 2 parents) ... (2011)



3. méthyl filtration puis séquençage (HiSeq2500; Illumina) de 60 banques de pois :
 - a. 12 lignées pures de pois, dont les 2 parents d'1 pop de RILs, pour recherche de SNP (2012-2013)
 - b. 48 RILs pour ordonnancement des SNP (bin mapping) (2013-2014)



Test de la stratégie par séquençage d'1 banque sur 1 ligne de 454 GS FLX

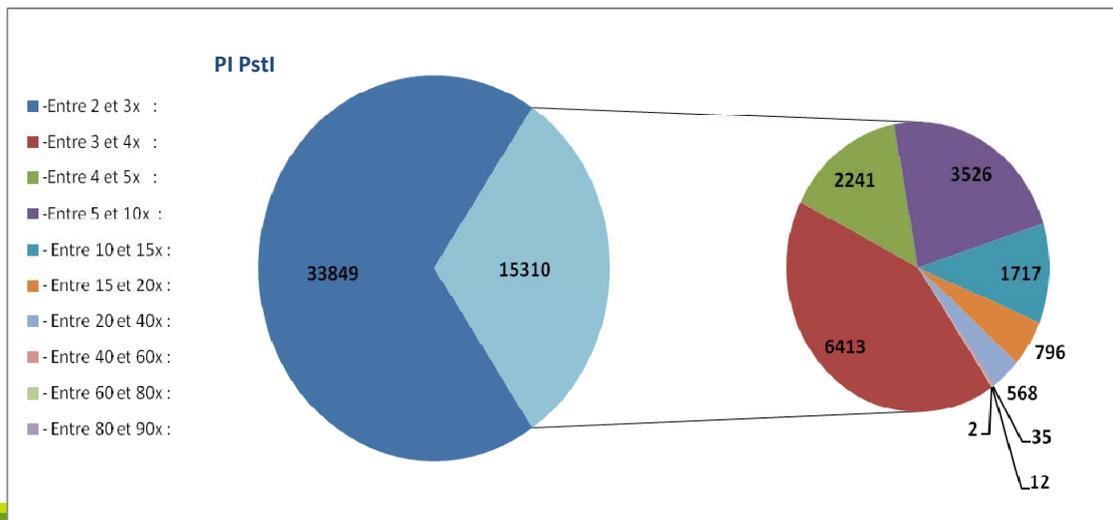
- ❖ Le process « methyl-filtration » sur gDNA CTAB a été globalement validée sur 454 pour les deux enzymes

Malgré des points restant à préciser pour la suite sur projet:

- Séquençage concomitant des banques « PstI » et « HpaII » pour les 60 lignées ?
 - Ou séquençage d'une banque (« PstI » ?) pour les 60 lignées, analyse ... et séquençage de l'autre (« HpaII » ?) seulement si nécessaire ?
 - Intérêt (?) d'une extraction de noyaux préalable pour HpaII ?
 - Adaptation aux spécificités du process Hiseq2000 du protocole validé sur 454 = Verrous techniques à lever pour la suite sur projet
- ❖ **Feu vert pour tester la manip de séquençage sur Hiseq2000 Illumina sur 4 banques Baccara/PstI, Baccara/HpaII, PI/PstI, PI/HpaII**

Séquençage Hiseq2000 4 banques/lanes : Problème de profondeur des contigs

- ❖ Couverture des Contigs Trop faible pour chaque lignée :
 - ❖ seulement 1500 à 3000 contigs ont une couverture > 10x
 - ❖ 15 000 contigs ont une couverture > 3x
 - ❖ 40 000 à 50 000 contigs ont une couverture > 2x
- ❖ Or, Il faudrait AU moins 20k contigs de couverture > 10x pour la recherche et validation de SNP ; et Au moins 20k contigs de couverture > 4x pour le bin-mapping



Recherches de SNP Bacc/PI avec Read2SNP :

faible nombre de SNP trouvés
cohérent avec la faible profondeur
des contigs :

Tests avec cov_min=5:

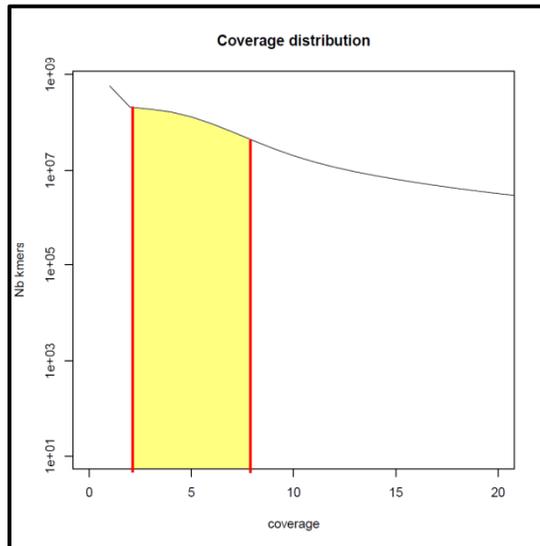
Nb SNP read2SNP filtrés = 3919

Orientations Développements 02/2013

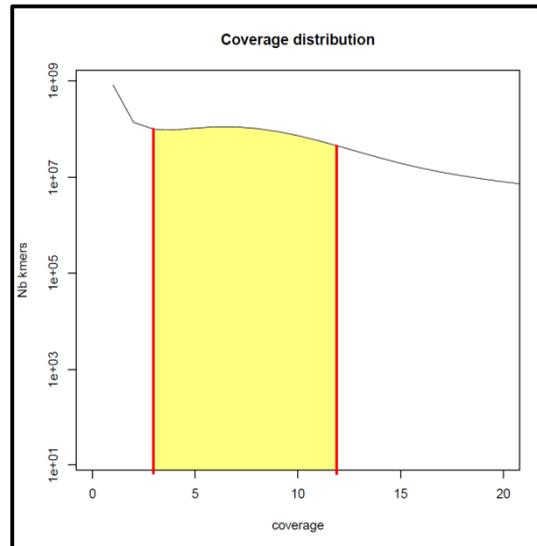
- ❖ Ne plus travailler qu'avec l'enzyme PstI
- ❖ Lancer en pilote les manip sur 4 lignées : Baccara, PI, Champagne et Terese
- ❖ Réduire la fenêtre de sizing pour obtenir une plage identique à celle des manip 454(plage 250- 2000 pb)
- ❖ Conserver l'étape de sonication après sizing : l'idéal aurait été de la supprimer mais les fragments sizés restent trop long pour être séquencés tels quels
- ❖ Augmenter la profondeur en passant 1 banque par lane (contre 4 précédemment)
... tout en sachant qu'on n'aura les moyens de ne passer au mieux que 2 banques par lane pour la pop de 48 Rils
- ❖ Analyse de ces 4 banques avant de décider si/comment on passe à grande échelle sur les 48 Rils et/ou sur 4 à 8 lignées supplémentaires

Résultats Couverture

Zoom des graphes 2 banques / lane et 1 banque / lane.

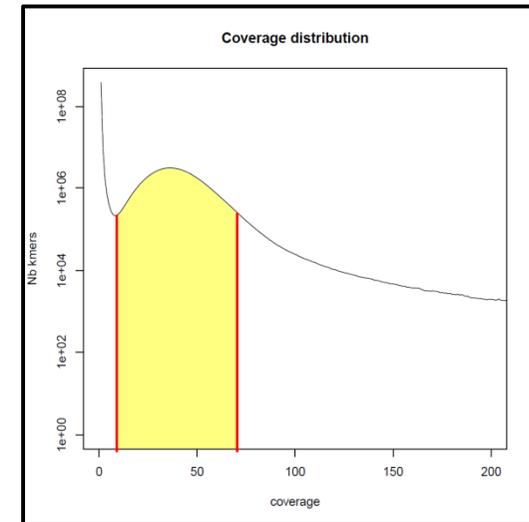


2 banques / 1 lane :
Couverture moyenne
~ 5x



1 banque / 1 lane :
Couverture moyenne
~ 7,5x

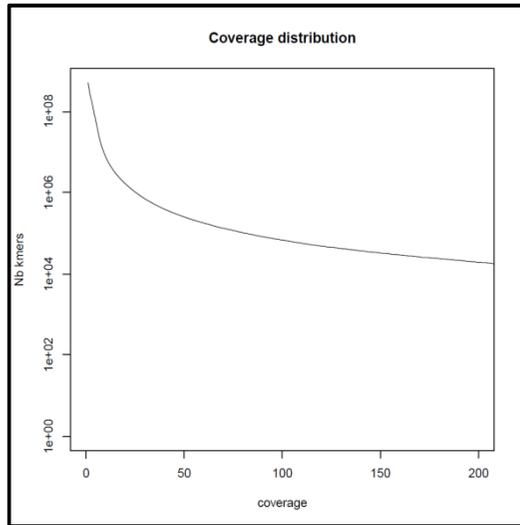
Couvertures des kmers calculer avec le logiciel DSK (kmer = 27)



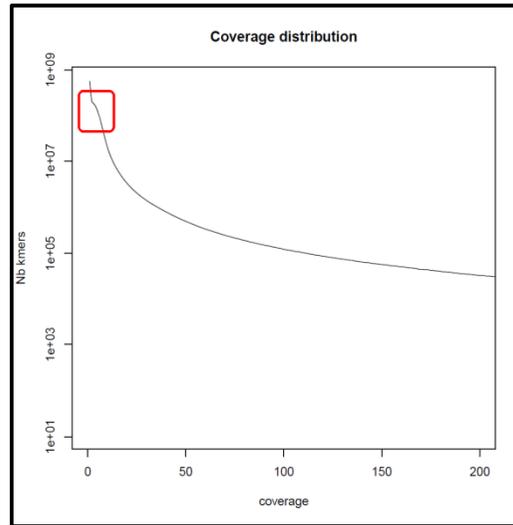
Exemple de graphe que l'on peut obtenir

- Axe Y : nombre de kmer couvert à une profondeur donnée
- Axe X : couverture (profondeur)

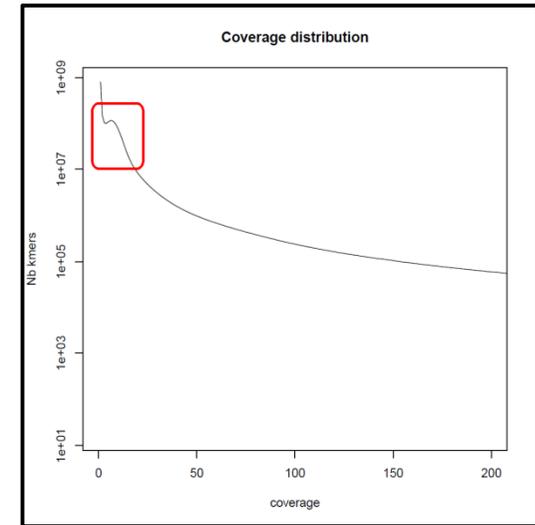
Résultats Couverture



4 banques / 1 lane :
Aucun pic de
couverture.
La couverture est
vraiment trop faible.



2 banques / 1 lane :
Un petit pic de
couverture.
La couverture reste
néanmoins faible.



1 banque / 1 lane :
Pic de couverture.
Couverture correct
en théorie
“whole genome”

Read2SNP

Available :

sequenced reads

1 to n sets (replicates, strains, individuals, ...)

Non available:

reference genome (close or good)

Need:

SNPs with their coverage/quality in each set

(No physical location needed)

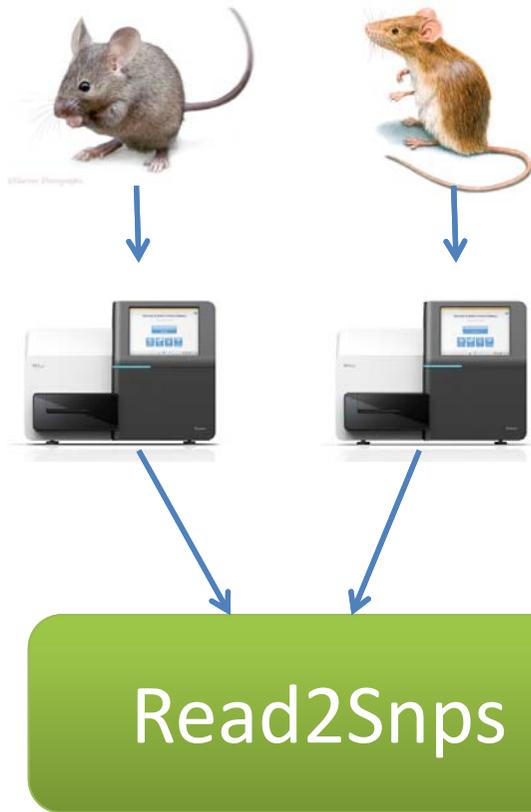
<http://colibread.inria.fr/read2snps/>



Colib'read



Read2SNP Overview



```
>SNP_higher_path_99994|low|left_contig_length_52|right_contig_length_70|C1_14|C2_0|Q1_64|Q2_0|rank_1.00000  
AGTCCTTTAAATATATAAAAAATTAACAAAACGGGAGAAAAGAAATATGGAAAAGGGCTTG  
>SNP_lower_path_99994|low|left_contig_length_52|right_contig_length_70|C1_0|C2_27|Q1_0|Q2_64|rank_1.00000  
AGTCCTTTAAATATATAAAAAATTAACAAAATGGGAGAAAAGAAATATGGAAAAGGGCTTG  
>SNP_higher_path_999949|high|left_contig_length_351|right_contig_length_54|C1_11|C2_0|Q1_63|Q2_0|rank_1.00000  
ACCAGAAGACCAAGGCCAAGGGTTCACTAAATTGATTTCTCCTGAGCCTTTTATCTGTGAC  
>SNP_lower_path_999949|high|left_contig_length_351|right_contig_length_54|C1_0|C2_38|Q1_0|Q2_62|rank_1.00000  
ACCAGAAGACCAAGGCCAAGGGTTCACTAAAGTTGATTTCTCCTGAGCCTTTTATCTGTGAC  
>SNP_higher_path_999916|high|left_contig_length_424|right_contig_length_7|C1_27|C2_0|Q1_57|Q2_0|rank_1.00000  
GGAAATTTGTTCTGCCTAATGGACTATAGCCATTATTACTTTTTCTCAGATACACCA  
>SNP_lower_path_999916|high|left_contig_length_424|right_contig_length_7|C1_0|C2_30|Q1_0|Q2_64|rank_1.00000  
GGAAATTTGTTCTGCCTAATGGACTATAGCCATTATTACTTTTTCTCAGATACACCA  
>SNP_higher_path_999871|high|left_contig_length_9|right_contig_length_131|C1_27|C2_0|Q1_60|Q2_0|rank_1.00000  
CTTGAGACTCGGGTTCCCAACCGTTCCAGAAGCAACGCCGATACTCCATAGGAGGTCAA  
>SNP_lower_path_999871|high|left_contig_length_9|right_contig_length_131|C1_0|C2_53|Q1_0|Q2_58|rank_1.00000  
CTTGAGACTCGGGTTCCCAACCGTTCCAGAAGCAACGCCGATACTCCATAGGAGGTCAA  
>SNP_higher_path_99978|high|left_contig_length_1025|right_contig_length_397|C1_16|C2_0|Q1_59|Q2_0|rank_1.00000  
AAAAGACAAGGATCTTTAGAAGTTTCACTCTAAAGGGAACATATGTTGATCTCTT  
>SNP_lower_path_99978|high|left_contig_length_1025|right_contig_length_397|C1_0|C2_26|Q1_0|Q2_60|rank_1.00000  
AAAAGACAAGGATCTTTAGAAGTTTCACTCTAAAGGGAACATATGTTGATCTCTT  
>SNP_higher_path_999759|low|left_contig_length_21|right_contig_length_10|C1_0|C2_20|Q1_0|Q2_67|rank_1.00000  
AATATAAGTTAGGATATTAATAATTCGCATAAGATACAGGAATTAATAAAAAATTTCTTTA  
>SNP_lower_path_999759|low|left_contig_length_21|right_contig_length_10|C1_12|C2_0|Q1_66|Q2_0|rank_1.00000  
AATATAAGTTAGGATATTAATAATTCGCATTAGATACAGGAATTAATAAAAAATTTCTTTA  
>SNP_higher_path_999758|low|left_contig_length_10|right_contig_length_25|C1_0|C2_37|Q1_0|Q2_64|rank_1.00000  
AAGGAAAAGAGGAGAAAAGAAAAGAAAGCTAGAAATAGTTGTACTCAAGACCTCAAGACC  
>SNP_lower_path_999758|low|left_contig_length_10|right_contig_length_25|C1_16|C2_0|Q1_61|Q2_0|rank_1.00000  
AAGGAAAAGAGGAGAAAAGAAAAGAAAGCTGGAATAGTTGTACTCAAGACCTCAAGACC
```

Parameters:

- k size of the used k -mers. Generated paths are of length $2k-1$
- c minimal coverage

Read2SNP Output

```
>SNP_higher_path_99994|low|left_contig_length_52|right_contig_length_70|C1_14|C2_0|Q1_64|Q2_0|rank_1.00000
AGTCCTTTAAATATATAAAAATTAACAAAACGGGAGAAAAGAAATAATGGAAAAGGGCTTG
>SNP_lower_path_99994|low|left_contig_length_52|right_contig_length_70|C1_0|C2_27|Q1_0|Q2_64|rank_1.00000
AGTCCTTTAAATATATAAAAATTAACAAAATGGGAGAAAAGAAATAATGGAAAAGGGCTTG
>SNP_higher_path_999949|high|left_contig_length_351|right_contig_length_54|C1_11|C2_0|Q1_63|Q2_0|rank_1.00000
ACCAGAAGACCAAGGCCAAGGGTTCCTAAATTGATTTCTCCTGAGCCTTTTATCTGTGAC
>SNP_lower_path_999949|high|left_contig_length_351|right_contig_length_54|C1_0|C2_38|Q1_0|Q2_62|rank_1.00000
ACCAGAAGACCAAGGCCAAGGGTTCCTAAAGTTGATTTCTCCTGAGCCTTTTATCTGTGAC
>SNP_higher_path_999916|high|left_contig_length_424|right_contig_length_7|C1_27|C2_0|Q1_57|Q2_0|rank_1.00000
GGAAATTTGTTCTTGCCTAATGGACTATAGCCCATTATTACTTTTTCTTCAGATACACCA
>SNP_lower_path_999916|high|left_contig_length_424|right_contig_length_7|C1_0|C2_30|Q1_0|Q2_64|rank_1.00000
GGAAATTTGTTCTTGCCTAATGGACTATAGTCCATTATTACTTTTTCTTCAGATACACCA
>SNP_higher_path_999871|high|left_contig_length_9|right_contig_length_131|C1_27|C2_0|Q1_60|Q2_0|rank_1.00000
CTTGAGACTCGGGTCCCCAACCGTTCAGAAGCAACGCCGATACTCCCATAGGAGGTCAA
>SNP_lower_path_999871|high|left_contig_length_9|right_contig_length_131|C1_0|C2_53|Q1_0|Q2_58|rank_1.00000
CTTGAGACTCGGGTCCCCAACCGTTCAGAGGCAACGCCGATACTCCCATAGGAGGTCAA
>SNP_higher_path_99978|high|left_contig_length_1025|right_contig_length_397|C1_16|C2_0|Q1_59|Q2_0|rank_1.00000
AAAAGACAAAGGATACTTTAGAAGTTTCACCTCTAAAGGGAACATATGTATTGCATCTCTT
>SNP_lower_path_99978|high|left_contig_length_1025|right_contig_length_397|C1_0|C2_26|Q1_0|Q2_60|rank_1.00000
AAAAGACAAAGGATACTTTAGAAGTTTCACCTCTAAAGGGAACATATGTATTGCATCTCTT
>SNP_higher_path_999759|low|left_contig_length_21|right_contig_length_10|C1_0|C2_20|Q1_0|Q2_67|rank_1.00000
AATATAAGTTAGGATATTAATTTCTGCATAAGATACAGGAATTAATAAATTATTCTTA
>SNP_lower_path_999759|low|left_contig_length_21|right_contig_length_10|C1_12|C2_0|Q1_66|Q2_0|rank_1.00000
AATATAAGTTAGGATATTAATTTCTGCATTAGATACAGGAATTAATAAATTATTCTTA
>SNP_higher_path_999758|low|left_contig_length_10|right_contig_length_25|C1_0|C2_37|Q1_0|Q2_64|rank_1.00000
AAGCAAAAGCAGCAAAAAGAAAAGAAAGCTAGAAATAGTTGTACCTCAAGACCTCAAGACC
```

Read2SNP Output

```
>SNP_higher_path_999949|high|left_contig_length_351|right_contig_length_54|C1_11|C2_0|Q1_63|Q2_0|rank_1.00000  
ACCAGAAGACCAAGGCCAAGGGTTCACTAATGATTTCTCCTGAGCCTTTTATCTGTGAC  
>SNP_lower_path_999949|high|left_contig_length_351|right_contig_length_54|C1_0|C2_38|Q1_0|Q2_62|rank_1.00000  
ACCAGAAGACCAAGGCCAAGGGTTCACTAGTGTGATTTCTCCTGAGCCTTTTATCTGTGAC
```

- Complexity: high/low
- Length left and right extensions (not shown here)
- C1: coverage of the path in the first read set
- C2: ...
- Q1: phred quality of the polymorphic base
- Q2: ...
- Rank:
 - 1 = SNP discriminative between read sets
 - 0 = SNP non discriminative between read sets

Filtres « post Read2SNPs »

Paramètres Read2SNP :

- Kmer27
- Couverture minimale pour 1 lignée c3

Résultats « Bruts » read2SNP :

- 4 banques/lanes : 874 056 SNPs
- 2 banques/lanes : 2 129 420 SNPs
- 1 banque/lane : ? (analyse en cours)

... nécessité de filtrer !!!

Filtre1/ Le SNP est au minimum couvert à **5x** pour au moins l'une des 4 lignées.

Filtre2/ Suppression des snps hétérozygote (j'autorise **1 erreur par tranche de 10x**) :

	Bacc	Champ	PI	Ter
A	14	15	10	14
T	10	12	15	11

	Bacc	Champ	PI	Ter
A	0	54	14	1
T	15	2	0	32

Filtres « post Read2SNPs »

Paramètres Read2SNP :

- Kmer27
- Couverture minimale pour 1 lignée c3

Résultats « Bruts » read2SNP :

- 4 banques/lanes : 874 056 SNPs
- 2 banques/lanes : 2 129 420 SNPs
- 1 banque/lane : ? (analyse en cours)

... nécessité de filtrer !!!

Filtre1/ Le SNP est au minimum couvert à **5x** pour au moins l'une des 4 lignées.

Filtre2/ Suppression des snps hétérozygote (j'autorise **1 erreur par tranche de 10x**) :

	Bacc	Champ	PI	Ter
A	14	15	10	14
T	10	12	15	11

	Bacc	Champ	PI	Ter
A	0	54	14	1
T	15	2	0	32

Filtres « post Read2SNPs »

Paramètres Read2SNP :

- Kmer27
- Couverture minimale pour 1 lignée c3

Résultats « Bruts » read2SNP :

- 4 banques/lanes : 874 056 SNPs
- 2 banques/lanes : 2 129 420 SNPs
- 1 banque/lane : ? (analyse en cours)

... nécessité de filtrer !!!

Filtre3/ couverture allèle minoritaire : Ok si $\max1 \leq (\max2/2)$ ou $\max2 \leq (\max1/2)$.

	Bacc	Champ	PI	Ter
A	1	4	1	8
T	34	48	0	1

	Bacc	Champ	PI	Ter
A	1	2	21	31
T	32	44	2	1

Filtres « post Read2SNPs »

Paramètres Read2SNP :

- Kmer27
- Couverture minimale pour 1 lignée c3

Résultats « Bruts » read2SNP :

- 4 banques/lanes : 874 056 SNPs
- 2 banques/lanes : 2 129 420 SNPs
- 1 banque/lane : ? (analyse en cours)

... nécessité de filtrer !!!

Filtre3/ couverture allèle minoritaire : Ok si $\max1 \leq (\max2/2)$ ou $\max2 \leq (\max1/2)$.

	Bacc	Champ	PI	Ter
A	1	4	1	8
T	34	48	0	1

	Bacc	Champ	PI	Ter
A	1	2	21	31
T	32	44	2	1

Résultats SNP Read2SNPs après filtrage

	4 banques / lane			2 banques / lane			1 banque / lane		
pour chaque SNP, données de séquence disponible pour chacune des 4 lignées									
	Total SNPs	Designable "≥ 50"	Designable "≥ 100"	Total SNPs	Designable "≥ 50"	Designable "≥ 100"	Total SNPs	Designable "≥ 50"	Designable "≥ 100"
Total SNP Polymorphes "génotypage 4 lignées"	6861	5341	3083	43512	29549	13829			
SNP Polymorphes Baccara-Pi	4545	3548	2047	30013	20399	9398			
SNP Polymorphes Baccara-Champagne	5130	3981	2272	33380	22646	10554			
SNP Polymorphes Baccara-Terese	1337	1038	625	8520	5753	2747			
SNP Polymorphes Champagne-Pi	2805	2195	1283	16261	11085	5304			
SNP Polymorphes Champagne-Terese	5363	4163	2367	34010	23071	10729			
SNP Polymorphes Pi-Terese	4764	3712	2124	30699	20824	9631			

43k SNP (29k « designables ») dont **30k SNP** (20k « designables ») polymorphes entre les parents de la pop de RILs à cartographier par séquençage
...Avec un effort de séquençage à 2 banques par lane

Conclusion - Perspectives

1. Malgré les efforts sur le sizing des banques, la **« filtration »** parait peu efficace, avec une profondeur moyenne plus proche de celle d'un séquençage « whole genome » que d'un séquençage ciblé.
2. Les premières analyses (toujours en cours) nous indiquent des **banques simili « whole genome »** avec malgré tout un **enrichissement 2 à 5x de la portion « methyl-filtrée »** qui devrait représenter **1 à 10% du génome**.
3. Dans ces conditions, un **effort de séquençage de 24 lanes Hiseq2000** apparait **nécessaire** (et suffisant) **pour la cartographie par séquençage de la population de 48 RILs** et ainsi répondre aux objectifs initiaux du projet (développement et bin mapping environ 15K SNP génomiques).

Dans l'état actuel, la technique ne semble donc pas complètement au point et n'est probablement pas la plus élégante ni la moins onéreuse, mais elle permet de répondre aux objectifs du projet.

MERCI

INRA UMR 1349 IGEPP :

- Susete Alves Carvalho (Genscale / Genouest) – analyses bioinfo
- Raluca Uricaru (équipe GenScale INRIA) – développement read2SNP
- Alain Baranger



INRA GeT-PlaGe :

- Emeline Lhuillier – développement Methyl-filtration / Banques Hiseq2000
- Jérôme Luch – développement 454
- Olivier Bouchez

INRIA-IRISA :

- Olivier Collin (plate-forme Genouest)
- Delphine Naquin (plate-forme Genouest) – analyses bioinfo
- Pierre Peterlongo (équipe Genscale)



INRA CNRGV:

- Arnaud Bellec – tests préliminaires Methyl-filtration

BIOGEMMA Auvergne :

- Jorge Duarte
- Nathalie Rivière

