

UNIVERSITÉ GRENOBLE-ALPES

THÈSE

Pour obtenir le grade de

DOCTEUR DE L'UNIVERSITÉ DE GRENOBLE-ALPES

Spécialité : **Modèles, méthodes et algorithmes en biologie, santé et environnement**

Arrêté ministériel : ?

Présentée par

Thomas Karaouzene

Thèse dirigée par **Pierre Ray**

Thèse co-dirigée par **Nicolas Thierry-Mieg**

préparée au sein du laboratoire
et de l'école doctorale "**Ingénierie de la Santé, de la Cognition et Environnement**" (EDISCE)

Écrire le titre de la thèse ici

Thèse soutenue publiquement le 31 octobre 2017,
devant le jury composé de :



**Université
Grenoble
Alpes**

Préface

This is an example of a thesis setup to use the reed thesis document class (for LaTeX) and the R bookdown package, in general.

Table des matières

Chapitre 1 : Delete line 6 if you only have one advisor	1
Remerciements	3
Résumé	5
Chapitre 2 : Introduction	7
Chapitre 3 : Investigation génétique et physiologique de la globo- zoospermie	9
Chapitre 4 : Mise en place d’une stratégie pour l’analyse des données exomiques – application en recherche clinique	11
4.1 Intro	12
4.2 Résultats	13
4.2.1 Description de la pipeline	13
4.2.2 Utilisation du pipeline dans des cas familiaux :	17
Résultats des exomes	19
Discussion	37
4.2.3 Etude d’une large cohorte de patients MMAF	40
Description de la cohorte	40
Application de la pipeline - Résultats	41
Analyse des listes de gènes	43
Discussion	78
4.3 Conclusion	81
Chapitre 5 : MutaScript	83
Conclusion	85
Chapitre 6 : Tables des analyses d’exomes	87
References	89

Liste des tableaux

4.1	Liste simplifiée des conséquences prédites par VEP avec leur description et impact associée	16
4.2	Tableau récapitulatif des familles séquencées et de leur phénotype . .	18
4.3	Récapitulatif des variants ayant passé l'ensemble des filtres pour chaque famille	32
4.4	Liste des différents projets de séquençages effectués	40
4.5	Liste des gènes prédits comme faisant partie du ciliome humain retrouvés mutés à l'état homozygote chez nos patients	46
4.6	Analyse n°1 : Liste des patients portant au moins un variant homozygote tronquant sur le gène *DNAH1*	51
4.7	Analyse n°1 : Liste des patients portant au moins un variant homozygote non tronquant sur le gène *DNAH1*	52
4.8	Analyse n°1 : Liste des patients portant au moins deux variant hétérozygotes sur le gène *DNAH1*	53
4.9	Analyse n°2 : List des gènes présents dans la liste ciliome sur lesquels au moins deux patients portent une mutation tronquante homozygote	57
4.10	Analyse n°2 : Liste des patients portant un variant non troquant homozygote sur un des gènes suivant : ARMC2, CCDC146, CFAP44 et TTC29	58
4.11	Analyse n°3 : List des gènes sur lesquels au moins deux patients portent une mutation tronquante non présents dans la liste ciliome	63
4.12	Analyse n°3 : Liste des patients portant un variant non troquant homozygote sur un des gènes suivant : *BAZ1A*, *CCDC129*, *CFAP43*, *FSIP2*, *ICA1*, *NACA*, *SART3* et *TRAV26-1*	64
4.13	Analyse n°3 : Liste des patients portant un variant non troquant homozygote sur un des gènes suivant : *BAZ1A*, *CCDC129*, *CFAP43*, *FSIP2*, *ICA1*, *NACA*, *SART3* et *TRAV26-1*	66
4.14	Analyse n°4 : List des gènes présents dans la liste ciliome sur lesquels un seul patient portent une mutation homozygote tronquante	70
4.15	Analyse n°4 : Liste des patients portant au moins deux variants hétérozygotes sur un des gènes suivant : *C21orf59*, *C6orf118*, *CCDC65* et *SPEF2*	71
4.16	Analyse n°5 : Liste des patients portant au moins un variant homozygote tronquant sur le gène sur l'un des 59 genes identifiés dans cette analyse	75

4.17 Analyse n°5 : Liste des patients portant au moins un variant homozygote non tronquant sur le gène sur l'un des 59 genes identifiés dans cette analyse	77
--	----

Table des figures

4.1	Listes des différentes conséquences prédites par VEP et leur positionnement sur le transcrit	14
4.2	Processus simplifié du contrôle qualité des *reads*	20
4.3	Contrôle qualité des variants appelés	23
4.4	Annotation des variants	25
4.5	Filtrage des transcrits jugés "non pertinents" et des variants les chevauchant	27
4.6	Nombre d'individus composant la cohorte contrôle de chaque famille .	28
4.7	Comparaison de l'efficacité de chacun des six filtres utilisés	30
4.8	Expression tissulaire des gènes *SPINK2* et *GUF1*	33
4.9	Expression tissulaire du gène *PLCZ1*	34
4.10	Expression tissulaire des gènes retenus pour la famille MMAF3	35
4.11	Expression tissulaire du gène *TGIF2*	36
4.12	Nombre de gènes passant l'ensemble des filtres par famille	39
4.13	Résultats de l'appel des variants par individus et par projet de séquençage	42
4.14	Résultats de l'étape de filtrage	43
4.15	Répartition des gènes retrouvés mutés chez nos patients dans les différentes classes de la liste du cillione	45
4.16	Représentation schématique desTODO	48
4.17	Analyse du gène *DNAH1*	50
4.18	Analyse des gènes sélectionnés dans l'Analyse n°1	56
4.19	Analyse des gènes sélectionnés dans l'Analyse n°2	62
4.20	Analyse des gènes sélectionnés dans l'Analyse n°4	69

Chapitre 1

Delete line 6 if you only have one advisor

Remerciements

Résumé

Chapitre 2

Introduction

Chapitre 3

Investigation génétique et physiologique de la globozoospermie

Chapitre 4

Mise en place d'une stratégie pour
l'analyse des données exomiques –
application en recherche clinique

4.1 Intro

En 2011, les bases moléculaires d'environ 3700 pathologies à transmission Mendélienne avaient été élucidées cependant, pour une quantité équivalente de pathologies Mendéliennes (ou suspectées de l'être) cette cause reste un mystère (Amberger, Bocchini, & Hamosh, 2011). Avec plusieurs centaines de pathologies caractérisées depuis 2010 (S. B. Ng et al., n.d.), les séquençages WGS et WES ont, depuis leur émergence, révolutionnés les méthodes de recherche dans le cadre d'étude phénotype-génotype en permettant de manière rapide et à moindre coup le séquençage de la quasi totalité des gènes humains. Dès lors, le défi de ces analyses n'est plus le séquençage de l'ADN mais l'interprétation des données massives produites. En effet, l'un des plus grand challenge des analyses phénotype-génotype réalisées par WES réside dans la quantité importante de variant portés par chaque individu s'élevant à plusieurs dizaine de milliers lorsque l'on compare avec le génome de référence. Même après avoir retiré les variants retrouvés fréquemment dans la population générale, des méthodes additionnelles sont nécessaires pour prédire, parmi les variants restant, lesquels induisent des conséquences fonctionnelles sérieuses afin de le prioriser (Pelak et al., 2010). De nombreux logiciels tel que Variant effect predictor (W. McLaren et al., 2016), snpEFF ou encore Annovar permettent d'identifier quels sont les variants qui ont un effet tronquant sur la protéine. Cependant, avec en moyenne 165 variants homozygotes ayant un effet tronquant sur la protéine retrouvés dans chaque exome (Pelak et al., 2010) ces méthodes, bien qu'efficaces sont souvent insuffisantes.

D'autres logiciels tel que Exomiser vont, à partir d'une liste de variant **déjà** appelés effectuer les étapes d'annotation, de filtrage et de priorisation. Malgré l'efficacité de ces logiciels, aucun d'entre eux ne couvrent l'ensemble des étapes allant de l'alignement des *reads* à la priorisation des variants. La plupart ayant pour point de départ une liste de variants appelés en amont. Ils ne contrôlent donc en aucune manière les étapes d'alignement et d'appel des variants. Or, comme il a été dit plus tôt, ces deux étapes constituent la bases de l'analyse [].

Ce chapitre décrit à la fois la constitution d'un pipeline d'analyse des données de séquençage exomique et son utilisation dans le cadre de la recherche de mutations entraînant différents phénotypes d'infertilité d'une part de cas familiaux composés de duos ou trio et pour finir d'une large cohorte d'individus non apparentés présentant tous le même phénotype.

4.2 Résultats

Dans cette partie, nous allons, après avoir décrit notre pipeline, détailler les résultats de l'analyse des données de WES de 75 patients tous atteints d'un phénotype d'infertilité. Ces études seront séparées en deux parties distinctes, la première se concentrera sur l'étude de 6 familles incluant 13 de ces patients. La seconde portera sur l'analyse des 62 patients restant étant tous non-apparentés et présentant un phénotype MMAF.

4.2.1 Description de la pipeline

Après avoir été séquencés, les données recueillies pour ces patients sont traitées au sein de la même pipeline d'analyse qui comprend quatre étapes allant de l'alignement des *reads* au filtrage des variants :

1. **L'alignement** : L'alignement des *reads* le long du génome de référence (hg19 / GRCh37) est effectué par le logiciel MAGIC (Su et al., 2014). Afin d'écartier toute ambiguïté au moment de l'interprétation de l'alignement, l'intégralité des *reads* dupliqués et / ou s'alignant à plusieurs zones du génome seront filtrés et ne seront donc pas pris en compte pour l'ensemble des analyses en aval. Suite à cela, MAGIC va produire quatre comptages pour chaque position couverte du génome : R+, V+, R- et V- :
 - a. **R+ et R-** : Ces deux comptages correspondent au nombre de *reads forward* (+) et *reverse* (-) sur lesquels est observé l'allèle de **référence** (R) à une position donnée.
 - b. **V+ et V-** : À l'inverse de R+ et R-, ces comptages correspondent au nombre de *reads forward* et *reverse* sur lesquels est observé un allèle de **variant** (V) à une position donnée.
2. **L'appel des variants** : Comme nous l'avons vu plus tôt, il est fortement conseillé d'effectuer l'appel des variants en tenant compte de l'aligneur choisi (Nielsen, Paul, Albrechtsen, & Song, 2011, M. A. DePristo et al. (2011), Lunter & Goodson (2011)). C'est pourquoi, nous avons développé notre propre algorithme d'appel des variants spécialement conçu pour l'analyse des données de MAGIC. Ainsi, l'appel des variants sera directement basé sur les quatre comptages vus précédemment. Tout d'abord, les positions ayant une couverture < 10 sur l'un des deux *strands* seront considérées comme de faible qualité, celles ayant une couverture < 10 sur les deux *strands* seront exclues. Ensuite pour chaque variant, des appels indépendants seront effectués pour chaque *strand*. L'appel final sera une synthèse de ces deux appels où seul les cas où ces deux appels sont concordants seront considérés comme de bonne qualité.

3. **L'annotation** : Chaque variant retenu sera ensuite annoté tout d'abord par le logiciel *variant effect predictor* (VEP) (W. McLaren et al., 2016) qui nous indiquera pour chaque variant la conséquence que celui-ci aura sur la séquence codante de l'ensemble des transcrits Ensembl qu'il chevauche (**Figure : 4.1**) (**Table : 4.1**). Ensuite, nous ajouterons récupérerons pour chaque gène son expression tissulaire en nous basant sur les données Ensembl (Aken et al., 2017) générées par le projet Illumina BodyMap qui recense les données RNAseq des gènes humains pour 16 tissus différents. Suite à cela nous ajoutons, lorsque celle-ci est disponible, la fréquence du variant dans les bases de données ExAC (Lek et al., 2016), ESP600 (Exome Variant Server, NHLBI GO Exome Sequencing Project (ESP), Seattle, WA) et 1000Genomes (1000 Genomes Project Consortium et al., 2015) donnant ainsi une estimation de sa fréquence dans la population générale. De même, la particularité de ce pipeline est qu'elle conserve l'ensemble des variants identifiés dans les études effectuées précédemment permettant d'ajouter aux annotations la fréquence d'un variant chez les individus déjà séquencé et donc la fréquence d'un variant dans chaque phénotype étudié créant ainsi une base de données interne qui pourra servir de contrôle dans les études ultérieures.

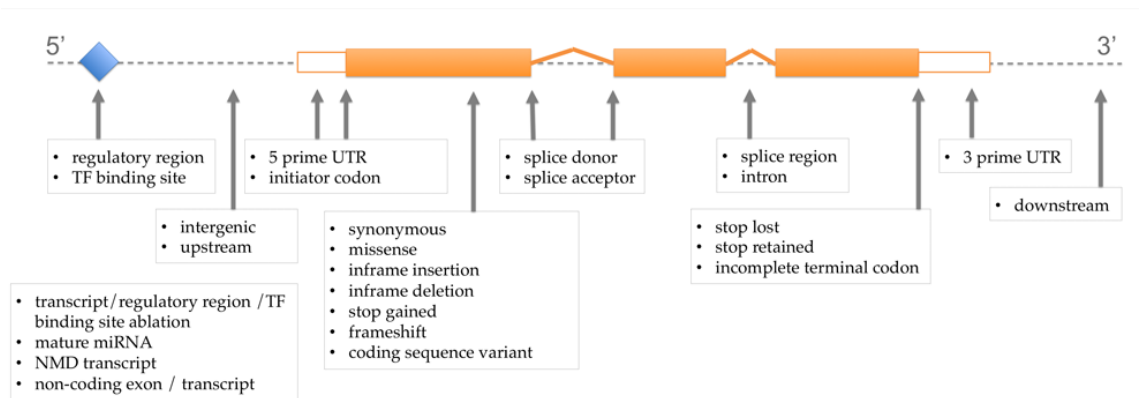


Figure 4.1 – Listes des différentes conséquences prédites par VEP et leur positionnement sur le transcript d'après [VEP site](<http://www.ensembl.org/info/genome/variation/consequences.jpg>)

4. **Le filtrage des variants** : L'étape de filtrage est extrêmement importante si l'on souhaite analyser de manière efficace les données provenant de WES. C'est pourquoi elle occupe une place importante dans notre pipeline. L'intégralité des paramètres de cette étape peuvent être modifiés par l'utilisateur de sorte à faire correspondre les critères de filtre aux besoins de l'étude. Afin de rendre son utilisation la plus efficace possible, nous avons souhaité définir des paramètres par défauts pertinents dans la plupart des études de séquençage exomique de sorte que à moins que le contraire ne soit spécifié, seuls les variants impactant les transcrits codant pour une protéine sont conservés. De même les variants synonymes ou affectant les séquences UTRs sont filtrés ainsi que les variants ayant une fréquence $\geq 1\%$ dans les bases dans l'une des bases de données (ExAC, ESP6500).

ou 1KH). Aussi, pour un phénotype donné, l'ensemble des variants homozygotes observés chez les individus étudiés présentant un phénotype différent sont de même enlevés de la liste finale.

Table 4.1 – Liste simplifiée des conséquences prédites par VEP avec leur description et impact associée

VEP consequence	VEP impact	Description
Splice acceptor / donor	HIGH	A splice variant that changes the 2 base region at the 3' / 5' end of an intron
Stop gained	HIGH	A sequence variant whereby at least one base of a codon is changed, resulting in a premature stop codon, leading to a shortened transcript
Frameshift	HIGH	A sequence variant which causes a disruption of the translational reading frame, because the number of nucleotides inserted or deleted is not a multiple of three
Stop lost	HIGH	A sequence variant where at least one base of the terminator codon (stop) is changed, resulting in an elongated transcript
Start lost	HIGH	A codon variant that changes at least one base of the canonical start codon
Inframe insertion / deletion	MODERATE	An inframe non synonymous variant that inserts / deletes bases into in the coding sequence
Missense	MODERATE	A sequence variant, that changes one or more bases, resulting in a different amino acid sequence but where the length is preserved
Splice region	LOW	A sequence variant in which a change has occurred within the region of the splice site, either within 1-3 bases of the exon or 3-8 bases of the intron
Stop retained	LOW	A sequence variant where at least one base in the terminator codon is changed, but the terminator remains
Synonymous	LOW	A sequence variant where there is no resulting change to the encoded amino acid
5' / 3' prime UTR	MODIFIER	A UTR variant of the 5' / 3' UTR
Intron	MODIFIER	A transcript variant occurring within an intron
NMD transcript	MODIFIER	A variant in a transcript that is the target of NMD
Non coding transcript	MODIFIER	A transcript variant of a non coding RNA gene

4.2.2 Utilisation du pipeline dans des cas familiaux :

Dans cette partie, je me concentre sur l'analyse bioinformatique des résultats des séquençages exomiques effectués entre 2012 et 2014 de 13 individus infertiles provenant de 6 familles différentes. Parmi celles-ci, 3 phénotypes différents ont été observés :

1. **L'Azoospermie** : Comme nous avons pu le voir, l'azoospermie est un phénotype d'infertilité masculine caractérisé par l'absence de spermatozoïde dans l'éjaculât
2. **Échec de fécondation** : Ce phénotype d'infertilité se caractérise par l'incapacité des spermatozoïdes à féconder l'ovocyte.
3. **MMAF** : Le syndrome MMAF (*multiple morphological abnormalities of the sperm flagella*) caractérise comme son nom l'indique les patients présentant une majorité de spermatozoïdes atteints par une mosaïque d'anomalie morphologique du flagelle.

Parmi ces 6 chacune composée de 2 à 3 frères, les familles AZ, FF et MMAF2 présentent un historique de consanguinité, les parents étant soit cousins germains, soit cousins au second degré. La consanguinité favorisant la transmission de variants à l'état homozygote, nous avons décidé, dans un premiers temps de concentrer nos analyses uniquement sur les variants (SNVs et indels) homozygotes pour l'ensemble des familles. Pour les 3 familles n'ayant pas d'historique de consanguinité, ce choix nous permet de réduire la liste des variants candidats de sorte à faciliter les analyses. L'études des variants hétérozygotes sera effectuée *a posteriori* pour les familles dont la cause génétique du phénotype n'a pas pu être identifiée en se limitant aux variants homozygotes. Un récapitulatif des familles et de leur phénotype est disponible dans la table 4.2.

Table 4.2 – Tableau récapitulatif des familles séquencées et de leur phénotype

Family	Consanguinity	Individuals	Phenotype	Year	Place
AZ	Yes	AZ1, AZ2	Azoospermia	2012	Mount Sinai Institut
FF	Yes	FF1, FF2	Fertilization failure	2014	Genoscope (Evry)
MMAF1	No	MMAF1.1, MMAF1.2	MMAF	2014	Genoscope (Evry)
MMAF2	Yes	MMAF2.1, MMAF2.2	MMAF	2014	Genoscope (Evry)
MMAF3	No	MMAF3.1, MMAF3.2	MMAF	2014	Genoscope (Evry)
MMAF4	No	MMAF4.1, MMAF4.2, MMAF4.3	MMAF	2014	Genoscope (Evry)

Résultats des exomes

Résultat de l'alignement Pour rappel, l'alignement consiste à repositionner l'ensemble des *reads* générés au cours de l'étape de séquençage le long d'un génome de référence.

La quantité de *reads* composant les exomes de chaque individu peut varier en fonction de plusieurs paramètres et n'est donc pas égale pour chaque patient bien que l'ordre de grandeur reste le même avec une médiane de 91438630 *reads*. Seuls les deux frères AZ1 et AZ2 se distinguent avec près de 3 fois plus de *reads* que les autres patients. Cette différence peut être expliquée car ces deux patients sont les deux seuls à avoir été séquencés au Mount Sinai Institut or leur protocole d'amplification précède le séquençage contient un nombre de cycles de PCR supérieur à ceux appliqués au Génomex d'Évry où ont été séquencés les autres patients. Il faut noter que ce nombre plus important de *reads* n'est en rien le reflet d'une meilleure qualité. En effet, celui-ci est causé par une grande quantité de *reads* dupliqués qui seront pour la plupart filtrés au cours des analyses ultérieures (**Table : 4.2, Figure : 4.2 - A**).

L'ensemble de nos exomes ayant été réalisés en *paired-end*, les deux extrémités de chaque fragment sont séquencées. Chaque *end* d'un même *read* peut donc être considérée comme un *read* à part entière qui sont alignées **indépendamment** le long du génome de référence. L'information fournie par le *paired-end* n'étant utilisée qu'à *posteriori* en tant que critère qualité. La première étape du contrôle qualité des *reads* consiste à filtrer les *reads* ne s'étant pas alignés sur le génome. Ces *reads* sont extrêmement minoritaires puisqu'ils ne représentent qu'entre 1.2 et 5.5 % des *reads* de nos individus (**Figure : 4.2 - B**).

Une fois cela fait, nous vérifions la "compatibilité" des deux *ends* composant chacun des *reads* s'étant correctement alignés. Un *reads* est dit compatible lorsque les deux *ends* qui le composent s'alignent face à face (une sur le *strand* + et l'autre sur le *strand* -) et couvrent une zone ne faisant pas plus de 3 fois la taille médiane de l'insert. Les *reads* dont les deux *ends* se sont alignés mais ne remplissant pas ces conditions seront dit "Non compatible", ceux dont une seule des deux *ends* s'est alignée seront appelés "orphelins". Dans nos analyses, seuls les *reads* compatibles sont conservés, c'est à dire environ 89.5 % des *reads* s'étant correctement alignés. (**Figure : 4.2 - C**).

La dernière étape de ce contrôle-qualité consiste à analyser le nombre de sites auxquels se sont alignés les *reads*. En effet, certaines zones du génome étant dupliquées, l'une des problématiques des *short-reads* est qu'il est possible que ceux-ci s'alignent à plusieurs régions différentes du génome. Afin d'éviter toute ambiguïté, seuls ceux s'étant alignés sur un site unique sont conservés pour la suite des analyses. Ces *reads* représentent entre 92.3 et 96.9 % des *reads* ayant passé les précédents filtres (**Figure : 4.2 - C**).

Les *reads* ayant passé l'ensemble des critères qualité mentionnés précédemment seront ensuite utilisés pour effectuer l'appel des variants.

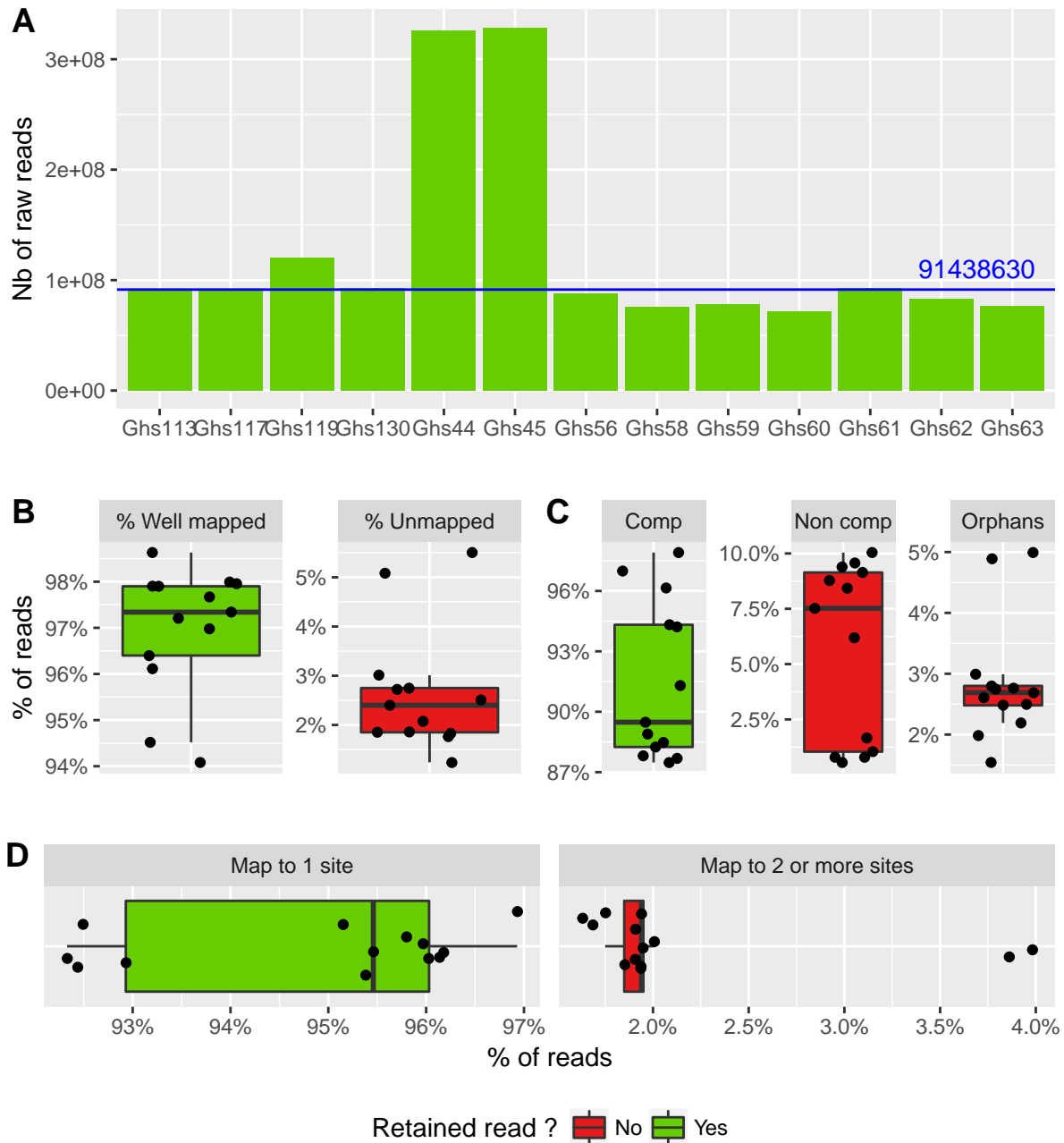


Figure 4.2 – Processus simplifié du contrôle qualité des *reads* : Pour chacun des graphiques, les *reads* représentés en vert sont conservés tandis que ceux en rouge sont filtrés. **A** : Quantité de *reads* bruts générés pour chaque patient au cours de l'étape de séquençage. La médiane des *reads* est représentée en bleue. **B** : Pourcentage pour chaque individu de *reads* s'étant aligné correctement et ne s'étant pas alignés sur le génome de référence. **C** : Distribution pour chaque patient des *reads* compatibles (Comp), non compatibles (Non comp) et orphelins (Orphans). **D** : Présentation pour chaque *reads* du nombre de site auxquels ils s'alignent

Résultat de l'appel des variants Comme dit précédemment, l'appel des variants fait suite à l'alignement et consiste à comparer la séquence d'un individu avec celle d'un génome de référence afin d'en relever les différences. La particularité de notre algorithme d'appel est d'effectuer pour chaque position deux appels indépendants. Le premier sera effectué en utilisant uniquement les *reads forward* et le second le *reads reverse*. Encore une fois, plusieurs filtres sont appliqués de sorte à conserver uniquement les variants les plus qualitatifs.

Tout d'abord, nos appels sont classés en trois catégories :

1. **Les appels *double strand* (DS) :** Qualifie les positions ayant une couverture ≥ 10 sur les deux *strands*. Ces appels sont ceux ayant la meilleure qualité
2. **Les appels *single strand* (SS) :** Ces appels définissent les positions pour lesquels **un des deux *strands*** présentent une couverture ≤ 10 . Dans ce cas, ce *strand* est ignoré et l'appel est effectué uniquement en utilisant le second *strand*.
3. **Les appels *non strand* (NS) :** Les positions NS sont celles pour lesquelles la couverture est ≤ 10 sur les deux *strands*. Aucun appel n'est effectué à ces positions.

Dans nos données, les appels SS sont majoritaires et représentent environ 48.1 % de nos appels (contre 35.6 % d'appels DS). Au vu de l'importance de ces appels, nous avons fait le choix de les conserver afin de ne pas filtrer une quantité trop importante de données. Ces appels seront cependant considérés comme étant de faible qualité, de fait, leurs analyses et interprétation seront plus précautionneuses. En revanche, au vu de la trop grande incertitude de l'appel des variants NS, ceux-ci sont systématiquement filtrés éliminant ainsi entre 10.3 et 18.7 % des positions appelées pour chaque patient (**Figure : 4.3 - A**).

Un second filtre est appliqué aux variants ayant été précédemment appelés DS. Celui-ci consiste à comparer les appels effectués indépendamment sur chacune des deux *ends* et à vérifier leur concordance, c'est à dire que les deux appels soit identique. Les appels discordant et ambigus sont filtrés, soit environ 86.3 % des variants DS. Il est intéressant de noter que bien que les variants *single strand* (SS) soient conservés, on peut s'attendre à ce qu'environ 13.7 % de ceux-ci soient aberrants, ceux-ci n'ayant pu subir le même contrôle que les SS (**Figure : 4.3 - B**).

Pour l'ensemble des variants ayant passé les filtres énoncés ci-dessus, c'est à dire les variants SS et les variants DS avec appels concordants, le génotype est déterminé en fonction du pourcentage de *reads* portant le variant à cette position. Par exemple, si à une position donnée, 0% des *reads* portent un variant, l'individu sera appelé "Homozygote référence", si 50% des *reads* sont porteurs d'un variant, l'appel sera "hétérozygote" et si 100% des *reads* portent un variant, l'appel sera "Homozygote variant". Ainsi, pour chaque individu nous avons pu établir une liste de SNVs et d'indels avec leur génotype associé. Pour chacun de nos 13 patients les ordres de

grandeur du nombre de variants appelés sont identique. Ainsi pour chaque patient nous avons appelés environ 43670 variants hétérozygotes (41044 SNVs et 2626 indels) et 65040 variants homozygotes (32520 SNVs et 1809 indels) (**Figure : 4.3 - C**).

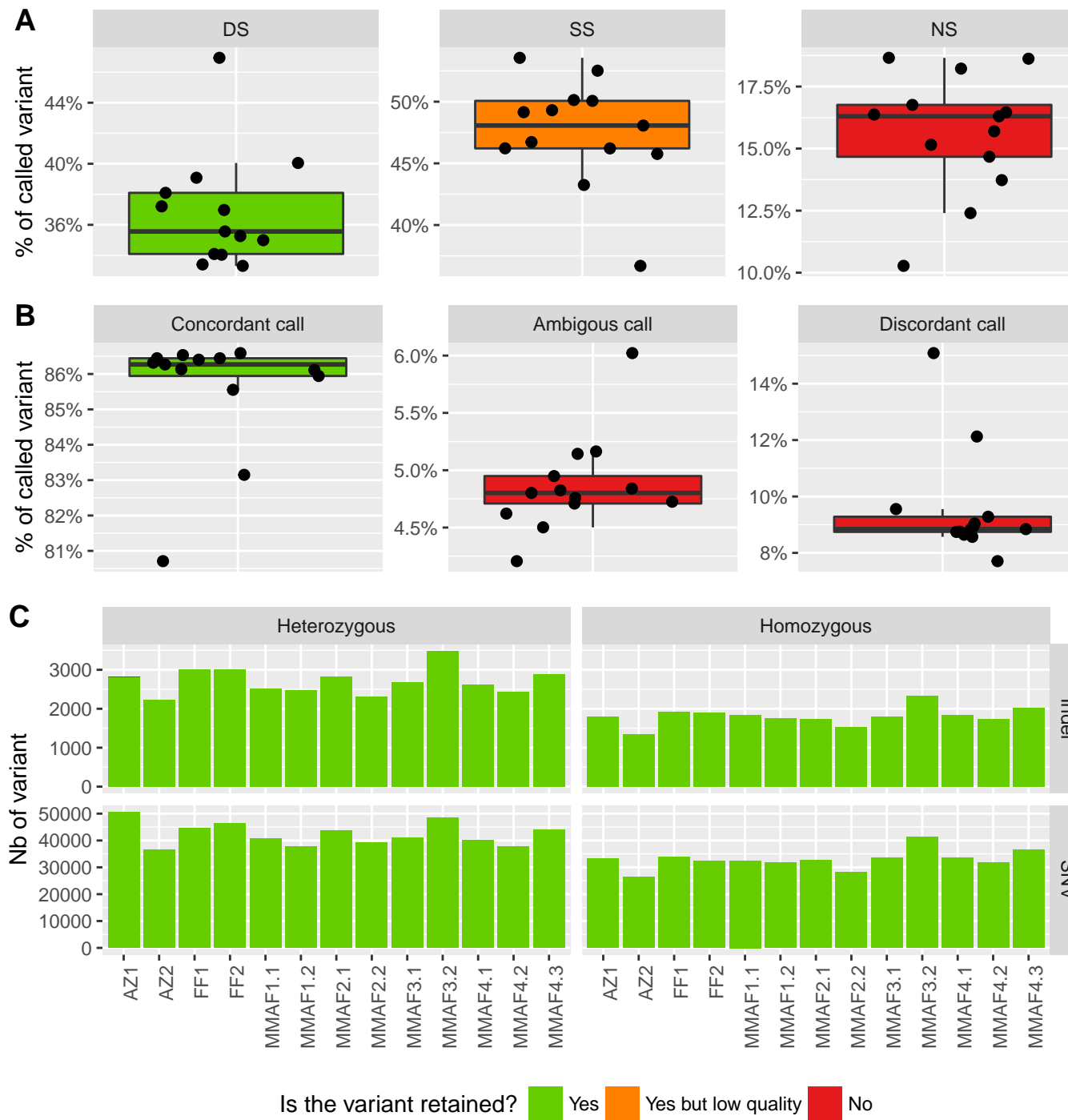


Figure 4.3 – Contrôle qualité des variants appelés : Pour chacun des graphiques, les variants représentés en vert et en orange sont conservés tandis que ceux en rouge sont filtrés. ****A**** : Distribution du *stranding* des appels pour chaque patient. ****B**** : Comparaison des appels entre les deux *ends* des variants appelés DS. ****C**** : Distribution des SNVs et indels en fonction de leur génotype pour chaque patients (représentés par une barre

Résultats de l'annotation L'annotation des variants appelés consiste à ajouter un maximum d'informations sur les variants. Ces informations seront ensuite utilisées afin de filtrer et / ou prioriser les variants. Dans ces analyses nous avons utilisé le logiciel *Variant Effect Predictor* (VEP) (W. McLaren et al., 2016) va prédire l'effet qu'auront ces variants sur l'ensemble des transcrits (et gènes) qu'ils chevauchent. Dans le cas de substitution faux-sens, c'est à dire entraînant le changement d'un seul acide-aminé de la séquence protéique, nous utiliserons les prédictions fournies par SIFT et PolyPhen afin d'estimer leur pathogénicité. Pour finir nous ajoutons, lorsqu'elle est disponible, la fréquence de chacun de ces variants dans les bases de données ExAC, 1000Genomes (1KG) et ESP6500.

Après avoir annoté nos variants, nous avons pu constater que pour chaque patient 24975 gènes sont en moyenne affecté par au moins un variant pour en moyenne 122735 transcrits (soit environ 5 transcrits par gènes). Il faut noter que parmi ces gènes se trouvent à la fois des gènes codant pour des protéine **et** d'autres non codant (**Figure : 4.4 - A**).

Chaque variant affectera l'ensemble des transcrits qu'il chevauche, ainsi un même variant pourra impacter plusieurs transcrits. Ces impacts sont ensuite classés par VEP en quatre catégories qui sont, de la plus délétère à la moins délétère : *HIGH*, *MODERATE*, *LOW*, *MODIFIER* (**Table :4.1**). Comme attendu, les variants ayant un impact tronquant se retrouvent être les moins fréquent chez chacun de nos patients. Ceci est d'autant plus flagrant pour l'impact *HIGH* qui regroupe, entre autres, les variants créant un codon stop ou encore ceux causant un décalage du cadre de lecture (**Table :4.1**), se retrouvent en quantité extrêmement faible puisqu'ils ne représentent en moyenne que 0.15 % des variants, soit une moyenne de 466 hétérozygotes et 370 homozygotes par patient) (**Figure : 4.4 - B**).

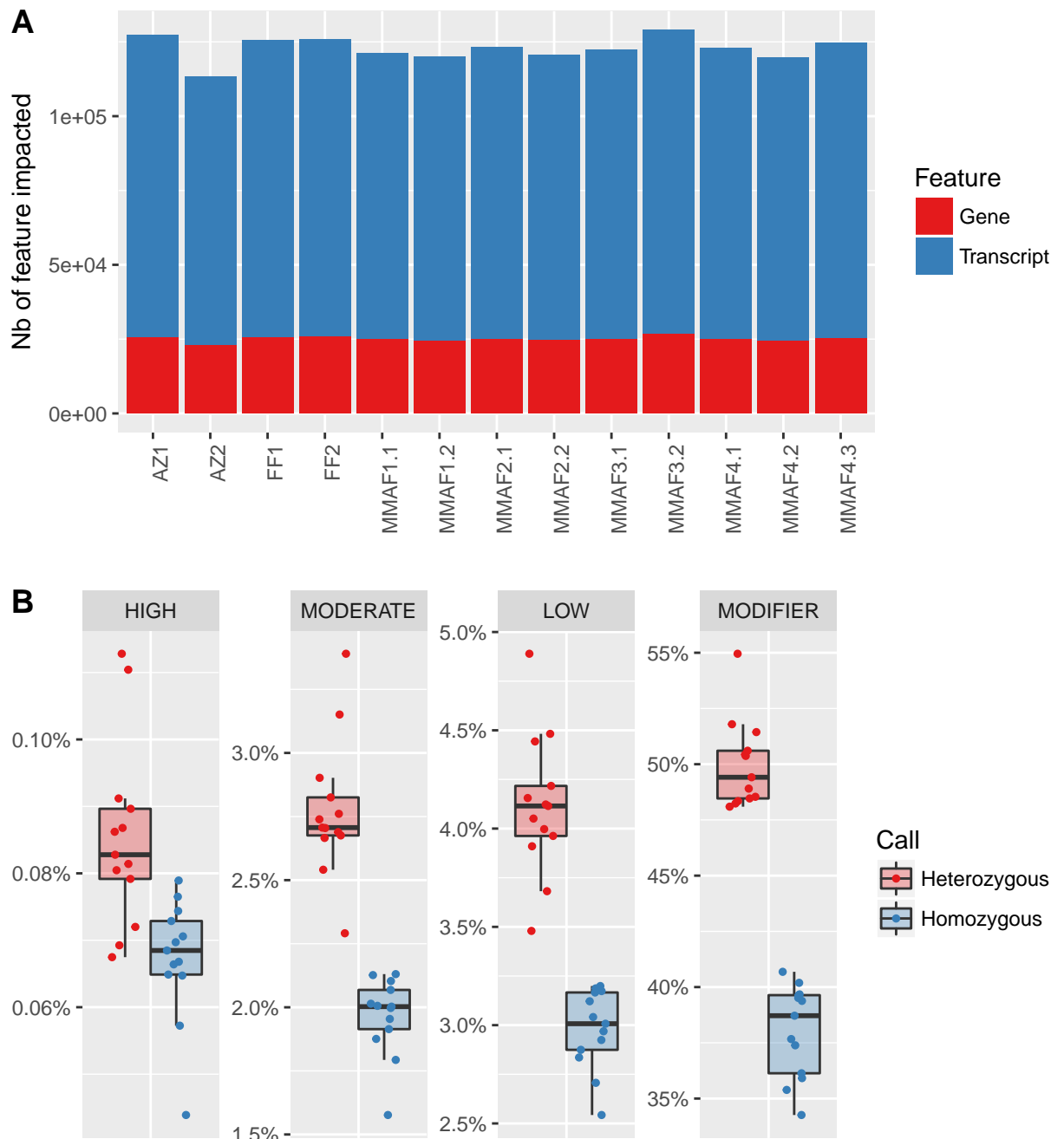


Figure 4.4 – Annotation des variants : ****A**** : Quantification du nombre de gènes (en bleu) / transcrits (en rose) impactés par au moins un variant pour chaque patient chacun représentés par une barre. ****B**** : Distribution des impacts HIGH MODERATE LOW et MODIFIER en fonction des patients et du génotype du variant

Résultats du filtrage Les étapes précédentes nous ont permis de mettre en évidence pour chaque patient une liste de variants passant l'ensemble de nos critères qualités. Ces variants ont dès lors pu être annotés nous permettant notamment d'avoir connaissance de leurs impacts sur les différents transcrits qu'ils chevauchent ou encore leur fréquence dans la population générale. Désormais, afin de ne conserver que les variants ayant la plus forte probabilité d'être responsable du phénotype de ces patients, nous avons appliqué successivement six filtres basés à la fois sur les différentes annotations que nous avons ajoutées mais aussi sur nos connaissances du mode de transmission du phénotype :

1. **Filtre 1 : L'union des variants** : Dans ces différentes études, nous avons à chaque fois séquencé des frères (deux ou trois) présentant phénotype. Ainsi nous avons pu formuler l'hypothèse d'une cause génétique commune entre les différents patients d'une même famille et donc filtrer l'ensemble des variants qui ne sont pas partagés les deux ou trois frères atteints testés.
2. **Filtre 2 : Génotype des variants** : Dans ces études, nous avons émis l'hypothèse d'une transmission récessive du phénotype. Ainsi, seuls les variants homozygotes ont été conservés. (**Figure : ??, ??**).
3. **Filtre 3 : Impact du variant** : Afin de ne conserver que les variants ayant un effet potentiellement délétère sur la protéine, nous avons filtré les variants intronique et ceux tombant dans les séquences UTRs. De même les variants synonymes ne sont pas conservés (exceptés ceux se trouvant proches des régions d'épissage) car ceux-ci n'ont aucun effet sur la séquence protéique. Pour les variants faux sens (changement d'un seul aa de la séquence protéique) il est plus difficile de se trancher, nous avons donc utilisé les logiciels SIFT (Kumar, Henikoff, & Ng, 2009) et Polyphen (Adzhubei et al., 2010) et filtré l'ensemble des faux-sens prédits comme *tolerated* par SIFT et *benign* par Polyphen.
4. **Filtre 4 : Les transcrits "non pertinents"** : Au cours de nos analyses nous nous sommes concentré uniquement sur les transcrits codant pour une protéine. Ainsi, l'ensemble des transcrits annotés comme étant non codant furent filtrés. De même le mécanisme NMD (*nonsense-mediated decay*) a pour but de contrôler la qualité des ARNm cellulaires chez les eucaryotes (Y.-F. Chang, Imam, & Wilkinson, 2007) en éliminant les ARNm qui comportent un codon stop prématuré (K. E. Baker & Parker, 2004), pouvant être le résultat d'une erreur de transcription, d'une mutation ou encore d'une erreur d'épissage. Il est donc peu probable que les variants présents sur des transcrits annotés NMD soient responsables du phénotype. Dès lors, ces transcrits ont été également filtrés. Ainsi, nous avons pu retirer de nos listes de variants l'ensemble des mutations impactant **uniquement** des transcrits non codant et / ou annoté NMD. Cette étape de filtre permet à elle seule de filtrer systématiquement entre 13712 et 17992 transcrits différents par patients, soit une moyenne de 1834 variants par individus (**Figure : 4.5**).

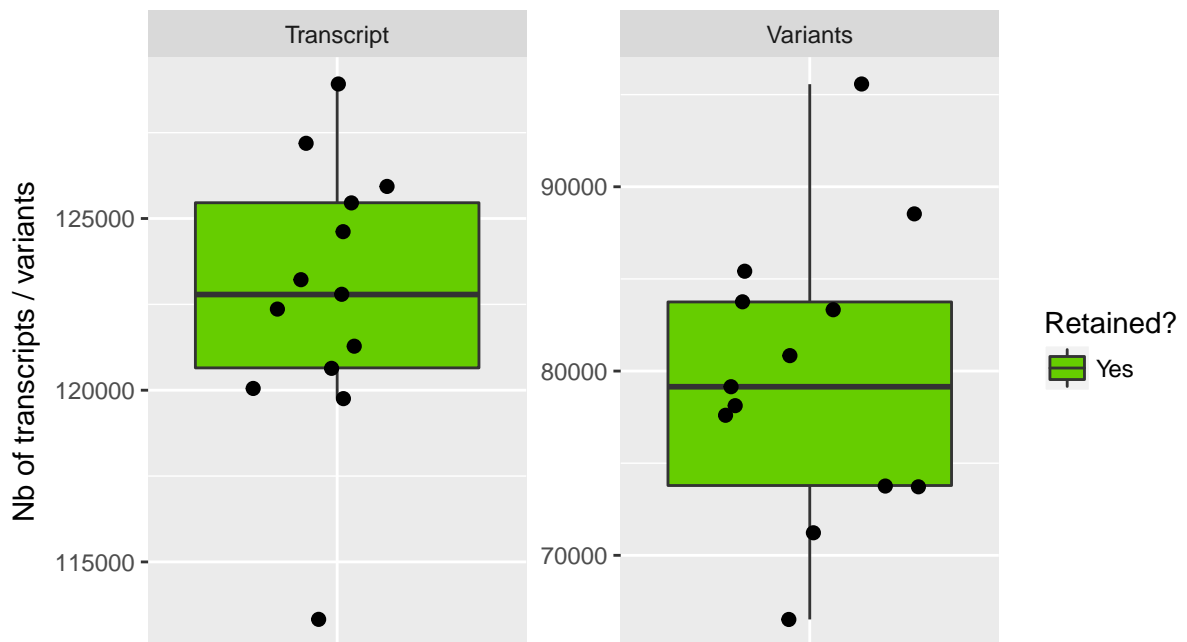


Figure 4.5 – Filtrage des transcrits jugés "non pertinents" et des variants les chevauchant : Pour chaque patients nous avons filtrer les transcrits jugés "non pertinents" pour l'analyse, c'est à dire ceux ne codant pas pour une protéine et ceux annoté NMD. Dès lors, l'intégralité des variants chevauchant uniquement des transcrits non pertinents ont put systématiquement être filtrés (boites rouges). les autres furent conservés (boites vertes)

5. **Fréquence des variants** : La fréquence d'un variant dans la population générale est un moyen rapide d'avoir une prédiction fiable de l'effet délétère ou non de celui-ci. En effet, il est peu probable qu'un variant retrouvé fréquemment dans la population générale soit causal d'une pathologie sévère. Ainsi nous avons filtré pour l'ensemble de nos patients l'ensemble des variants ayant une fréquence ≥ 0.01 dans l'une des trois bases de données que sont ExAC, ESP et 1KG.
6. **Présence des variants dans la cohorte contrôle** : Au cours de nos différentes études, nous avons été amenés à séquencer un total de 134 individus présentant un des 6 phénotypes que nous avons étudié (**Table : ??**). Ces phénotypes étant très différents, on peut émettre l'hypothèse que leurs causes génétiques soient également différentes. De même, les variants recherchés étant rares, il est peu probable qu'un individu porte les variants de deux phénotypes différents. Ainsi, pour chacune des 6 familles, nous avons pu constituer une cohorte contrôle composée dans l'ensemble des patients précédemment analysés et ne présentant pas le même phénotype que celui étudié dans la famille (**Figure : 4.6**). Dès lors, nous avons pu filtrer l'ensemble des variants retrouvés à la fois chez nos patients et observés à l'état homozygote dans la cohorte contrôle. Cette cohorte contrôle présente ainsi le même rôle que les bases de données publiques. Sont intérêt

principale par rapport à celles-ci est que les individus qui la composent ont pour la plupart la même origine ethnico-géographique que nos patients. De plus ceux-ci ont été séquencés en même temps dans les mêmes centres permettant ainsi d'identifier les artefacts dûs aux protocoles de séquençage.

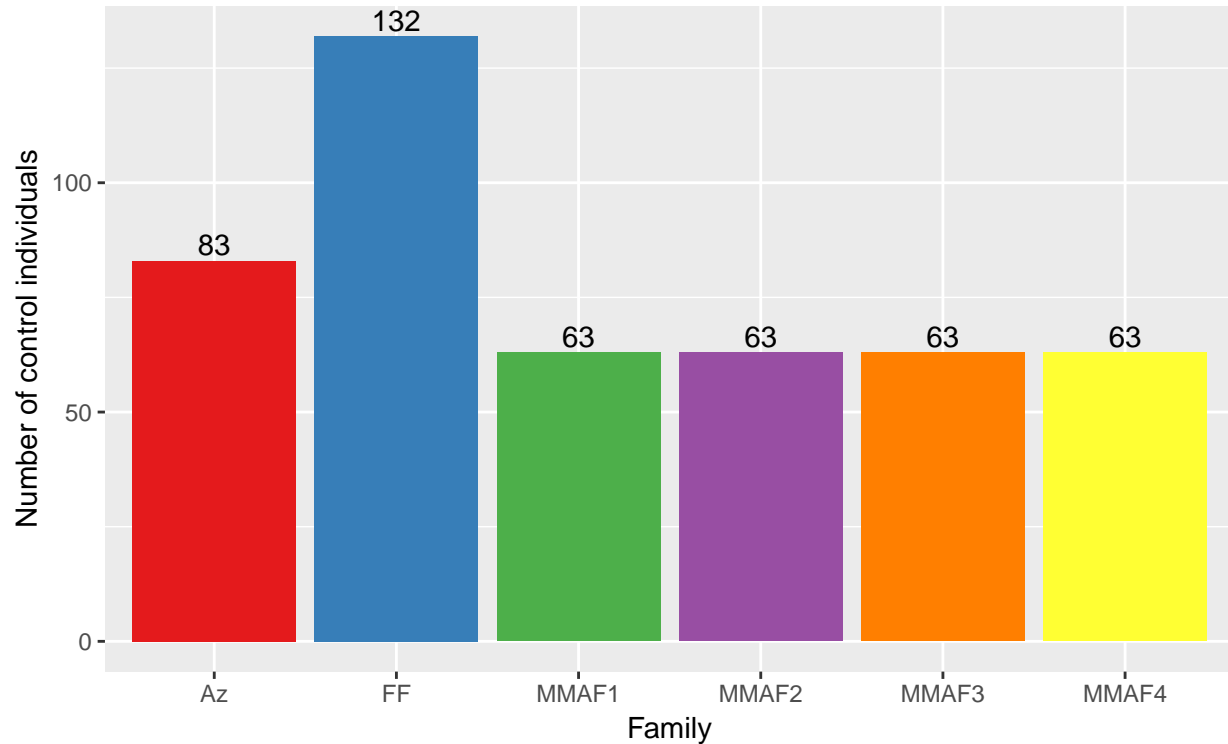


Figure 4.6 – Nombre d'individus composant la cohorte contrôle de chaque famille : Ici, chaque barre représente une famille et sa hauteur est déterminée par le nombre d'individus composant la cohorte contrôle à laquelle elle a été confrontée. Chaque individu de la cohorte contrôle a été séquencés en WES par notre équipe. Afin d'être considéré comme "contrôle" et intégrer cette cohorte, un individu doit être sain ou présenter un phénotype d'infertilité différent de la famille étudiée. Par exemple, un individus MMAF pourra servir de contrôle aux familles AZ et FF mais pas aux familles MMAF1-4

Comme on pouvait s’y attendre, ces six filtres ont un pouvoir discriminant extrêmement différent. En effet, tandis que le filtre “Transcript relevance” (filtre n°4) éliminer en moyenne 3.9 % des variants de chaque individu, le filtre “Variant impact” (filtre n° 3) élimine jusqu’à 90.1 % de ces mêmes variants. Cette différence n’est pas surprenante. En effet, comme nous l’avions vu plus tôt, les variants de la catégorie VEP *MODIFIER* qui regroupe entre autres les variants chevauchant les séquences UTRs et introniques (**Table : 4.1**) représentent en moyenne 88% des variants de nos patients. Ceux-ci étant tous filtrés, on s’attendait donc à une valeur aussi élevée. On peut également constater l’importance de la cohorte contrôle qui, je le rappelle, permet de filtrer l’ensemble des variants homozygotes observés en son sein, puisque ce filtre permet retirer entre 76.5 et 88.4% des variants de chaque individus (**Figure : 4.7 - A**).

Cependant, regarder uniquement le pourcentage de variants filtrés par chaque filtre révèle une information partielle. En effet, dans ce cas de figure, on observe la quantité de variant éliminé par chaque filtre indépendamment les uns des autres. Ainsi, un même variant peut donc être filtré par plusieurs filtres. Dès lors, il faut également analyser la quantité de variants filtrés **spécifiquement** par chaque filtre. Ainsi, on peut constater que le classement des filtres en fonctions de leur stringance reste quasiment identique. Il est tout de même intéressant de noter que désormais le filtre “Variant impact” apparaît moins efficace que les filtres “Ctrl” et “Genotype” en filtrant spécifiquement une moyenne de 253 variants par individu contre 423 pour le filtre génotype et 882 pour le filtre “Ctrl”. Ainsi, ce dernier devient celui filtrant spécifiquement le plus de variants avec entre 364 et 1060 variants spécifiquement filtrés par patients confirmant ainsi l’importance de ce filtre dans nos analyses. Aussi, les filtres “Transcript relevance”, “Union” et “Frequency” apparaissent désormais comme étant anecdotiques en comparaison aux trois autres filtres puisqu’ils filtrent au maximum 43 variants spécifiques (**Figure : 4.7 - B**).

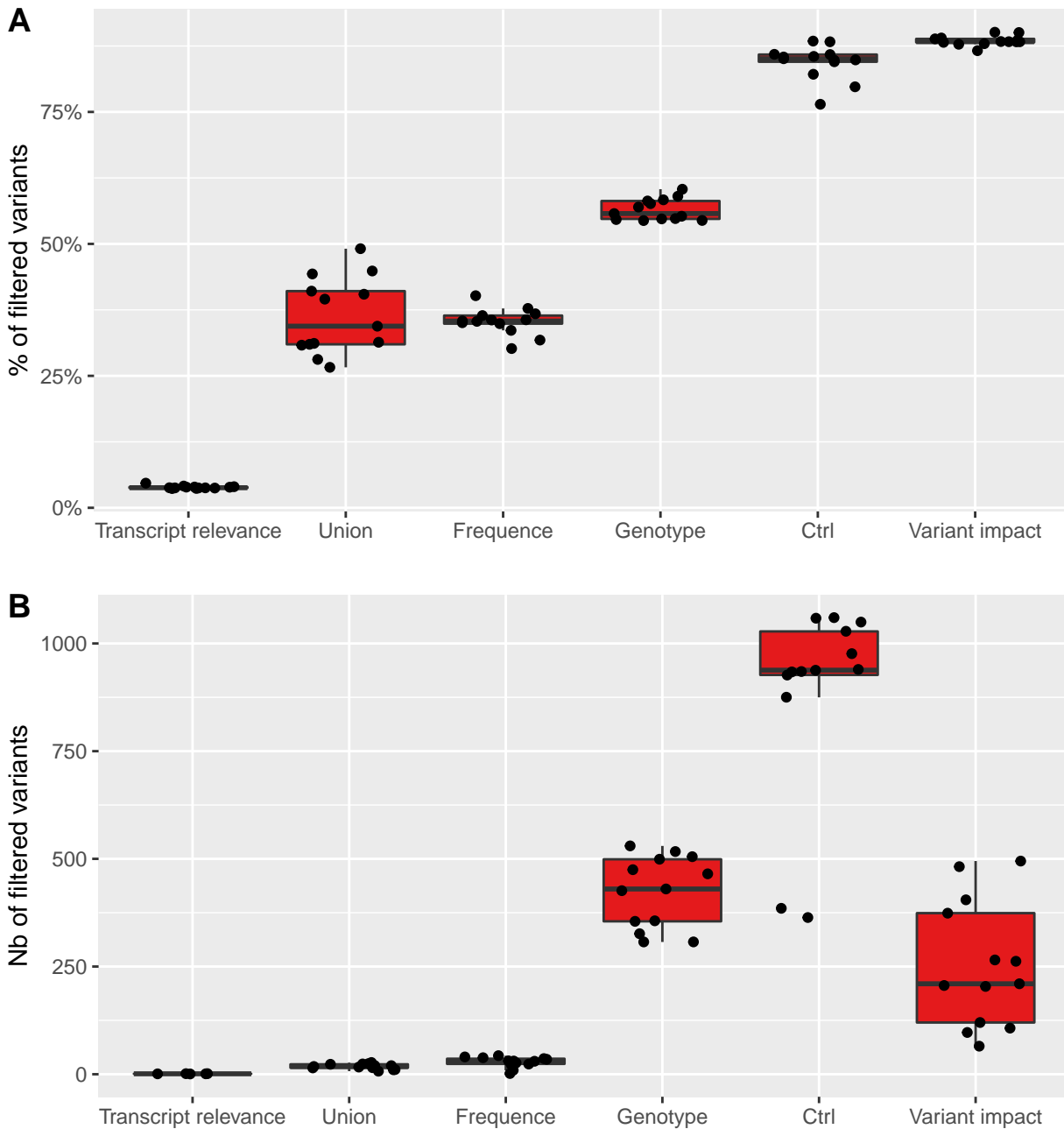


Figure 4.7 – Comparaison de l'efficacité de chacun des six filtres utilisés : ****A**** : Comparaison du pourcentage de variants filtrés par chacun des six filtres indépendamment les uns des autres pour chaque patient (représenté par les points. Dès lors, un même variant peut être filtré par plusieurs filtres. ****B**** : Comparaison du nombre de variant filtré spécifiquement par chacun des filtres. Ici, un variant ne peut-être filtré que par un seul filtre

Après avoir appliqué l'ensemble de ces filtres, seuls quelques variants subsistent nous permettant d'obtenir une liste de gènes restreinte pour chaque famille et ainsi de tirer des conclusions quant au variant responsable du phénotype de chacune d'entre elles (**Table : 4.3**).

Table 4.3 – Récapitulatif des variants ayant passé l'ensemble des filtres pour chaque famille

Family	Gene	HGVSc, HGVSp	Variant impact			Variant frequency		
			Consequence	SIFT	PolyPhen	ExAC	ESP	1KG
AZ	GUF1	c.443A>T; p.Ser148Ile	missense	deleterious	proba damaging	0.00207	0.0028	9e-04
AZ	SPINK2	c.56-3C>G; .	splice region
FF	PLCZ1	c.1465G>T; p.Ile489Phe	missense	deleterious	possib damaging	8.24e-06	.	.
MMAF1	PLA2G4B	c.1710-6delA; .	splice region
MMAF2	DNAH1	.; .	splice acceptor
MMAF2	MYH11	c.4625G>A; p.Arg1542Gln	missense	.	proba damaging	0.00234	0.0016	5e-04
MMAF3	WEE2	.; p.Pro92Leu	missense	deleterious	benign	0.000372	2e-04	.
MMAF3	ZFYVE28	c.1729C>A; p.Val577Met	missense	deleterious	benign	0.000998	2e-04	.
MMAF3	FCGR3A	c.133T>C; p.Ala45Pro	missense	deleterious	proba damaging	.	.	.
MMAF3	GBP2	c.412T>A; p.Ala138Thr	missense	deleterious	proba damaging	0.00176	0.0012	5e-04
MMAF4	TGIF2	c.496C>A; p.Leu166Met	missense	tolerated	proba damaging	.	.	.

1. **Famille AZ** : Parmi les 2 gènes restant pour cette famille, *SPINK2* est apparu comme étant un candidat évident. Notamment son expression étant spécifique au testicule tandis que celle de *GUF1* est ubiquitaire (**Figure : 4.8**). De plus, des mutations du gène *Spink2* chez la souris avait déjà été identifiée comme induisant des défauts de la spermatogenèse (B. Lee et al., 2011).

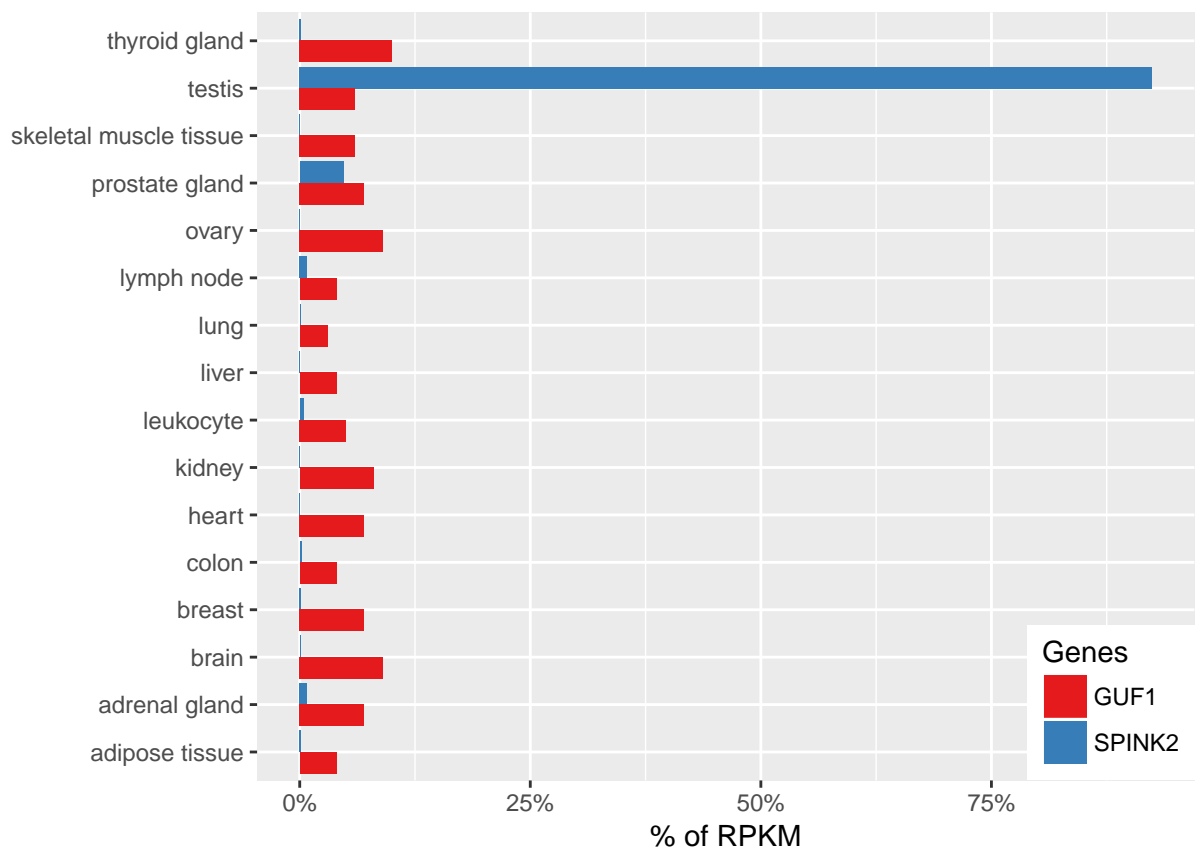


Figure 4.8 – Expression tissulaire des gènes *SPINK2* et *GUF1* :
Données provenant du projet de transcriptome Illumina bodyMap

2. **Famille FF** : Pour cette famille, seul le gène *PLCζ1* a passé l'ensemble des filtres. Nos connaissances sur la fonction de ce gène et notamment son rôle dans l'activation ovocytaire (Amdani, Jones, & Coward, 2013) ainsi que sa forte expression testiculaire ont fait de ce gène le candidat idéal pour expliquer le phénotype d'échec de fécondation de ces deux frères (**Figure : 4.9**).

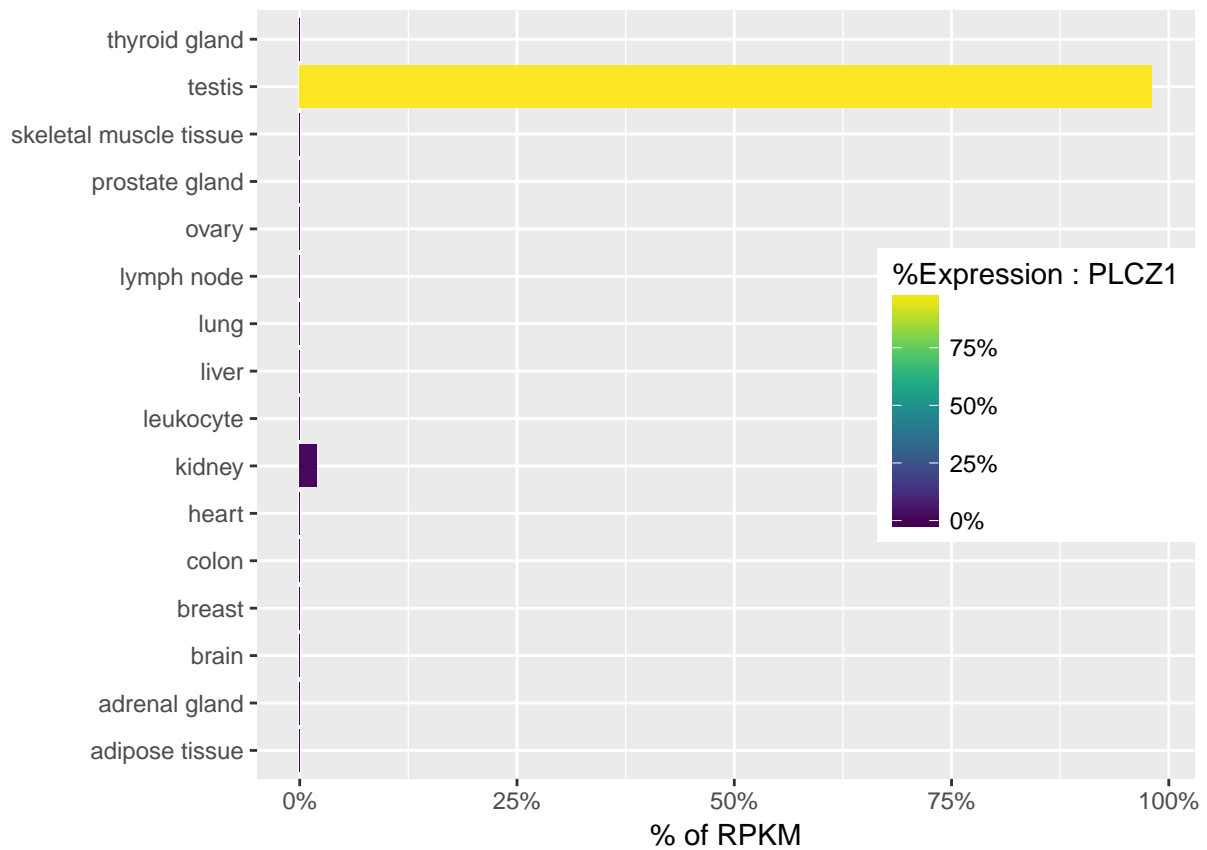


Figure 4.9 – Expression tissulaire du gène *PLCZ1* : D'après les données du Illumina BodyMap

3. **Famille MMAF1** : L'analyse bibliographique des 1 gènes ayant passé l'ensemble des filtres n'a pas pu nous permettre de d'affirmer que l'un de ces gènes étaient responsable du phénotype MMAF de ces 2 frères.
4. **Famille MMAF2** : À l'issue des filtres, 2 gènes ressortaient chez ces deux frères : *MYH11* et *DNAH1*. Or, notre équipe ayant déjà établi le lien entre des mutations du gène *DNAH1* et le syndrome MMAF (Ben Khelifa et al., 2014) ce gène s'est révélé être un candidat idéal pour expliquer le phénotype de ces 2 frères. De plus, l'implication de *MYH11* dans le phénotype de dissection aortique (Imai et al., 2015) l'ont écarté des candidats pour le phénotype MMAF.
5. **Famille MMAF3** : Comme pour les gènes de la famille MMAF2, l'analyse bibliographique des 4 gènes ayant ici passé les filtres de même que l'étude de leurs expressions ne nous a pas permis de conclure que l'un d'entre eux étaient responsable du phénotype MMAF de ces 2 frères.

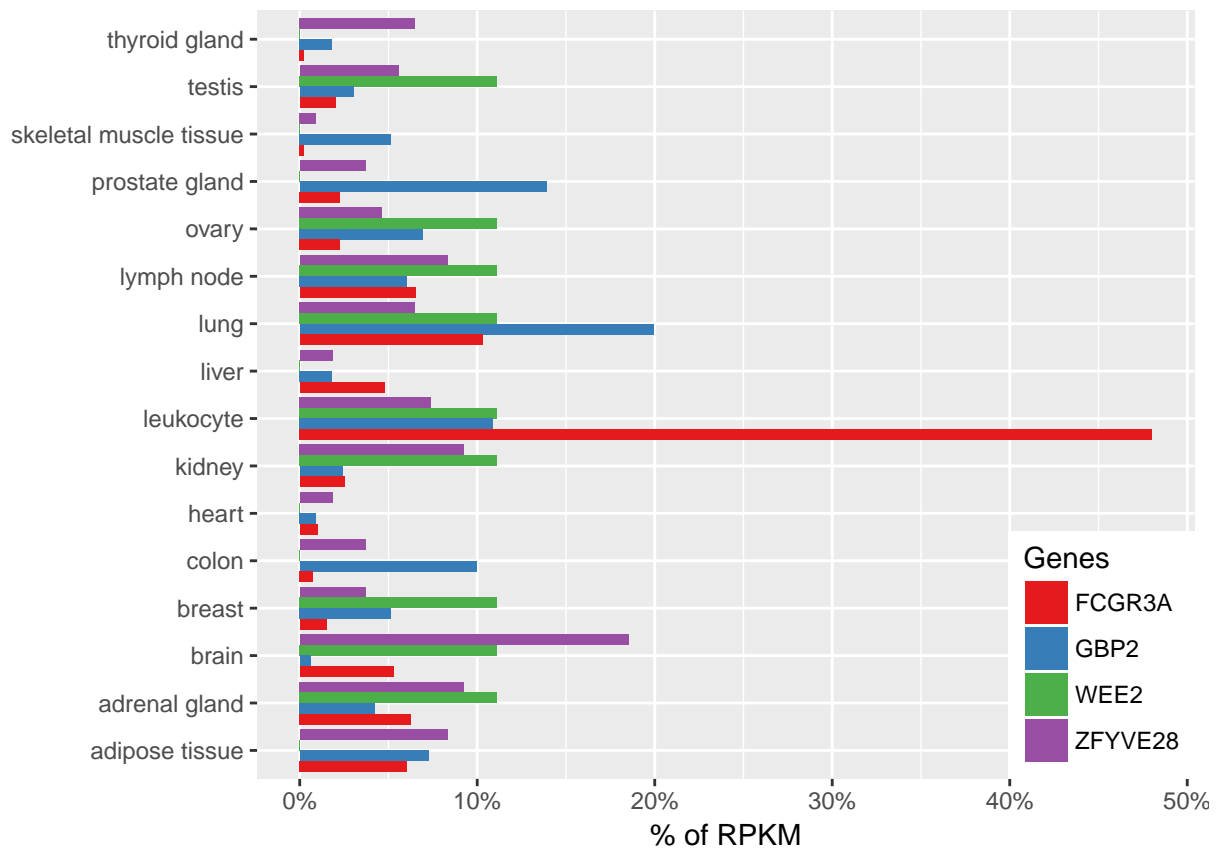


Figure 4.10 – Expression tissulaire des gènes retenus pour la famille MMAF3 : Données provenant du projet de transcriptome Illumina bodyMap

6. **Famille MMAF4** : Seul le gène *TGIF2* a passé l'ensemble des filtres pour la famille MMAF4. L'expression ubiquitaire de ce gène n'en font pas un candidat idéal. Cependant une étude de 2011 effectuée sur le wallaby décrit que la protéine TGIF2 est localisée spécifiquement dans le cytoplasme du spermatide, ainsi que dans le corps résiduel et la pièce intermédiaire du flagelle du spermatozoïde mature (Hu, Yu, Shaw, Renfree, & Pask, 2011). Ces données pourraient corrélérer avec le phénotype MMAF de ces 3 frères bien que l'expression de ce gène soit ubiquitaire (**Figure** : 4.11).

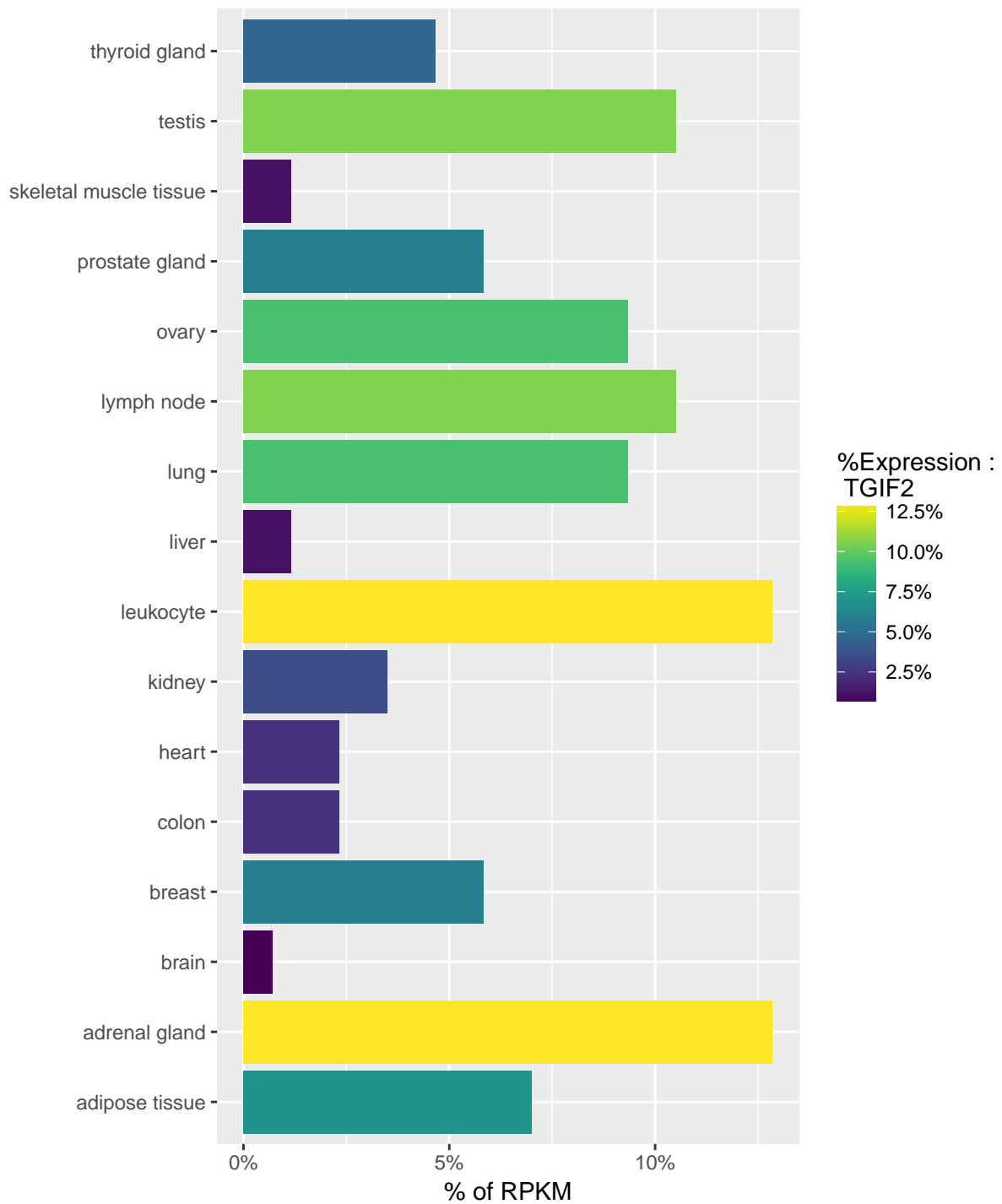


Figure 4.11 – Expression tissulaire du gène **TGIF2** : D'après les données du Illumina BodyMap

Discussion

L'analyse de ces 6 familles nous a permis de mettre en évidence l'efficacité de notre pipeline d'analyse puisque pour 3 d'entre elles (soit 50%) le variant causal a pu être identifié avec certitude (**Figure : 4.12**) et les résultats publiés dans trois revus dont je suis co-auteur :

1. **Famille AZ : SPINK2 deficiency causes infertility by inducing sperm defects in heterozygotes and azoospermia in homozygotes** : Dans cet article j'ai effectué non seulement l'intégralité des analyses bioinformatiques des données d'exomes de deux frères infertiles présentant un phénotype d'azoospermie mais j'ai aussi séquencé en Sanger les séquences codantes du gène *SPINK2* pour une partie des 611 individus analysés ainsi que contribué à l'extraction de l'ARN testiculaire des souris pour l'analyse fonctionnelle du gène *Spink2* sur le modèle murin.
2. **Famille FF : Homozygous mutation of PLCZ1 leads to defective human oocyte activation and infertility that is not rescued by the WW-binding protein PAWP** : Dans cet article j'ai, effectué l'intégralité des analyses bioinformatiques des données d'exomes effectuées sur deux frères infertiles présentant des échecs de fécondation.
3. **Famille MMAF2 : Whole-exome sequencing of familial cases of multiple morphological abnormalities of the sperm flagella (MMAF) reveals new DNAH1 mutations** : Dans cet article j'ai, comme précédemment, effectué l'ensemble des analyses bioinformatiques des données d'exomes effectuées sur deux frères infertiles présentant des échecs de fécondation.

Pour une d'entre elle, un candidat potentiel a pu être mis en évidence avec le gène *TGIF2* et notre équipe travaille actuellement sur la caractérisation de ce gène afin de savoir s'il peut effectivement expliquer le phénotype MMAF de cette famille (**Figure : 4.12**).

TODO : Il faut aller plus loin dans l'analyse et les arguments pour convaincre qu'il s'agit d'un bon candidat : quel type de mutation, ce gène est-il bien conservé, son expression n'est pas spécifique au testicule et ce gène serait impliqué dans un phénotype d'holoprocéphaly...

Pour les 2 familles restantes, aucun variant n'a pu pour l'instant expliquer leur phénotype. L'explication la plus vraisemblable est que le variant ait été filtré par l'un de nos six filtres, probablement celui consistant à filtrer l'ensemble des variants hétérozygotes. En effet, l'hypothèse d'un variant causal homozygote était extrêmement crédible pour les familles AZ, FF et MMAF2 étant donné l'historique consanguin de ces 3 familles dont les parents sont à chaque fois apparentés. En revanche rien ne laisse supposer une telle chose pour les familles restantes. Cependant, le filtre des variants hétérozygotes pour l'ensemble des patients de ces 3 familles a été maintenu en

première intention afin de faciliter les analyses en réduisant au maximum le nombre de variant. Au vu des résultats il apparaît clair que les variants responsables de leur phénotype aient été filtrés pour au moins 2 de ces familles. Dès lors, l'ensemble des analyses effectuées lors de l'étape de filtrage doivent être refaites en changeant les paramètres de filtrage. Cette fois-ci, les variants hétérozygotes seront conservés et les gènes sur lesquels au moins deux variants hétérozygotes seront recensés seront analysés en priorité. En effet, bien que les analyses exomiques nous fournissent en l'état pas d'informations suffisante pour savoir si ces deux variants sont présent sur le même allèle ou bien sur deux allèles différents, cela pourrait-être la signature de variants hétérozygotes composites. C'est donc sur ces analyses que se concentre actuellement notre équipe.

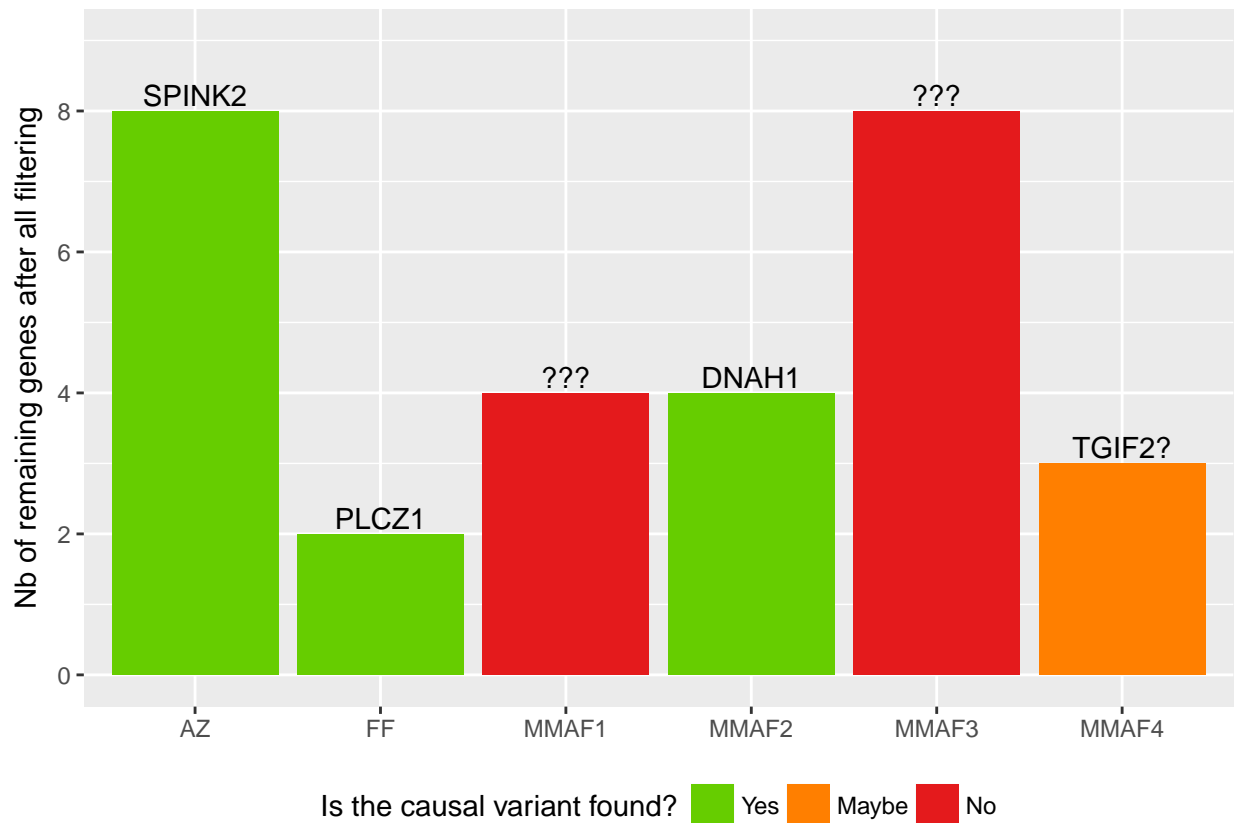


Figure 4.12 – Nombre de gènes passant l'ensemble des filtres par famille : Chaque barre représente une des familles analysées. La hauteur de cette barre correspond au nombre de gènes ayant passé l'ensemble des filtres pour chaque famille. Les barres vertes caractérisent les familles pour lesquelles le gène responsable de la pathologie a été identifié parmi la liste de gène (dans ce cas le symbole du gène est écrit au-dessus de la barre). La barre orange caractérise la famille pour laquelle un candidat potentiel a été identifié (le symbole du gène est écrit au-dessus suivi d'un "?"). Les barres rouges indiquent qu'aucun des gènes ayant passé les filtres pour ne semble expliquer le phénotype (dans ce cas il est écrit "???" au-dessus de la barre)

4.2.3 Etude d’une large cohorte de patients MMAF

Description de la cohorte

Après avoir mis en évidence l’implication du gène *DNAH1* dans le phénotype MMAF notre équipe s’est en partie spécialisé dans la caractérisation ce syndrome. Ainsi, entre 2012 et 2015, notre équipe a effectué le séquençage de 62 individus présentant tous ce phénotype afin d’en établir la cause génétique. ces séquençage ont été effectué dans 3 centres différents que sont Genoscope, MountSinai et Strasbourg et sur une seule plateforme plateforme de séquençage, le Illumina HiSeq2000.

Table 4.4 – Liste des différents projets de séquençages effectués

Place	Year	Platform	Nb of individuals
MountSinai	2012	Illumina Hiseq2000	2
Strasbourg	2012	Illumina Hiseq2000	13
Genoscope	2013	Illumina Hiseq2000	13
Genoscope	2014	Illumina Hiseq2000	31
Genoscope	2015	Illumina Hiseq2000	6

Application de la pipeline - Résultats

Après avoir appelé les variants de nos 62 patients, nous avons obtenu un total de 4484558 variants différents comprenant 4160274 SNVs et 324284 indels. Ces variants étant réparti entre chaque patient qui portaient environs chacun 81618 SNV et 5148 indels dont 42.8 % étaient homozygote. Comme on peut le voir, la proportion de chaque appel est relativement homogène lorsque l'on compare les patients ayant été séquencés dans le même centre la même année. Cependant, il est possible de noter de grandes disparités lorsque l'on compare les données provenant de différents centres ou bien du même centre avec plusieurs années de différences. Ces écarts peuvent-être causés par plusieurs facteur, tel que les différents kits de capture d'exons qui on put être utilisés puisque ... (**todo lister les différents kits de capture dans une table**) en revanche nous pouvons écarter un effet dus à la plateforme de séquençage ou encore le modèle de séquenceur puisque tous ces projets ont été réalisés sur des Illumina HiSeq2000 (**Table : 4.4**) (**Figure : 4.13 - A**).

Le même constat peut être effectué lorsque l'on compare la qualité des appels puisque plus les projets de séquençage s'avèrent être récent, plus la proportion d'appel *Single Strand* s'avère être faible tandis que la proportion d'appel *Double Strand* (DS) est élevée. Ceci est une bonne chose, car, bien que ces deux appels soient conservés dans les analyses ultérieures, les appels DS sont de meilleure qualité que les appels SS. Cette augmentation des appels DS au cours du temps pourrait s'expliquer par une amélioration des protocoles de séquençage ainsi que des kits de capture. En revanche cela est à pondérer avec le taux croissant d'appels *No-strand* (NS) au fur et à mesure des années pour atteindre environs 21.3 % en 2015 avec un projet réalisé au Génoscope. Ces derniers appels étant systématiquement filtrés, ils n'altéreront en rien les résultats obtenus en aval hormis le fait qu'ils réduisent la quantité des données utilisées (**Figure : 4.13 - B et C**).

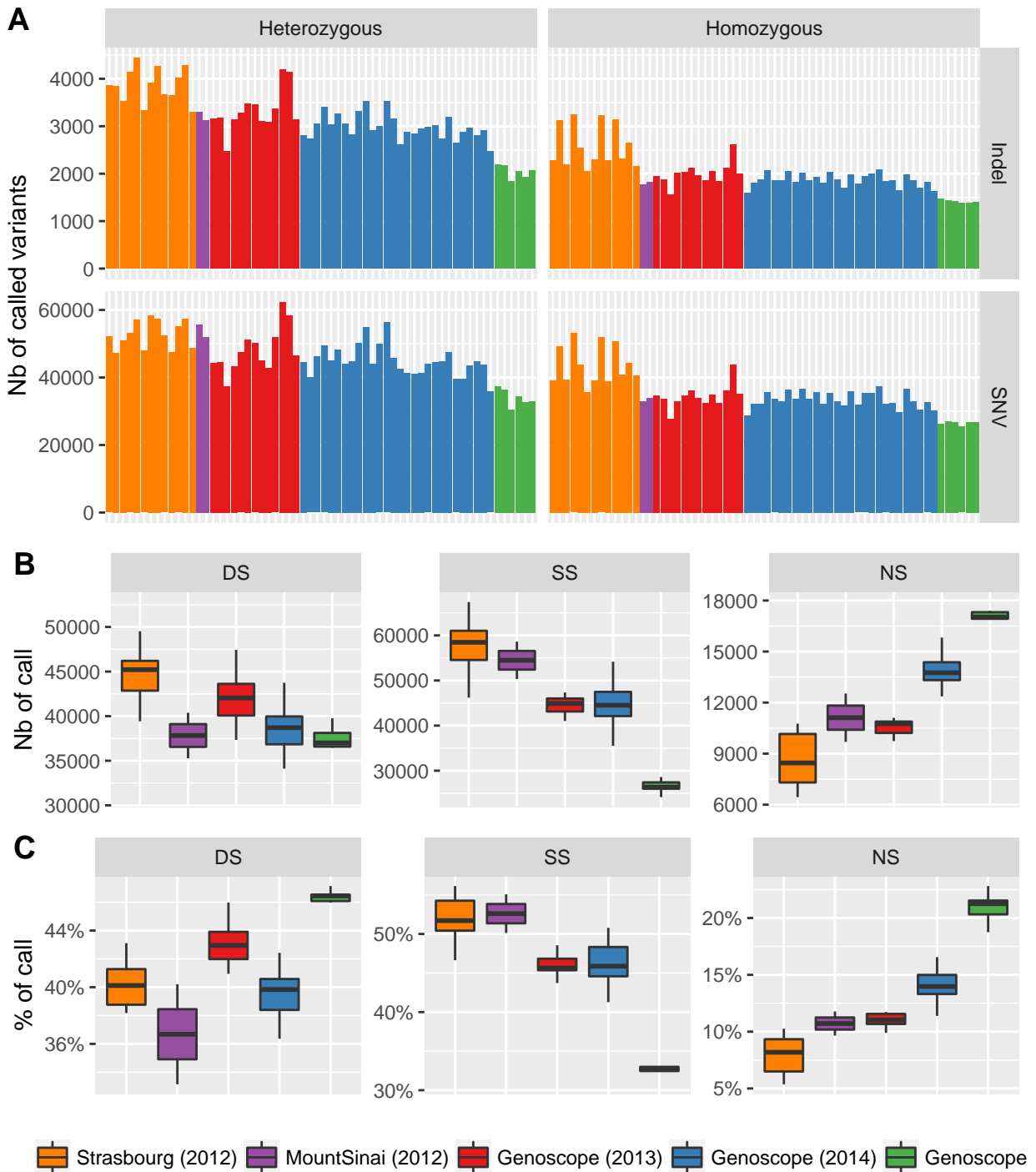


Figure 4.13 – Résultats de l'appel des variants par individus et par projet de séquençage : Chaque couleur définit un projet de séquençage caractérisé par un centre de séquençage et une année. ****A**** : Quantification pour chaque individus (représentés par les barres) du nombre de variants (SNVs et Indels) appelés homozygotes et hétérozygotes. ****B**** : Quantification des appels **Double Strand** (DS), **Single Strand** (SS) et **No strand** (NS) pour chaque projet de séquençage. ****C**** : Même chose en pourcentage

Analyse des listes de gènes

Après avoir appliqué les mêmes filtres que ceux décrit précédemment à l'exception du filtre n°1 "Union" puisqu'ici nous avons uniquement des individus non apparentés, nous avons pu obtenir une liste de 1711 variants différents composés de 1470 SNVs et 241 indels et impactant un total de 1432 gènes distincts. Ces variants étant répartis sur l'ensemble de nos 62 patients ceux-ci portaient en moyenne 27 SNVS et 5 indels, de sorte que chacun d'entre eux avaient entre 1 et 86 gènes impactés par au moins un variants homozygote (**Figure : 4.14 - A et B**).

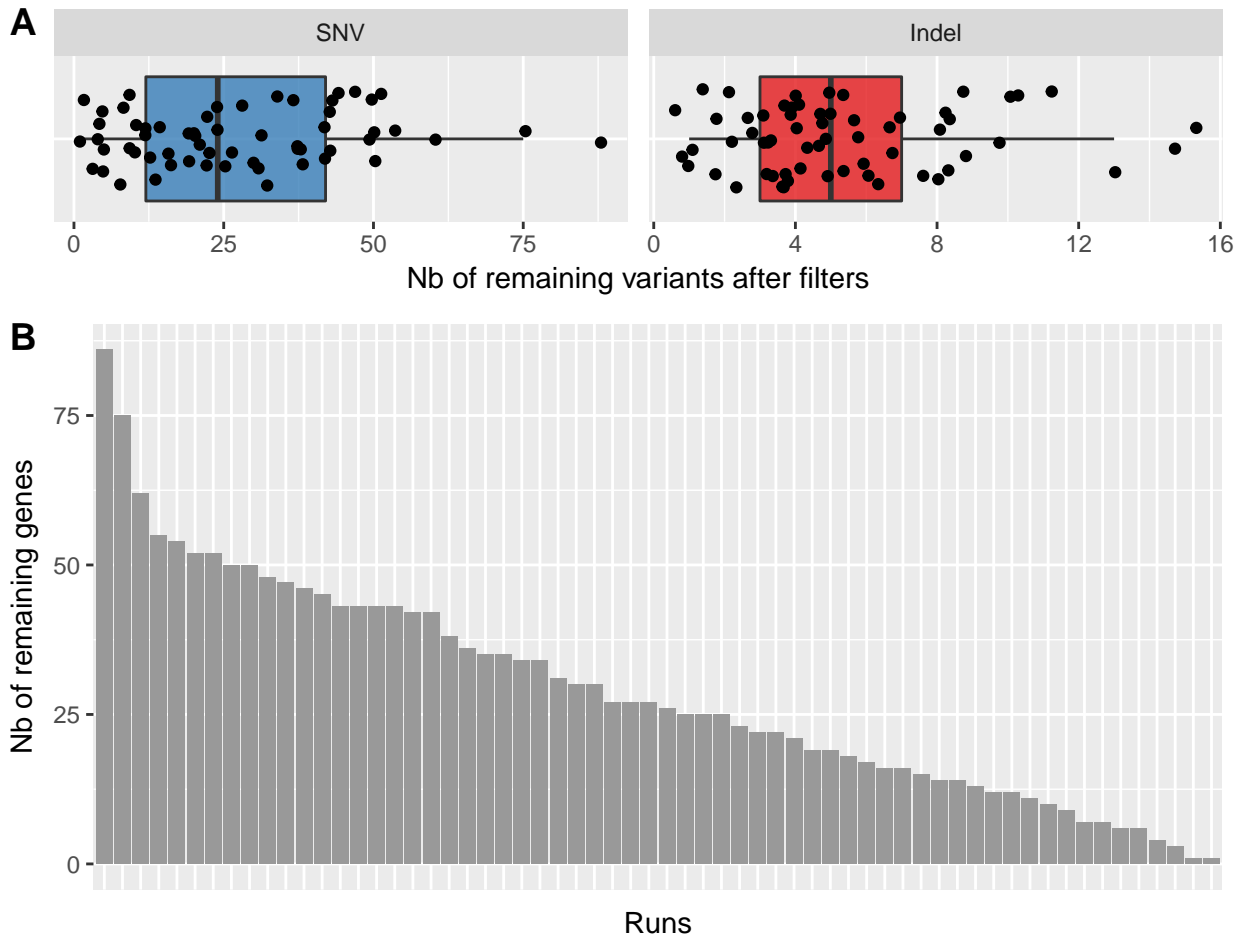


Figure 4.14 – Résultats de l'étape de filtrage : ****A**** : Quantification du nombre de SNVs et indels ayant passé l'ensemble des filtres pour chaque patient. ****B**** : Nombre de gènes impactés par au moins un variant ayant passé les filtres pour chaque individu représenté par les barres. ****C**** : Présentation

Afin de nous orienter dans nos recherches, nous nous sommes basés sur une étude de 2012 qui établissait une liste des gènes humains pouvant être impliqués dans cillio-me, c'est à dire (todo def cillio-me) (Ivliev, 't Hoen, Roon-Mom, Peters, & Sergeeva, 2012). La constitution de cette liste se basait à la fois sur les données de CilDB (Arnaiz, Cohen, Tassin, & Koll, 2015) et de MEDLINE [ref?] mais aussi des analyse *in silico* permettant d'effectuer des prédiction. Ainsi, chaque gène était classé dans l'une des 3 catégories suivantes en fonction des preuves déjà existante (au moment de l'étude) permettant de lier un gène au cillio-me humain : **Strong evidence from previous studies** (Strong), **Weak evidence from previous studies** (Weak) et **No evidence from previous studies** (Novel). L'utilisation de cette liste nous a permis d'ajouter une nouvelle annotation pertinente à nos gènes. En effet, le spermatozoïde humain est une cellule ciliée, et le flagelle en est le cil. Nous pouvons donc attendre à ce qu'une partie des gènes responsables du phénotype MMAF soit présents dans cette liste de 371 gènes.

Ainsi, 33 de nos 1432 gènes retenus faisaient partis de cette liste dont 22 présentaient des preuves fortes de leur appartenance au cillio-me. Il faut tout de même noter que bien que cette liste soit un bon outil pour orienter les recherches et prioriser certains gènes, elle ne peut constituer un critère suffisant pour filtrer les gènes n'en faisant pas partie. Par exemple le gène *DNAH1*, de par son expression ubiquitaire n'a pas été intégré à cette liste (**Figure : 4.17 - A**), or on connaît désormais son implication dans le phénotype MMAF (**Table : 4.5, Figure : 4.15 - A**).

Parmis l'ensemble des gènes retrouvés mutés, on peut constater que ... d'entre eux, soit ...%, ne sont observés muté, à l'état homozygote, uniquement chez un seul de nos patients tandis que ... sont retrouvés muté à l'état homozygote chez au moins 2 patients. Ainsi, analyser en priorité ces ... gènes permet à la fois de réduire considérablement notre liste, et donc de faciliter les analyses, mais aussi de pouvoir identifier en priorité les gènes affectant le plus grand nombre de patient et donc probablement les principaux acteurs du phénotype MMAF (**Figure : 4.15 - B**).

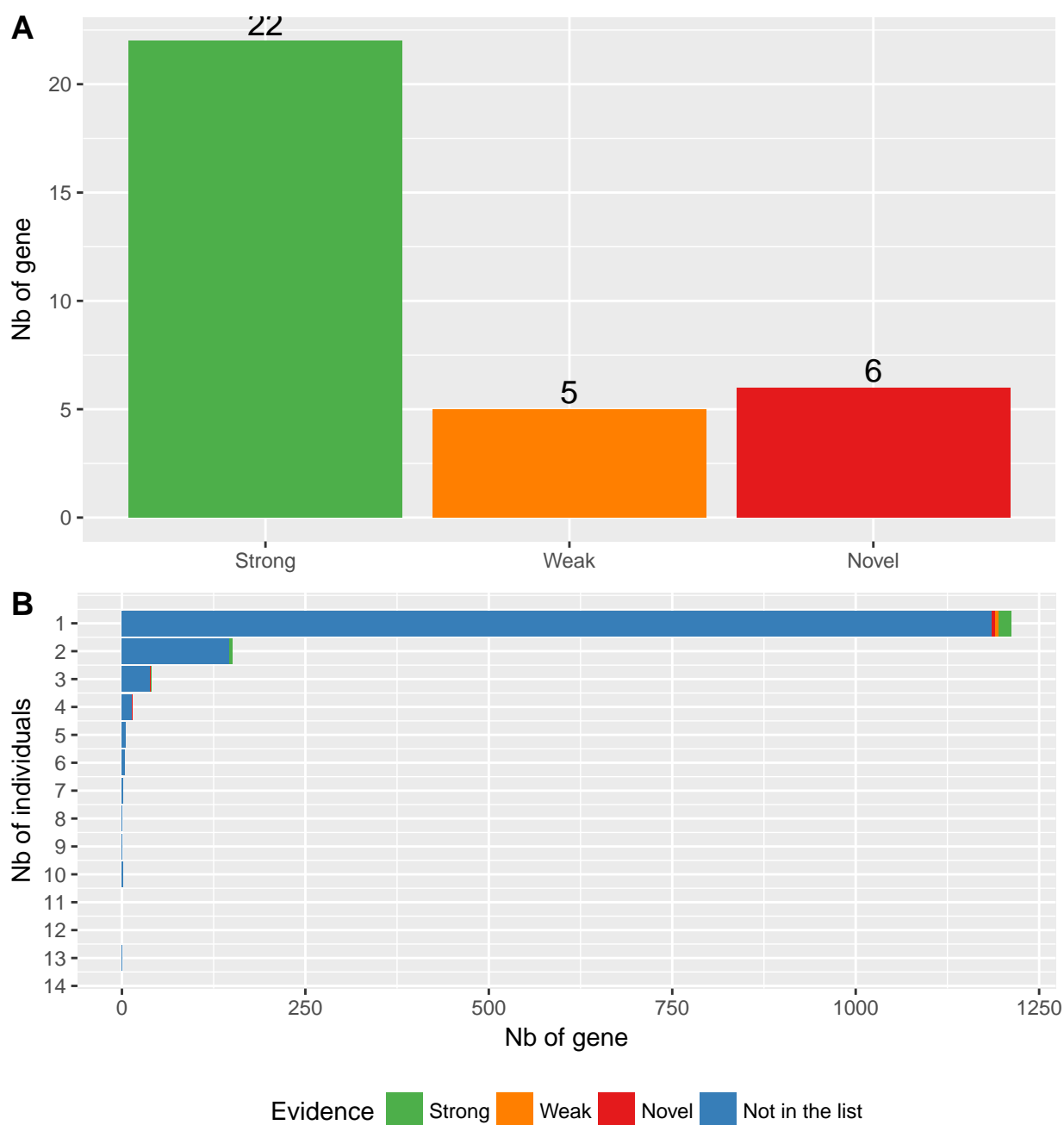


Figure 4.15 – Répartition des gènes retrouvés mutés chez nos patients dans les différentes classes de la liste du cillioime : Chaque couleur définit une classe de la liste des gènes du cillioime décrit dans [Ivliev2012]. Vert = *Strong evidence from previous studies* (*Strong*), Orange = *Weak evidence from previous studies* (*Weak*), rouge = *No evidence from previous studies* (*Novel*), bleu = Non présent dans la liste. ****A**** : Quantification du nombre de gène ayant passé les filtres au sein des 3 classes de la liste des gènes du cillioime. ****B**** : Quantification pour chacun des gènes du nombre d'individus présentant un variant homozygotes ayant passé l'ensemble des filtres

Table 4.5 – Liste des gènes prédits comme faisant partie du ciliome humain retrouvés mutés à l'état homozygote chez nos patients

Ciliome evidence	Gene
Strong	AK7
Strong	C21orf59
Strong	C6orf118
Strong	CCDC146
Strong	CCDC147
Strong	CCDC65
Strong	CFAP44
Strong	DLEC1
Strong	EFCAB6
Strong	FAM81B
Strong	HYDIN
Strong	KIF6
Strong	KIF9
Strong	NPHP1
Strong	PROM1
Strong	RSPH9
Strong	SPEF2
Strong	STK33
Strong	TTC26
Strong	TTC29
Strong	WDR16
Strong	ZMYND10
Weak	C6
Weak	FBXO15
Weak	MIPEP
Weak	SLFN13
Weak	SPAG17
Novel	ACYP1
Novel	ARMC2
Novel	C21orf58
Novel	KIAA0319
Novel	KIAA0556
Novel	LRRC43

Nous avons ensuite *designer* 5 analyses fonctionnant de manière pyramidales. C'est à dire que la première analyse permet d'identifier les candidats candidats le plus évidents tandis que la cinquième ceux pour lesquels il y a moins d'indices. L'intérêt de ce fonctionnement est qu'à chacune des analyses, l'ensemble des variants portés par les patients pour lesquels un candidat a put être déterminer avec certitudes seront retirés dans analyses suivantes. Ce procédés nous permettra ainsi après chacune de ces 5 analyses d'alléger notre liste de variants, et donc de gènes, rendant l'analyse suivante plus simple :

1. **Analyse n°1** : Au moment de nos analyses, le gène *DNAH1* était encore le seul décrit comme responsable du phénotype MMAF faisant de lui un candidat évident pour expliquer le phénotype MMAF de nos patients malgré son expression non spécifique au testicule (**Figure : 4.17 - A**). C'est pourquoi nous avons rechercher en priorité des patients portant des variants tronquants homozygotes sur ce gène.
2. **Analyse n°2** : Dans second temps, nous avons sélectionné uniquement les gènes **présents dans la liste cillio** sur lesquels **au moins deux** de nos patients présentaient au moins 1 variant tronquant à l'état **homozygote**.
3. **Analyse n°3** : Ensuite, nous avons étudié les gènes **absents dans la liste cillio** mais sur lesquels on trouvait toujours **au moins deux** de nos patients présentant au moins 1 variant tronquant à l'état **homozygote**.
4. **Analyse n°4** : Dans un troisième temps nous sommes revenus à étudier les gènes **présents dans la liste cillio** en considérant cette fois-ci les gènes sur lesquels **un seul** de nos patients présentaient au moins 1 variant tronquant à l'état **homozygote**.
5. **Analyse n°5** : Pour finir nous avons étudié les gènes **absents dans la liste cillio** sur lesquels **un seul** de nos patients présentaient au moins 1 variant tronquant à l'état **homozygote**.

Chacune de ces analyses s'est ensuite déroulée en 5 étapes :

1. **Étape n°1** : Cette première étape consistait à récupérer l'ensemble des patients répondant aux critères de l'analyse en question, par exemple : l'ensemble des patients portant un variant homozygote sur le gène *DNAH1* dans le cas de l'analyse n°1. Au vu de l'effet délétère des variants sélectionnés au cours de cette étape, les preuves génétiques seront considéré comme de forte confiance (*High trust*).
2. **Étape n°2** : Pour l'ensemble des gènes retenus dans l'étape n°1, nous recherchons ensuite des patients portant, **toujours à l'état homozygote**, des variant aux effets **non tronquants** tel que des variants faux-sens ou encore des variants intronique situés proches des sites d'épissage. Dans le cas des variants faux-sens, les logiciels SIFT et PolyPhen sont ensuite utilisés afin de nous orienter quant à l'effet délétère du variant, bien que comme nous l'avons déjà vu, ces logiciels son contredisent régulièrement [TODO : ref!!!]. Au vus de la difficulté à déterminer

l'effet délétère de ces variants, les preuves génétiques seront considérées comme de confiance modérée (*Moderate trust*).

3. **Étape n°3** : Cette étape consiste à rechercher des patients éventuellement hétérozygotes composites, c'est à dire des patients portant deux variants hétérozygotes différents sur chacun des deux allèles d'un même gène. Malheureusement, dans le cadre des séquençages WES et WGS, il est impossible de connaître le "phasage" des variants, c'est à dire que l'on ne peut déterminer si deux variants hétérozygotes sont situés sur le même allèle ou sur deux allèles différents (**Figure** : 4.16). Pour cela, des analyses de biologie moléculaire sont nécessaires. C'est pour cette raison que ces preuves génétiques seront considérées comme étant de faible confiance (*Low trust*).

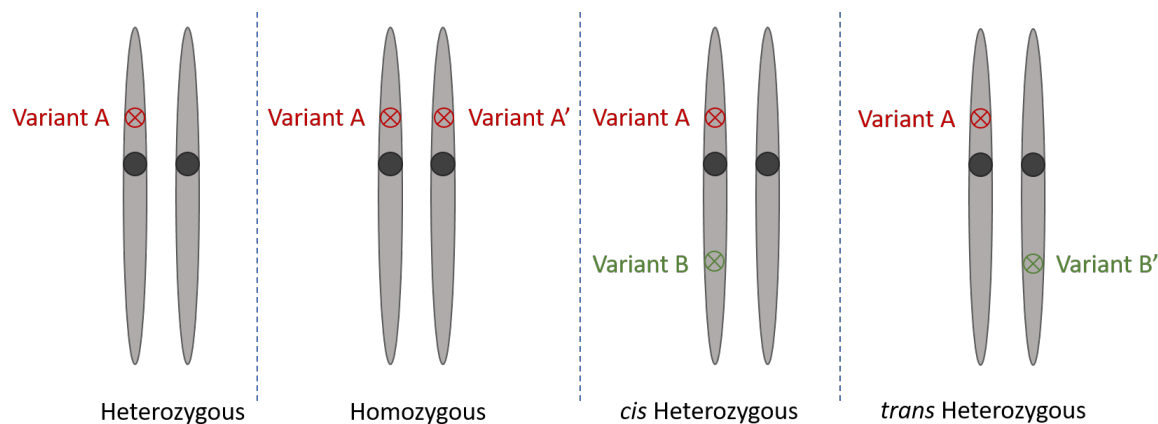


Figure 4.16 – Filtrage des transcrits jugés "non pertinents" et des variants les chevauchant : Un variant est dit homozygote lorsque le **même** variant est présent sur les deux allèles d'un gène et hétérozygote lorsqu'il est présent sur **un seul** des deux allèles. On parlera d'hétérozygotes **cis** lorsque deux variants hétérozygotes différents seront positionnés sur **le même** allèle et d'hétérozygote **trans** (ou composite) lorsque ces deux hétérozygotes seront positionnés sur **deux allèles différents**. En WES et en WGS il est impossible de différencier les hétérozygotes **cis** des hétérozygotes **trans**.

4. **Étape n°4** : Au cours de cette étape, nous allons étudier chacun des gènes sélectionnés et déterminer au vu des différentes preuves génétiques, des données d'expressions de ces gènes mais aussi des informations disponibles dans la littérature, si ceux-ci sont de bons candidats pour expliquer le phénotype MMAF de nos patients. Il est à noter que **cette étape est la seule à ne pas être automatisée**.
5. **Étape n°5** : Pour finir, les données des patients pour lesquels un gène candidat a été identifié et dont les preuves génétiques sont considérées comme de forte confiance (c'est à dire les patients identifiés lors de l'étape n°1) sont retirées de notre liste de variants allégeant ainsi celle-ci pour les analyses suivantes.

Analyse n°1 :

1. ***DNAH1*** : Parmi l'ensemble de nos 62 patients MMAF 1 portait un indel homozygote entrainant un décalage du cadre de lecture sur le gène *DNAH1* et répertorié dans aucune des 3 bases de données que nous avons utilisées. Le patient et *Ghs90* était porteur de 3 variants homozygotes successifs entrainant tous un faux-sens tandis que le patient et *Ghs95* portait lui 1 autre variant homozygote entrainant un faux-sens différents des 3 portés par et *Ghs90*. On peut noter qu'aucun de ces 4 variants n'est répertorié dans les bases de données, cependant tous sont prédits comme *benign* par PolyPhen tandis que SIFT ne donne aucune prédiction. La recherche de patients potentiellement hétérozygotes composites nous a permis de révéler 6 patients portant tous 2 variants hétérozygotes sur le gène *DNAH1*. On peut alors noter le patient *Ghs36* porteur d'un premier variant créant un codon stop de manière prématuré et d'un second variant entrainant un faux-sens absent des 3 bases de données et prédit comme *probably damaging* par PolyPhen. Ainsi, s'ils étaient situés sur deux allèles différents, ces deux variants pourraient être des bons candidats pour expliquer le phénotype du patient *Ghs36*. On peut noter également le patient *Ghs129* portant deux faux-sens hétérozygote tout deux prédits comme *probably damaging* par PolyPhen et dont un seul est retrouvé dans la base de données ExAC avec une très faible fréquence. Pour les 4 autres patients, il est plus difficile de se prononcer à ce stade compte tenu des fréquences parfois élevée des variants ou bien des prédictions fournies par PolyPhen.

Cette première analyse nous a permis de révéler que 9 des 62 patients de notre cohorte portaient au moins 1 variant sur le gène *DNAH1* et que pour 3 d'entre eux ce(s) variant(s) étaient présents à l'état homozygote. Cependant, il faut noter que du fait de son effet tronquant sur la protéine, seul le variant homozygote porté par le patient et *Ghs122* nous permet d'être certains de la causalité du phénotype MMAF. C'est pourquoi, seuls les données de ce patient seront retirées de notre liste de variants réduisant celle-ci à `rn_distinct_var2` variants chevauchant 1395 gènes (**Table : ??**). Pour les 8 autres patients, des analyses fonctionnelles complémentaires sont nécessaires.

Ainsi, les mutations du gène *DNAH1* seraient ainsi responsables de manière certaine de 2 % du phénotype MMAF de notre cohorte. Ce pourcentage monte jusqu'à 15 % si l'on considère l'ensemble des patients identifiés dans cette analyse.

Bien que ce pourcentage soit en deçà des 40% (TODO : à confirmer !) observés dans notre étude précédente (Ben Khelifa et al., 2014), ces résultats confirment néanmoins le rôle primordial de la protéine DNAH1 dans la structure du flagelle et l'implication majeure du gène *DNAH1* dans le phénotype MMAF.

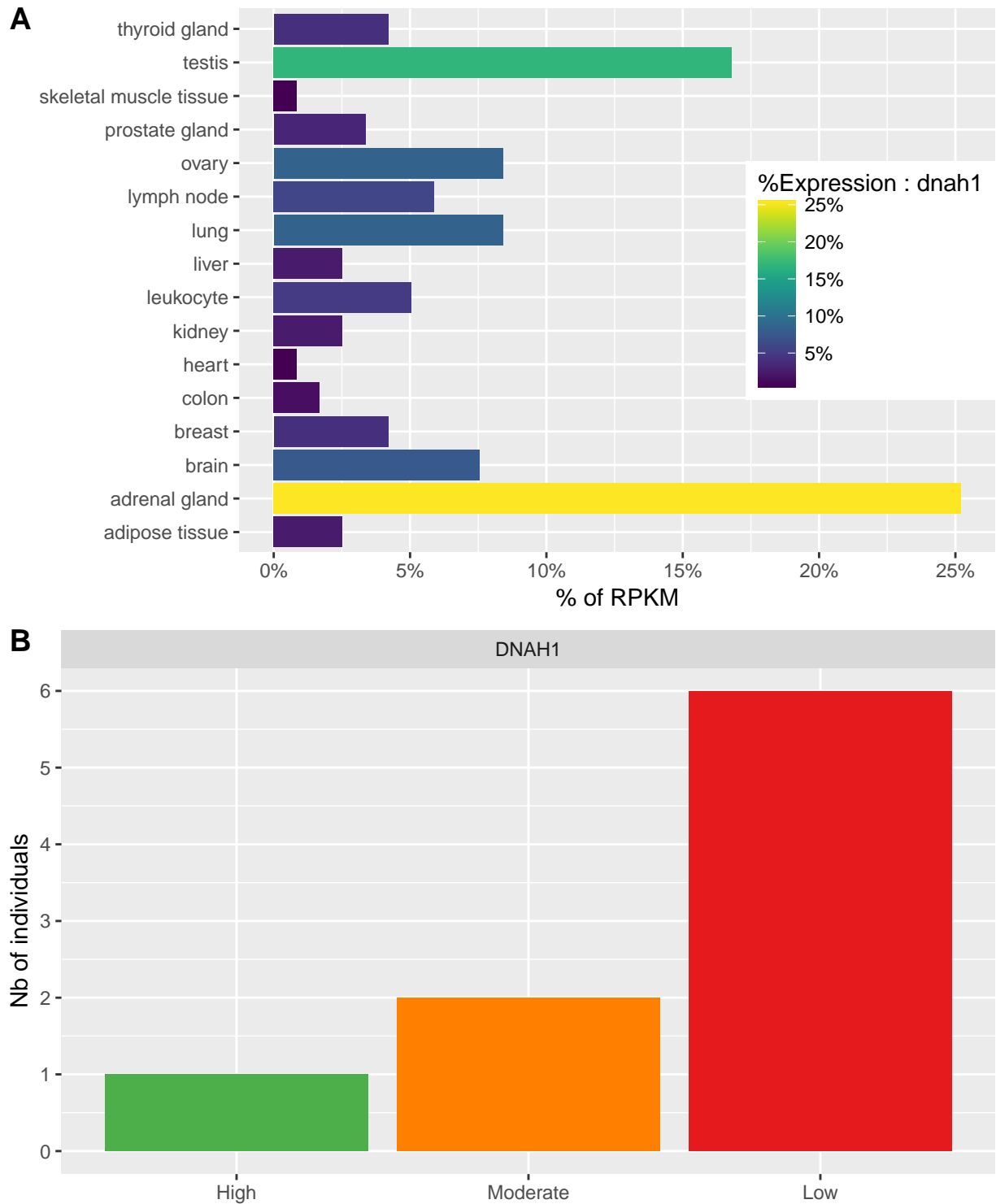


Figure 4.17 – Analyse du gène **DNAH1** : Expression tissulaire du gène **DNAH1** d'après les données du projet Illumina BodyMap. Quantification du nombre de patients portant au moins un variant sur le gène **DNAH1** pour chacun des 3 niveau de confiance

Table 4.6 – Analyse n°1 : Liste des patients portant au moins un variant homozygote tronquant sur le gène *DNAH1*

Patient	Gene	Evidence	Variant impact		Variant frequency		
			HGVSc, HGVSp	Consequence	ESP	1KG	ExAC
Ghs122	DNAH1	Not in the list	c.7533delC ; p.Gln2511HisfsTer27	frameshift	.	.	.

Table 4.7 – Analyse n°1 : Liste des patients portant au moins un variant homozygote non tronquant sur le gène *DNAH1*

Patient	Gene	Variant impact				Variant frequency			
		HGVSc, HGVSp	Consequence	SIFT	PolyPhen	ESP	1KG	ExAC	
Ghs90	DNAH1	c.2122G>C; p.Ile708Leu	missense	.	benign	.	.	.	
Ghs90	DNAH1	c.2123A>C; p.Ile708Thr	missense	.	benign	.	.	.	
Ghs90	DNAH1	c.2125A>C; p.Phe709Leu	missense	.	benign	.	.	.	
Ghs95	DNAH1	c.9278C>G; p.Ala3093Gly	missense	.	benign	.	.	.	

Table 4.8 – Analyse n°1 : Liste des patients portant au moins deux variant hétérozygotes sur le gène *DNAH1*

Patient	Gene	HGVS _c , HGVSp	Variant impact				Variant frequency			
			Consequence	SIFT	PolyPhen	ESP	1KG	ExAC		
Ghs28	DNAH1	c.752C>G ; p.Glu251Gly	missense	.	benign	1e-04	.	1.65e-05		
Ghs28	DNAH1	c.1172A>G ; p.Tyr391Cys	missense	.	benign	0.0027	0.0019	0.00233		
Ghs36	DNAH1	c.845A>G ; p.Leu282Trp	missense	.	proba damaging	.	.	.		
Ghs36	DNAH1	c.3931G>T ; p.Gln1311Ter	stop gained	8.29e-06		
Ghs42	DNAH1	c.6916G>A ; p.Ala2306Thr	missense	.	proba damaging	.	.	.		
Ghs42	DNAH1	.;.	splice region		
Ghs87	DNAH1	c.3431C>A ; p.Ser1144Asn	missense	.	benign	.	.	0.00024		
Ghs87	DNAH1	c.8885G>C ; p.Lys2962Thr	missense	.	proba damaging	7e-04	5e-04	0.000457		
Ghs88	DNAH1	c.2209C>A ; p.Val737Met	missense	.	benign	1e-04	.	0.000115		
Ghs88	DNAH1	c.3877T>A ; p.Asp1293Asn	missense	.	benign	1e-04	0.0019	0.000149		
Ghs129	DNAH1	c.7153G>A ; p.Trp2385Arg	missense	.	proba damaging	.	.	.		
Ghs129	DNAH1	.; p.Arg3169Gly	missense	.	proba damaging	.	.	8.26e-06		

Analyse n°2 Comme dit précédemment, dans cette analyse nous nous sommes concentrés sur les gènes **présents dans la liste cillio**me sur lesquels **au moins deux** de nos patients présentaient au moins 1 variant **homozygote** ayant un effet tronquant sur la protéine. Nous avons ainsi pu identifier les 4 gènes suivants : ARMC2, CCDC146, CFAP44 et TTC29 :

1. **ARMC2** : Sur ce gène, 2 patients portent un variant homozygote tronquant : *Ghs37* et *Ghs93*. Le patient et *Ghs37* portent un indel créant un décalage du cadre de lecture tandis que et *Ghs93* porte un variant créant un codon stop prématuré (**Table** : 4.9). Le patient et *Ghs107* porte quant à lui un variant faux-sens prédit comme *deleterious* par SIFT et *probably damaging* par PolyPhen. Nous pouvons également noter qu'aucun des variants portés par ces 3 sont absents des 3 bases de données. Les arguments génétiques mentionnés couplés à la forte expression testiculaire de ce gène (**Figure** : 4.18 - **B**) font de celui-ci un très bon candidat pour expliquer le phénotype de ces 3 bien que des analyses fonctionnelles soient nécessaires pour 1 d'entre eux.
2. **CFAP44** : Comme *ARMC2*, 2 patients portent un variant homozygote tronquant et 1 porte un variant homozygote non tronquant. et *Ghs22* porte un variant stop et et *Ghs34* un variant affectant un site donneur d'épissage. Le fait qu'aucun de ces variant ne soit répertorié dans aucune des bases de données laisse supposer qu'ils sont tout deux très rares. Le patient et *Ghs89* porte un variant faux-sens prédit comme *tolerated* par SIFT et *possibly damaging* par PolyPhen. Les preuves fortes impliquant ce gène dans le cillio
me humain ainsi que les effets délétères retrouvés chez 2 des 3 portant des mutations sur ce gène font de *CFAP44* un autre très bon candidat malgré son expression non spécifique au testicule (**Table** : 4.9 et 4.10, **Figure** : 4.18 - **B**).
3. **CCDC146** : Sur ce gène, seuls les patients *Ghs32* et *Ghs35* sont retrouvés mutés. Cependant, ces deux patients portent respectivement un variant induisant un codon stop prématuré et un décalage du cadre de lecture entraînant tout deux la production d'une protéine tronquée. Le premier de ces variants est retrouvé à une fréquence très faible à la fois dans ESP et ExaC, tandis que le second est totalement absent de l'ensemble des bases de données. On peut également ajouter que ce gène décrit comme faisant partie du cillio
me humain avec de fortes preuves présente également une expression testiculaire relativement élevée. Pour finir, la protéine CCDC146 codée par le gène *CCDC146* avait déjà été décrite comme composant du centrosome spermatique, un organe ayant un rôle dans l'orientation des cellules et étant à l'origine des cils et des flagelles (E. N. Firat-Karalar, Sante, Elliott, & Stearns, 2014) renforçant ainsi les arguments de l'implication de ce gène dans le phénotype MMAF nous permettent ainsi d'affirmer que ces variants sont responsables du phénotype MMAF de ces 2 patients (**Table** : 4.9 et 4.10, **Figure** : 4.18 - **B**).
4. **TTC29** : Sur ce gène, les patients *Ghs19* et *Ghs26* portent la même variation retrouvée à très faible fréquence dans les trois bases de données et impactant un site donneur d'épissage du transcrit induisant la production d'une protéine

aberrante. Ce gène à très forte expression testiculaire avait déjà été décrit en 2014 comme localisant au niveau de l'axonème du flagelle et qu'un *knock-down* entraînait des défauts de la cillio-génèse (Chung et al., 2014) (**Table : 4.9 et 4.10, Figure : 4.18 - B**).

Cette première analyse au cours de laquelle nous avons sélectionné les gènes retrouvés présents dans la liste cillio-génèse et sur lesquels au moins deux de nos patients présentaient au moins 1 variant homozygote ayant un effet tronquant sur la protéine nous a permis de mettre en évidence **rn_gene_grp1_high** nouveaux gènes candidats : *ARMC2*, *CCDC146*, *CFAP44* et *TTC29* retrouvés mutés à l'état homozygote chez 10 de nos patients dont 8 avec des variants tronquants soit 13.1 % des patients restant dans notre cohorte. Pour les 2 autres, des analyses fonctionnelles sont nécessaires afin de pouvoir être sûr que leurs variants sont bien responsables de leur phénotype. La cause génétique responsable du phénotype des patients *Ghs19*, *Ghs22*, *Ghs26*, *Ghs32*, *Ghs34*, *Ghs35*, *Ghs37* et *Ghs93* ayant été identifiées avec certitude, l'ensemble de leurs données de variants sont retirées de nos listes réduisant ainsi celle-ci à 1331 chevauchant 1164 et réparties sur 53. Les données des patients *Ghs107* et *Ghs89* sont elles conservées afin de voir si un meilleur candidat pourrait expliquer le phénotype de ces patients.

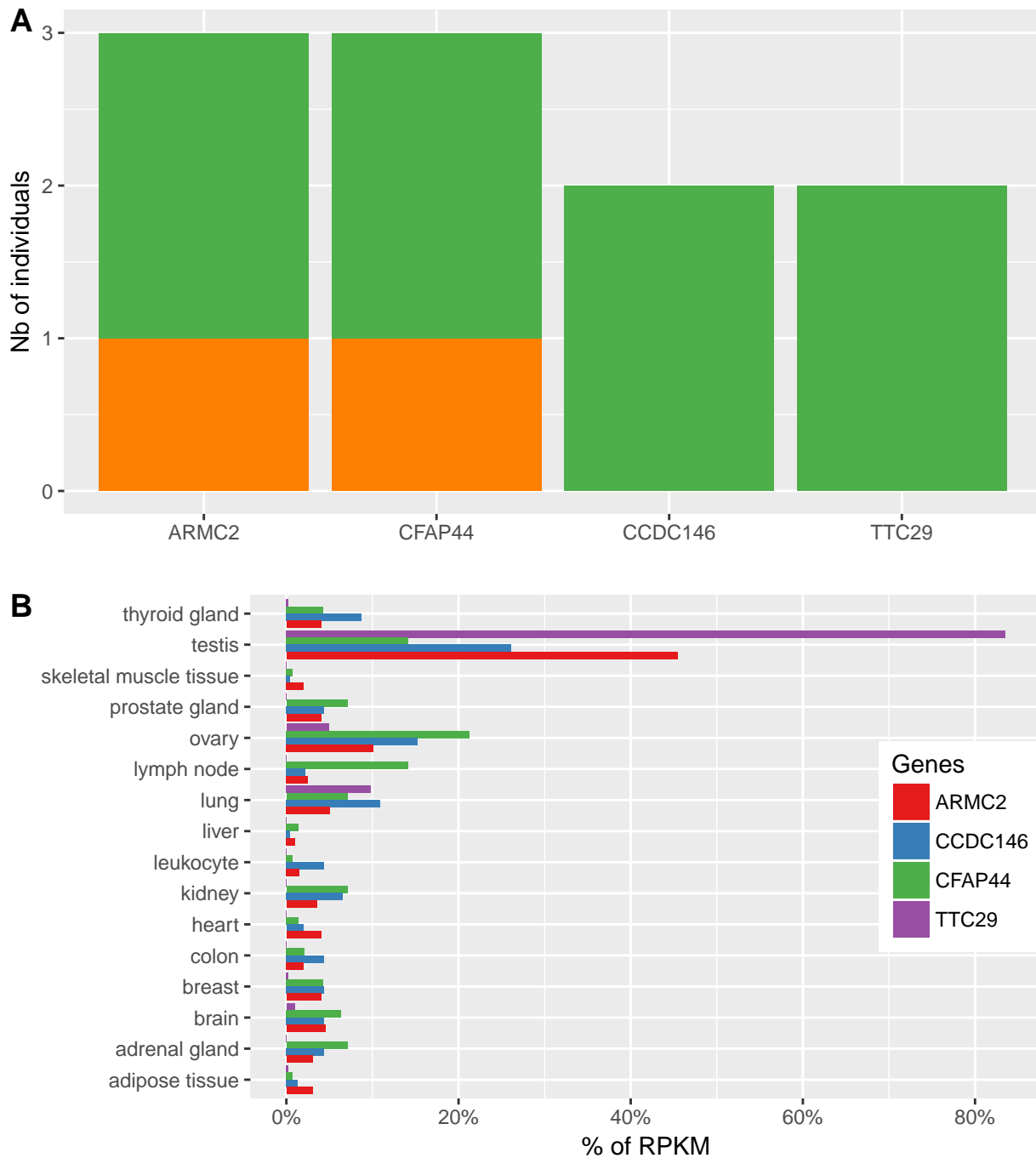


Figure 4.18 – Analyse des gènes sélectionnés dans l'Analyse n°1 : Expression tissulaire des gènes retenus d'après les données du projet de transcriptome Illumina bodyMap. Résumé de l'Analyse 1, quantification du nombre de patients retrouvés mutés sur chacun des gènes retenus ainsi que du degré de confiance accordé à la cause génétique

Table 4.9 – Analyse n°2 : List des gènes présents dans la liste ciliome sur lesquels au moins deux patients portent une mutation tronquante homozygote

Patient	Gene	Evidence	Variant impact		Variant frequency		
			HGVSc, HGVSp	Consequence	ESP	1KG	ExAC
Ghs37	ARMC2	Novel	c.2353_2354delTA ; p.Leu785MetfsTer5	frameshift	.	.	.
Ghs93	ARMC2	Novel	c.1023+1G>A ; .	splice donor	.	.	.
Ghs32	CCDC146	Strong	c.1084G>T ; p.Arg362Ter	stop gained	1e-04	.	2.47e-05
Ghs35	CCDC146	Strong	c.2112delT ; p.Arg704SerfsTer7	frameshift	.	.	.
Ghs22	CFAP44	Strong	. ; p.Arg1059Ter	stop gained	.	.	.
Ghs34	CFAP44	Strong	. ; .	splice donor	.	.	.
Ghs19	TTC29	Strong	c.176+1T>A ; .	splice donor	0.0012	5e-04	0.000158
Ghs26	TTC29	Strong	c.176+1T>A ; .	splice donor	0.0012	5e-04	0.000158

Table 4.10 – Analyse n°2 : Liste des patients portant un variant non troquant homozygote sur un des gènes suivant : ARMC2, CCDC146, CFAP44 et TTC29

Patient	Gene	HGVS _c , HGVS _p	Variant impact			Variant frequency		
			Consequence	SIFT	PolyPhen	ESP	1KG	ExAC
Ghs89	CFAP44	c.1457A>T ; p.Ala486Val	missense	tolerated	possib damaging	0.0012	0.0014	0.000692
Ghs107	ARMC2	c.2279G>A ; p.Ile760Asn	missense	deleterious	proba damaging	.	.	.

Analyse n°3 Pour rappel, au cours de cette analyse, nous avons étudié les gènes **absents dans la liste cillio**me mais sur lesquelles on trouvait toujours **au moins deux** de nos patients présentant au moins 1 variant tronquant à l'état **homozygote**. 17 patients différents portaient ainsi au moins un variant homozygote tronquant sur l'un des 8 gènes suivants : *BAZ1A*, *CCDC129*, *CFAP43*, *FSIP2*, *ICA1*, *NACA*, *SART3* et *TRAV26-1*.

1. ***CFAP43*** : 7 patients portent au moins 1 variant sur le gène *CFAP43*. Parmi ceux-ci, les patients *Ghs102*, *Ghs105*, *Ghs126*, *Ghs17* et *Ghs41* portent une mutation tronquante à l'état homozygote soit absente des bases de données soit présentent avec une très faible fréquence. Le patient et *Ghs25* lui porte un variant homozygote intronique au situé au sein de la région d'épissage (et non sur un site donneur ou accepteur). Bien que ce type de variants puissent effectivement avoir un impact sur l'épissage, il pourrait également être sans effet, or, il est difficile de le prédire à ce stade. (TODO : : faire tourner un algo de prédiction). Le patient et *Ghs132* en revanche semble plus intéressant, puisque celui-ci porte deux variant hétérozygotes sur le gène *CFAP43* parmi lesquels un est un indel entraînant un décalage du cadre de lecture tandis que l'autre est un faux-sens prédit comme *possibly damaging* par Polyphen, bien qu'il soit annoté *tolerated* par SIFT. Malgré une expression ubiquitaire (**Figure : 4.19 - B**), le nombre important de patients portant des variant sur celui-ci, notamment 5 portant des variants tronquants homozygote font de ce gène un bon candidat.
2. ***FSIP2*** : Comme pour le gène *CFAP43*, 7 de nos patients portaient 1 variant sur le gène *FSIP2* cependant pour seulement 2 d'entre eux ce variant était tronquant à l'état homozygote. En effet, les patients *Ghs20* et *Ghs21* porte tout deux un indel entraînant un décalage du cadre de lecture dont aucun n'est répertorié dans les bases de données. Le patient et *Ghs131* porte lui un faux-sens homozygote prédit comme *benign* par PolyPhen. 4 autres patients portent au moins deux variants hétérozygote sur ce gène, cependant la plupart sont des faux-sens prédit également comme *benign* par PolyPhen. Bien que l'effet sur la protéine des variants portés par 5 des 7 patients portant au moins un variant sur ce gène soit incertains, les variants tronquant portés par les patients *Ghs20* et *Ghs21* ainsi que sa forte expression testiculaire et le fait que son implication dans la structure de la gaine fibreuse du flagelle spermatique ait été montrée en 2003 (Brown, Miki, Harper, & Eddy, 2003) font de ce gène un excellent candidat pour expliquer le phénotype d'au moins 2 patients. Pour les autres, des analyses fonctionnelles seront nécessaires.
3. ***SART3*** : Un total de 4 patients portaient **le même** variants homozygote. Ces patients étant tous issus du même projet de séquençage (Strasbourg 2012) ce variant pourrait parfaitement être un artefact dû au protocole de séquençage. Cette hypothèse est d'autant plus probable que l'ensemble des patients séquencés dans ce projet présentaient un phénotype MMAF, ainsi, aucun des individus de notre cohorte contrôle n'a donc pu servir à filtrer les variants artefactuels issus de ce projet. De plus, on peut constater que l'expression tissulaire de ce gène

est ubiquitaire. Pour ces différentes raisons, ce gène n'a pas été retenu comme candidat dans nos analyses.

4. **ICA1 et TRAV26-1** : Ces deux gènes sont chacun retrouvés présentant des variants homozygotes tronquant chez 6 de nos patients. Pour chaque gène on peut noter que c'est à chaque fois le même variant qui est partagé par les trois patients ce qui est assez étonnant étant donné le fait que ces patients ne sont pas apparentés et n'ont pas été séquencés lors du même projet de séquençage. De plus, on peut observer que le patient et *Ghs40* porte un variant sur chacun de ces deux gènes, de même qu'il en portait également un sur le gène *SART3*. L'expression testiculaire de ces deux gènes apparaît elle aussi très faible. De plus, la bibliographie du gène *ICA1* révèle que celui-ci est principalement exprimé dans le pancréas (Mally, Cirulli, Hayek, & Otonkoski, 1996, Stassi, Schloot, & Pietropaolo (1997)) et que celui-ci serait à la fois lié au diabète de type 1 (S. Martin et al., 1995, R. Gaedigk, Duncan, Miyazaki, Robinson, & Dosch (1994), S. Martin, Lampasona, Dosch, & Pietropaolo (1996)) et au syndrome de Gougerot-Sjögren (S. Winer et al., 2002). Malgré le nombre important de patients présentant des variants homozygotes tronquant sur ces deux gènes, ceux-ci n'ont pas été retenus comme candidat pour expliquer le phénotype MMAF de ces 3 patients.
5. **BAZ1A** : Sur ce gène, les patients *Ghs18* et *Ghs94* portent tous deux le même indel à l'état homozygote. Celui-ci entraînant un décalage du cadre de lecture induit la formation d'une protéine aberrante. L'absence de cette indel dans les trois bases de données laisse supposer que une faible fréquence de ce variant dans la population générale. Bien que ce gène présente une faible expression testiculaire, une étude de 2013 effectuée chez la souris a démontré son rôle majeur dans la spermatogénèse, les souris *KO Baz1a*^{-/-} présentaient de nombreux défauts tel que des spermatozoïdes non motiles, présentant une morphologie de la tête et du flagelle aberrante (Dowdle et al., 2013) faisant ainsi de ce gène un bon candidat pour expliquer le phénotype de ces 2 patients.
6. **CCDC129** : Ce gène est retrouvé muté à l'état homozygote chez 2 patients porteur du même indel, non répertorié dans les bases de données, causant un décalage du cadre de lecture. On peut également constater que ce gène possède une expression testiculaire forte et exclusive faisant ainsi de ce gène un bon candidat malgré une littérature pauvre à son sujet.
7. **NACA** : 2 patients partagent le même variant homozygote causant un décalage du cadre de lecture sur le gène *NACA*. On peut noter que les patients *Ghs37* et *Ghs37* portaient respectivement de variant homozygote tronquant sur les gènes *CCDC146* et *ARMC2* portaient également ce même variant. L'ensemble de ces patients provenant du même projet de séquençage (Strasbourg 2012), cela laisse supposer que comme pour *SART3*, ce variant est artefactuel. On peut d'ailleurs noter que le patient *Ghs41* était lui aussi porteur du variant sur *SART3* renforçant ainsi l'hypothèse de l'erreur de séquençage. Pour ces raisons, ce gène n'a pas été retenu en tant que candidat.

Dans cette analyse, 8 gènes ont dans un premier temps été identifiés. Cependant, une analyse plus approfondie a fait que seuls les gènes *BAZ1A*, *CCDC129*, *CFAP43* et *FSIP2* ont été retenus. Les variants présents sur les gènes *SART3* et *NACA* étant probablement des artefacts dus aux erreurs de séquençage. Dès lors, seuls 4 des 8 gènes identifiés présentaient des arguments suffisamment convainquants pour être considérés comme responsables du phénotype MMAF. Ainsi, 18 de nos patients sont porteurs de variants sur l'un de ces gènes dont 13 à l'état homozygote avec notamment les patients *Ghs102*, *Ghs105*, *Ghs126*, *Ghs132*, *Ghs17*, *Ghs18*, *Ghs20*, *Ghs21*, *Ghs41*, *Ghs91* et *Ghs94* porteurs de variants homozygotes tronquants. Les gènes *BAZ1A*, *CCDC129*, *CFAP43* et *FSIP2* étant de bons candidats et les données génétiques des 11 patients précédemment cités étant suffisamment fortes, l'ensemble de leurs variants furent ainsi retirés de notre liste contenant désormais 1093 variants et 962 gènes différents.

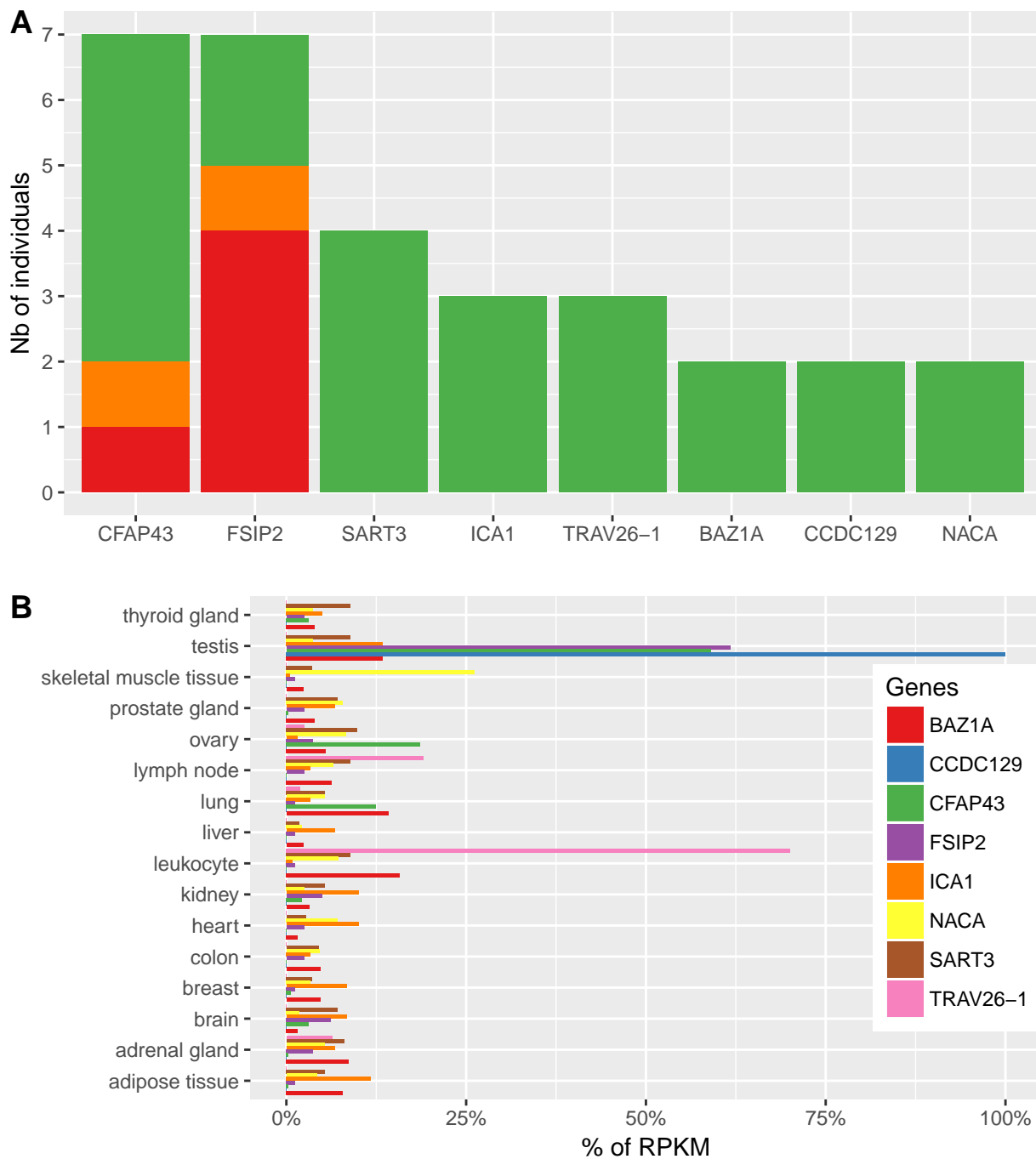


Figure 4.19 – Analyse des gènes sélectionnés dans l'Analyse n°2 : Expression tissulaire des gènes retenus d'après les données du projet de transcriptome Illumina bodyMap. Résumé de l'Analyse 2, quantification du nombre de patients retrouvés mutés sur chacun des gènes retenus ainsi que du degré de confiance accordé à la cause génétique

Table 4.11 – Analyse n°3 : List des gènes sur lesquels au moins deux patients portent une mutation tronquante non présents dans la liste ciliome

Patient	Gene	Variant impact		Variant frequency		
		HGVSc, HGVS _p	Consequence	ESP	1KG	ExAC
Ghs91	CCDC129	.; .	frameshift	.	.	.
Ghs132	CCDC129	.; .	frameshift	.	.	.
Ghs17	CFAP43	c.2658G>A ; p.Trp886Ter	stop gained	2e-04	.	9.88e-05
Ghs41	CFAP43	c.2680C>T ; p.Arg894Ter	stop gained	.	.	8.24e-06
Ghs102	CFAP43	c.3882delA ; p.Glu1294AspfsTer47	frameshift	.	.	.
Ghs105	CFAP43	c.3541-2A>C ; .	splice acceptor	.	.	.
Ghs126	CFAP43	c.3352C>T ; p.Arg1118Ter	stop gained	.	.	3.29e-05
Ghs20	FSIP2	c.2549_2550insA ; p.Asn850LysfsTer4	frameshift	.	.	.
Ghs21	FSIP2	c.1177delG ; p.Gln393LysfsTer13	frameshift	.	.	.
Ghs40	TRAV26-1	c.75_76insC ; p.Ser26LeufsTer6	frameshift	.	.	.
Ghs43	TRAV26-1	c.75_76insC ; p.Ser26LeufsTer6	frameshift	.	.	.
Ghs55	TRAV26-1	c.75_76insC ; p.Ser26LeufsTer6	frameshift	.	.	.

Table 4.12 – Analyse n°3 : Liste des patients portant un variant non troquant homozygote sur un des gènes suivant : *BAZ1A*, *CCDC129*, *CFAP43*, *FSIP2*, *ICA1*, *NACA*, *SART3* et *TRAV26-1*

Patient	Gene	Variant impact				Variant frequency		
		HGVSc, HGVSp	Consequence	SIFT	PolyPhen	ESP	1KG	ExAC
Ghs25	CFAP43	c.2141+5T>A ; .	splice region
Ghs131	FSIP2	. ; p.Ala86Val	missense	.	benign	.	.	0.00121

```
## Warning in add_footnote(., "ceci est un test"): Notation is set to 'number'  
## and other formats are not supported.
```

Table 4.13 – Analyse n°3 : Liste des patients portant un variant non troquant homozygote sur un des gènes suivant : *BAZ1A*, *CCDC129*, *CFAP43*, *FSIP2*, *ICA1*, *NACA*, *SART3* et *TRAV26-1*

Patient	Gene	HGVS _c , HGVSp	Variant impact			Variant frequency		
			Consequence	SIFT	PolyPhen	ESP	1KG	ExAC
Ghs40	FSIP2	.; p.Asn495Ile	missense	.	benign	.	0.0056	0.00157
Ghs40	FSIP2	.; p.Ile2960Met	missense	.	proba damaging	.	.	.
Ghs40	FSIP2	.; p.Lys5822Ile	missense	.	benign	0.0034	0.0056	0.0019
Ghs92	FSIP2	.; p.Asn495Ile	missense	.	benign	.	0.0056	0.00157
Ghs92	FSIP2	.; p.Lys5822Ile	missense	.	benign	0.0034	0.0056	0.0019
Ghs95	FSIP2	.; p.Asn495Ile	missense	.	benign	.	0.0056	0.00157
Ghs95	FSIP2	.; p.Lys5822Ile	missense	.	benign	0.0034	0.0056	0.0019
Ghs101	FSIP2	c.182T>C; p.Leu61Pro	missense	.	unknown	.	.	.
Ghs101	FSIP2	c.925A>T; p.Arg309Cys	missense	.	benign	.	.	.
Ghs132	CFAP43	c.1300_1301insT; p.Leu435SerfsTer26	frameshift
Ghs132	CFAP43	c.1040G>C; p.Val347Ala	missense	tolerated	possib damaging	2e-04	.	7.41e-05

Analyse n°4 Dans cette troisième analyse, nous avons sélectionné à nouveau les gènes **présents dans la liste ciliome** en conservant cette fois-ci ceux sur lesquels **un seul** de nos patients présentait au moins 1 variant tronquant à l'état **homozygote**. Les 4 gènes suivant correspondaient à ces critères : *C21orf59*, *C6orf118*, *CCDC65* et *SPEF2* tous présentant de fortes preuves indiquant leur appartenance au ciliome humain.

1. ***C6orf118*** : Ce gène a été retrouvé muté à l'état homozygote chez le patient et *Ghs40* qui portait une substitution entraînant la formation d'un codon stop prématuré. Il faut noter que ce même patient portait également des variants homozygotes tronquant sur les gènes *SART3*, *TRAV26-1*, et *ICA1* cependant la forte expression testiculaire de ce gène en fait un meilleur candidat. Le patient et *Ghs27* quant à lui porte deux variants faux-sens hétérozygote. Le premier étant prédit probablement *damaging* par PolyPhen et *tolerated low confidence* par SIFT tandis que le second est prédit *possibly damaging* et *tolerated*, il est difficile de se prononcer quant à l'effet délétère de ces deux variants sans effectuer d'analyses fonctionnelles. Il faut noter que *C6orf118* présente une forte expression dans le poulmon. De plus, ce gène a récemment été décrit comme étant associé au phénotype de tuberculose pulmonaire (E. P. Hong, Go, Kim, & Park, 2017). Cependant cela n'est en rien contradictoire avec le phénotype MMAF de ces 2 patients, le poulmon comprenant de nombreuses cellules ciliées, notamment au niveau de l'épithélium respiratoire, il n'est donc pas surprenant que des gènes du flagelle aient également une fonction au sein d'autres organes ciliés.
2. ***C21orf59* et *CCDC65*** : Les patients *Ghs127* et *Ghs88* portent tous deux des indels induisant un décalage du cadre de lecture, le premier sur le gène *C21orf59*, le second sur *CCDC65* deux gènes déjà connus pour être impliqués dans la formation des cils. En effet, la protéine NYD-SP28 (ancien nom de CCDC65) avait déjà été caractérisée comme faisant partie du flagelle spermatique (Y. Zheng et al., 2006) alors que *C21orf59* contrôlerait à la fois la mobilité et la polarisation des cils chez *zebrafish* (Jaffe et al., 2016). On note aussi que ces deux gènes ont été associés à des pathologies ciliaires chez le *zebrafish*, *Chlamydomonas* ainsi que chez l'humain, pouvant entraîner un phénotype de dyskinésie ciliaire primaire, de même, la protéine FBB18 du gène orthologue de *C21orf59* régulerait la motilité du flagelle chez *Chlamydomonas* (Austin-Tse et al., 2013). Ainsi, les arguments génétiques associés aux implications déjà avérées de ces 2 gènes dans la ciliogenèse font de *CCDC65* et *c21orf59* d'excellents candidats pour expliquer le phénotype MMAF de nos 2 patients malgré l'expression non spécifique au testicule du gène *C21orf59*.
3. ***SPEF2*** : Ce gène est retrouvé muté à l'état homozygote chez le patient et *Ghs131* qui porte un indel créant un décalage du cadre de lecture. Malgré son expression non spécifique au testicule, plusieurs études ont démontré le rôle important de la protéine SPEF2 chez la souris dans la formation et la structure du flagelle spermatique de même que l'absence de cette protéine entraînait la formation d'un flagelle court et désorganisé (Sironen et al., 2010,

Sironen et al. (2011)) faisant ainsi de ce gène un excellent candidat pour le phénotype MMAF du patient et *Ghs131*.

Analyser les gènes de la liste ciliome sur lesquels *un seul* patient portait un variant tronquant à l'état homozygote nous a permis d'identifier 4 nouveaux gènes : *C21orf59*, *C6orf118*, *CCDC65* et *SPEF2*. Ainsi, 4 patients portaient un variant homozygote tronquant sur l'un d'eux tandis qu'un autre patient portait deux variant faux-sens hétérozygotes sur le gène *C6orf118*. Parmi ces 4, seuls *CCDC65* et *C6orf118* présentaient une forte expression testiculaire, cependant tous avaient déjà été associés à la formation du cil chez l'humain ou chez une autre espèce animal. Bien que des analyses supplémentaires soient nécessaires afin de pouvoir affirmer que les deux variants hétérozygotes portés par le patients et *Ghs27* sont bien responsables de son phénotype, les arguments génétiques ainsi que les informations fournies par la littérature nous ont permis d'affirmer que ces gènes sont responsables des phénotypes des patients *Ghs127*, *Ghs131*, *Ghs40* et *Ghs88*. Les données de variants de ces derniers ont ainsi été retirés réduisant à nouveau notre liste à 934 variants et 821 gènes différents.

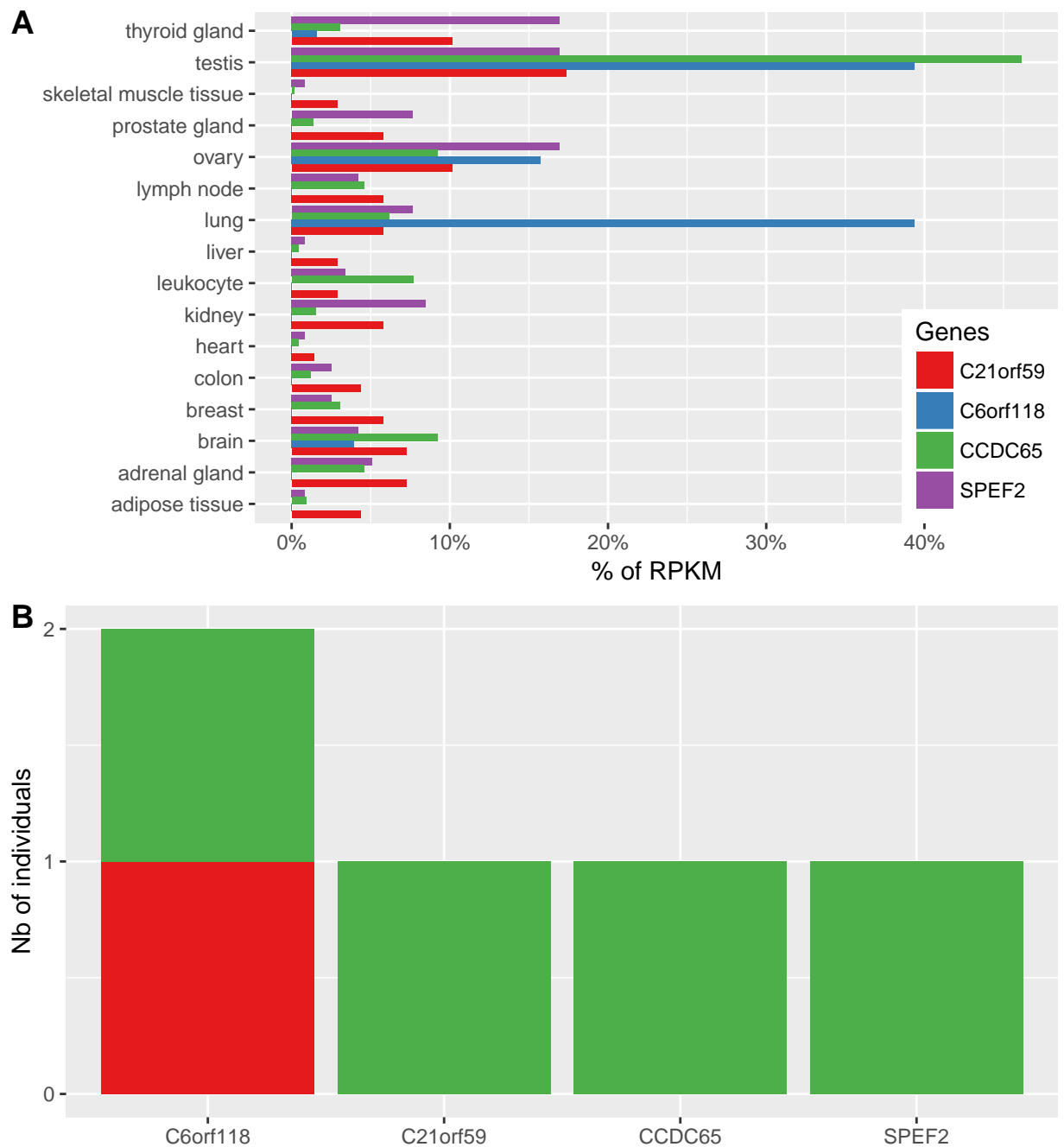


Figure 4.20 – Analyse des gènes sélectionnés dans l’Analyse n°4 :
****A**** : Expression tissulaire des gènes retenus dans cette analyse d’après les données du projet de transcriptome Illumina bodyMap. ****B**** : Résumé de l’Analyse 3, quantification du nombre de patients retrouvés mutés sur chacun des gènes retenus ainsi que du degré de confiance accordé à la cause génétique

Table 4.14 – Analyse n°4 : List des gènes présents dans la liste ciliome sur lesquels un seul patient portent une mutation homozygote tronquante

Patient	Gene	Evidence	Variant impact			Variant frequency		
			HGVSc, HGVSp	Consequence	ESP	1KG	ExAC	
Ghs40	C6orf118	Strong	c.1114C>T ; p.Arg372Ter	stop gained	.	.	8.24e-06	
Ghs127	CCDC65	Strong	c.1208delG ; p.Asn403IlefsTer9	frameshift	.	.	.	
Ghs131	SPEF2	Strong	c.3240delG ; p.Phe1080LeufsTer2	frameshift	.	.	.	

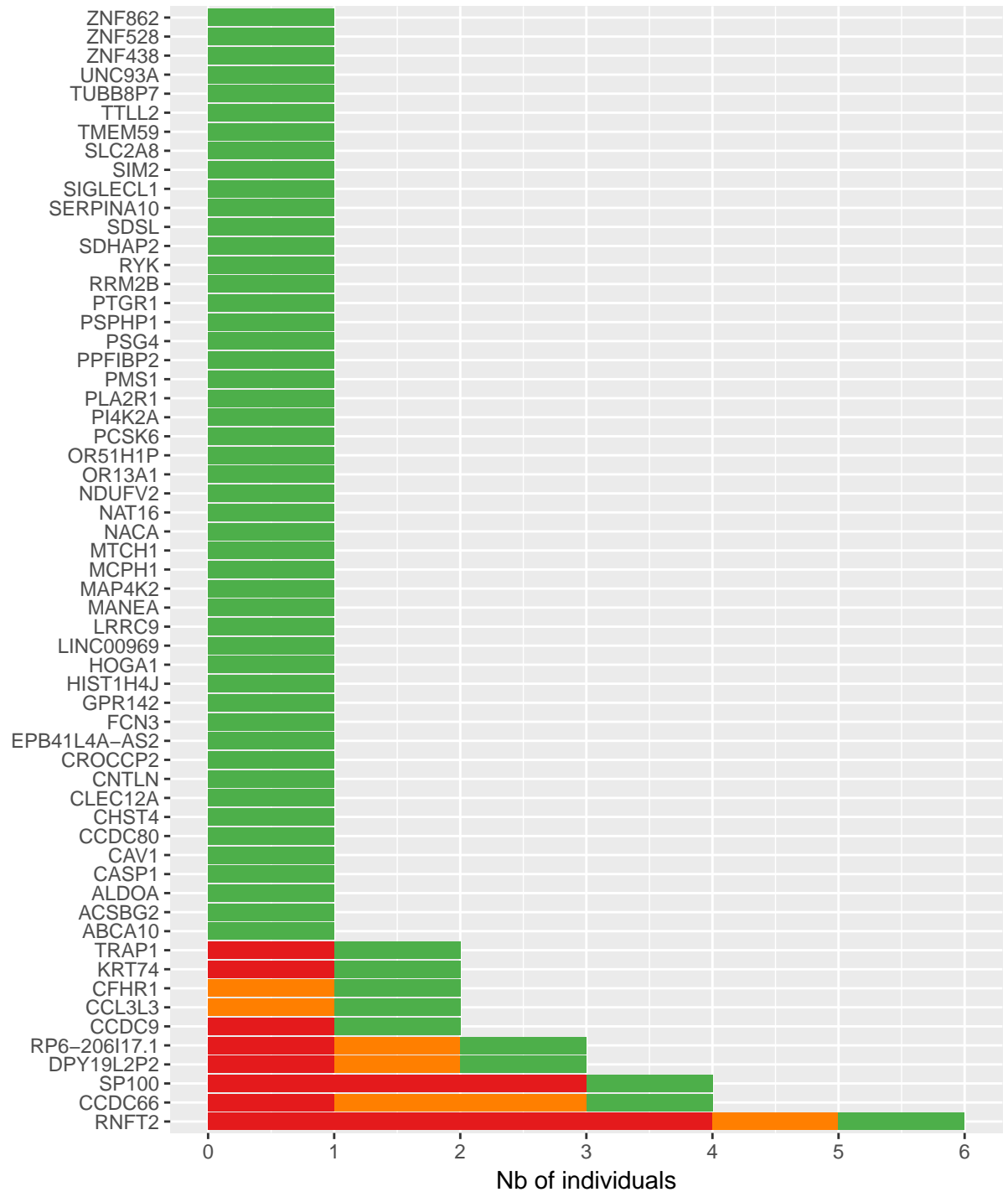
Table 4.15 – Analyse n°4 : Liste des patients portant au moins deux variants hétérozygotes sur un des gènes suivant : *C21orf59*, *C6orf118*, *CCDC65* et *SPEF2*

Patient	Gene	HGVS _c , HGVS _p	Variant impact				Variant frequency			
			Consequence	SIFT	PolyPhen	ESP	1KG	ExAC		
Ghs27	C6orf118	c.1372G>C ; p.Tyr458His	missense	tolerated	proba damaging	8e-04	0.0042	0.000297		
Ghs27	C6orf118	c.98C>T ; p.Pro33Leu	missense	tolerated	possib damaging	0.0038	9e-04	0.000684		

Analyse n°5 Pour finir nous avons dans cette analyse sélectionné l'ensemble des variants chevauchant des gènes **absents dans la liste cillio** sur lesquels **un seul** de nos patients présentaient au moins 1 variant tronquant à l'état **homozygote**. Cela nous a permis d'obtenir une liste de 59 gènes différents retrouvés mutés chez 30 de nos patients.

En raison de la grande quantité de gène, nous nous sommes tout d'abord concentré sur ceux retrouvés mutés à l'état homozygote chez au moins deux patients. Ainsi nous avons obtenus une liste de 6 gènes : *CCDC66*, *CCL3L3*, *CFHR1*, *DPY19L2P2*, *RNFT2* et *RP6-206I17.1*

1. ***CCDC66*** : Le gène *CCDC66* a été retrouvé muté à l'état homozygote chez ... patients parmi lesquels ... portait un variant tronquant entraînant un décalage du cadre de lecture et retrouvé dans aucune des bases de données. Malgré une expression testiculaire faible, l'implication récente de ce gène dans dans la ciliogénèse (Conkar et al., 2017) fait de ce gène un bon candidat bien que certaines études aient déjà démontré l'implication de celui-ci dans des pathologies rétiniennes à la fois chez le chien (Dekomien et al., 2010) mais aussi chez l'humain (Gerding et al., 2011, Khan et al. (2017)). En effet, ces anomalies rétini



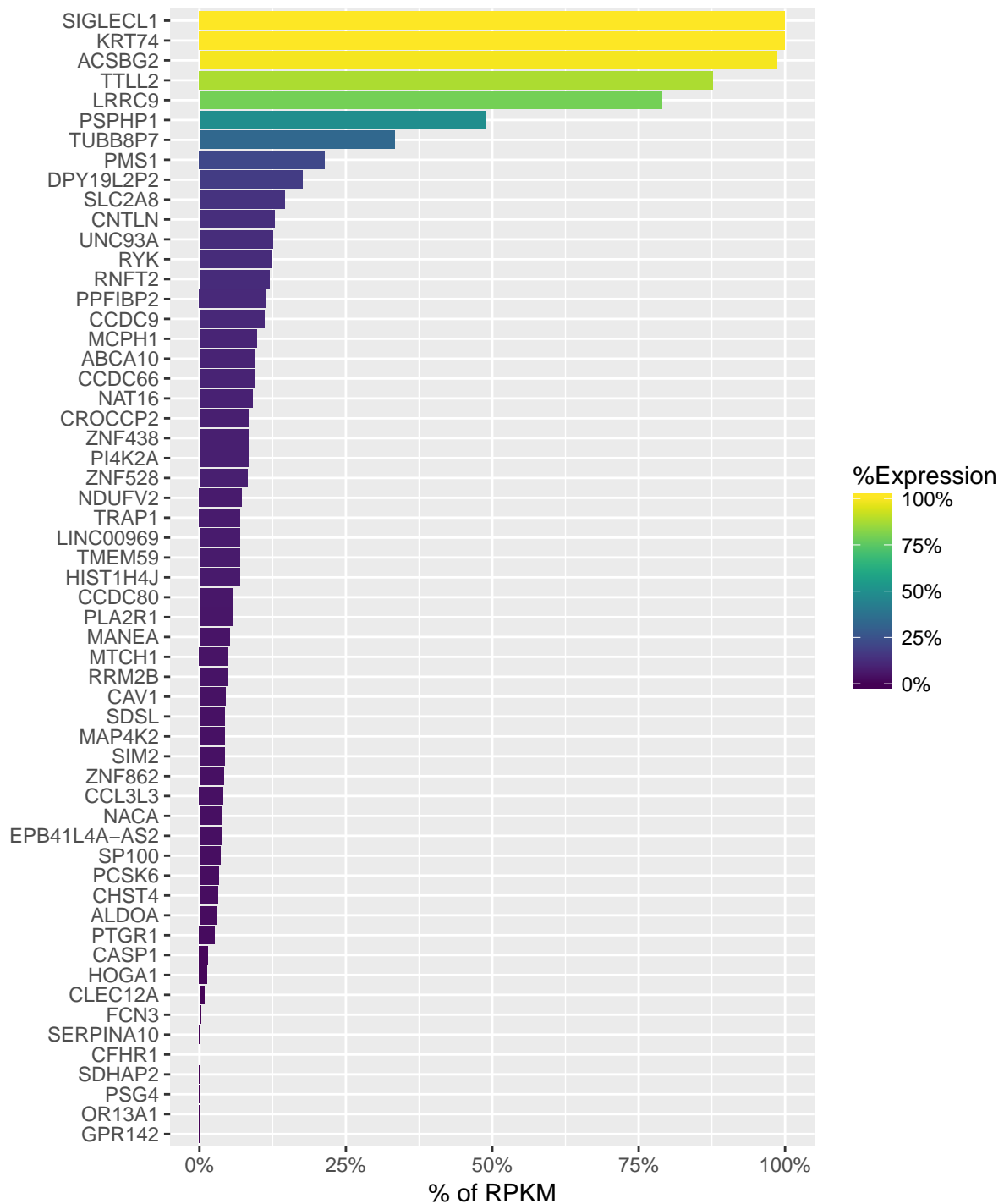


Table 4.16 – Analyse n°5 : Liste des patients portant au moins un variant homozygote tronquant sur le gène sur l'un des 59 genes identifiés dans cette analyse

Patient	Gene	Variant impact		Variant frequency		
		HGVSc, HGVS p	Consequence	ESP	1KG	ExAC
Ghs95	ABCA10	c.1356_1357delGG ; p.Ile453LeufsTer2	frameshift	.	.	.
Ghs52	ACSBG2	.;.	splice acceptor	0.0028	0.0019	0.00084
Ghs92	ALDOA	c.124_125delAG ; p.Gln42AspfsTer30	frameshift	.	.	.
Ghs47	CCDC9	c.720+2C>G ; .	splice donor	.	.	0.00343
Ghs96	CFHR1	c.790+1T>A ; .	splice donor	0.0093	0.0042	0.00231
Ghs38	CHST4	c.247C>T ; p.Lys83Ter	stop gained	.	.	1.65e-05
Ghs87	CLEC12A	c.121+1T>A ; .	splice donor	0.0013	5e-04	0.00119
Ghs42	CNTLN	c.3115-2delT ; .	splice acceptor	.	.	.
Ghs28	DPY19L2P2	.;.	splice acceptor	.	.	.
Ghs55	EPB41L4A-AS2	c.378delG ; p.Arg127GlyfsTer172	frameshift	.	.	.
Ghs42	FCN3	c.347delA ; p.Leu117SerfsTer65	frameshift	.	.	.
Ghs134	GPR142	c.215delG ; p.Pro72HisfsTer29	frameshift	.	.	.
Ghs133	HIST1H4J	c.78dupA ; p.Ile27AsnfsTer85	frameshift	.	.	.
Ghs25	HOGA1	c.208A>T ; p.Arg70Ter	stop gained	.	.	2.47e-05
Ghs87	KRT74	c.748-2T>G ; .	splice acceptor	0.0017	0.0014	0.000997
Ghs55	LRRC9	c.3256A>T ; p.Arg1086Ter	stop gained	.	.	.
Ghs134	MCPH1	c.2221G>T ; p.Arg741Ter	stop gained	1e-04	.	2.48e-05
Ghs38	MTCH1	c.28delT ; p.Trp11GlyfsTer82	frameshift	.	.	.
Ghs96	NAT16	c.808delT ; p.Ser270AlafsTer87	frameshift	.	.	.
Ghs125	OR13A1	c.804dupT ; p.Tyr269LeufsTer66	frameshift	.	.	.
Ghs90	OR51H1P	c.89delG ; p.Leu33TrpfsTer16	frameshift	.	.	.

Ghs25	PI4K2A	.; p.Arg70Ter	stop gained	.	.	2.47e-05
Ghs24	PLA2R1	c.1953delA ; p.Ala653GlnfsTer35	frameshift	.	.	.
Ghs31	PPFIBP2	.; .	frameshift	.	.	2.47e-05
Ghs33	PSG4	c.931C>T ; p.Glu311Ter	stop gained	0.0024	0.0014	0.0054
Ghs51	PSPHP1	.; .	splice donor	.	0.0051	.
Ghs87	PTGR1	.; .	splice donor	.	.	0.000231
Ghs43	RP6-206I17.1	n.211-2N>G ; .	splice acceptor	.	.	.
Ghs128	RYK	c.56_57insC ; p.Arg20GlnfsTer59	frameshift	.	.	.
Ghs87	SDHAP2	n.777+1C>A ; .	splice donor	.	.	0.00896
Ghs92	SDSL	c.685delT ; p.Leu229TrpfsTer30	frameshift	.	.	.
Ghs97	SERPINA10	c.262G>T ; p.Arg88Ter	stop gained	0.0053	0.0032	0.00749
Ghs23	SIGLECL1	c.141delG ; p.Val48TrpfsTer10	frameshift	.	.	.
Ghs96	SLC2A8	c.802G>T ; p.Gln268Ter	stop gained	7e-04	9e-04	0.000948
Ghs90	TRAP1	c.139G>T ; p.Arg47Ter	stop gained	1e-04	.	8.24e-05
Ghs87	TTL2	c.47+1T>C ; .	splice donor	.	.	.
Ghs107	UNC93A	c.676C>T ; p.Arg226Ter	stop gained	0.0024	0.0028	0.00236
Ghs87	ZNF438	c.244dupT ; p.Met83AspfsTer33	frameshift	.	.	.
Ghs23	ZNF528	c.381delG ; p.Ile129PhefsTer32	frameshift	.	.	.
Ghs38	ZNF862	.; .	splice donor	.	.	.

Table 4.17 – Analyse n°5 : Liste des patients portant au moins un variant homozygote non tronquant sur le gène sur l’un des 59 genes identifiés dans cette analyse

Patient	Gene	Variant impact				Variant frequency			
		HGVSc, HGVSnp	Consequence	SIFT	PolyPhen	ESP	1KG	ExAC	
Ghs28	DPY19L2P2	n.3419-4_3419-3insGC ; .	splice region
Ghs28	DPY19L2P2	n.3419-5_3419-4insAT ; .	splice region
Ghs36	CCL3L3	c.272G>C ; p.Leu91Pro	missense	deleterious	benign
Ghs90	CCDC66	c.11+6T>G ; .	splice region	4.03e-05	

Discussion

L'analyse de cette cohorte de 62 patients MMAF nous dans un premier temps permis de confirmer l'importance de l'implication du gène *DNAH1* dans ce phénotype grâce à 9 patients présentant des variants sur ce gène dont 3 à l'état homozygote (dont 1 tronquant). Elle nous a également permis d'identifier ... nouveaux gène candidats pouvant expliquer le phénotype de 9 de nos patients soit % de la cohorte. parmi ceux-ci, ... portaient au moins un variant homozygote sur un de ces gènes. Pour les autres des études sont nécessaires afin de déterminer si les différents variants hétérozygotes qu'ils portent sont situés sur leurs deux allèles différents faisant d'eux des hétérozygotes composites (**Figure : ?? - A**).

Parmi cet ensemble de patients, il faut noter que d'entre eux porte des variants pouvant expliquer leur phénotype sur plusieurs des gènes candidats que nous avons identifiés. En effet, n_run_2_candidat de nos patients portent des variants sur deux de nos gènes candidats et n_run_3_candidat sur 3 d'entre eux (**Figure : ?? - B**).

Cependant, parmi ces différents variants certains semblent plus probables pour expliquer le phénotype (**Table : ??**) avec par exemple :

1. Patient Ghs105 : Ce patient porte à la fois un variant homozygote affectant le site d'épissage du gène *WDR96* ainsi que deux variants hétérozygotes causant tous les deux un faux-sens sur le gène *EFCAB6*. Au vu du génotype homozygote et de l'effet tronquant du variant impactant *WDR96*, il paraît plus probable que celui-ci soit responsable du phénotype MMAF au détriment des deux variants hétérozygotes chevauchant *EFCAB6*.
2. Patient Ghs17 : Deux variants homozygotes ont été retenus pour ce patient. L'un causant un faux-sens sur le gène *EFCAB6* l'autre créant un codon stop prématuré sur *WDR96*. Ici aussi, au vu de l'impact délétère du variant chevauchant *WDR96*, il paraît plus probable que ce soit celui-ci qui soit la cause du phénotype de ce patient.
3. Patients Ghs32 et Ghs35 : Ces patients portent tous deux à la fois des variants sur le gène *ANKRD20A3* et sur le gène *CCDC146* cependant l'effet tronquant de leur variant impactant ce dernier nous laisse penser que ceux-ci soient la cause du phénotype MMAF de ces deux patients.

En procédant de la même manière pour les autres patients il est, dans la grande majorité des cas, possible de dégager un gène pour lequel l'implication dans le phénotype des patients paraît plus évidente bien que des analyses complémentaires soient nécessaires.

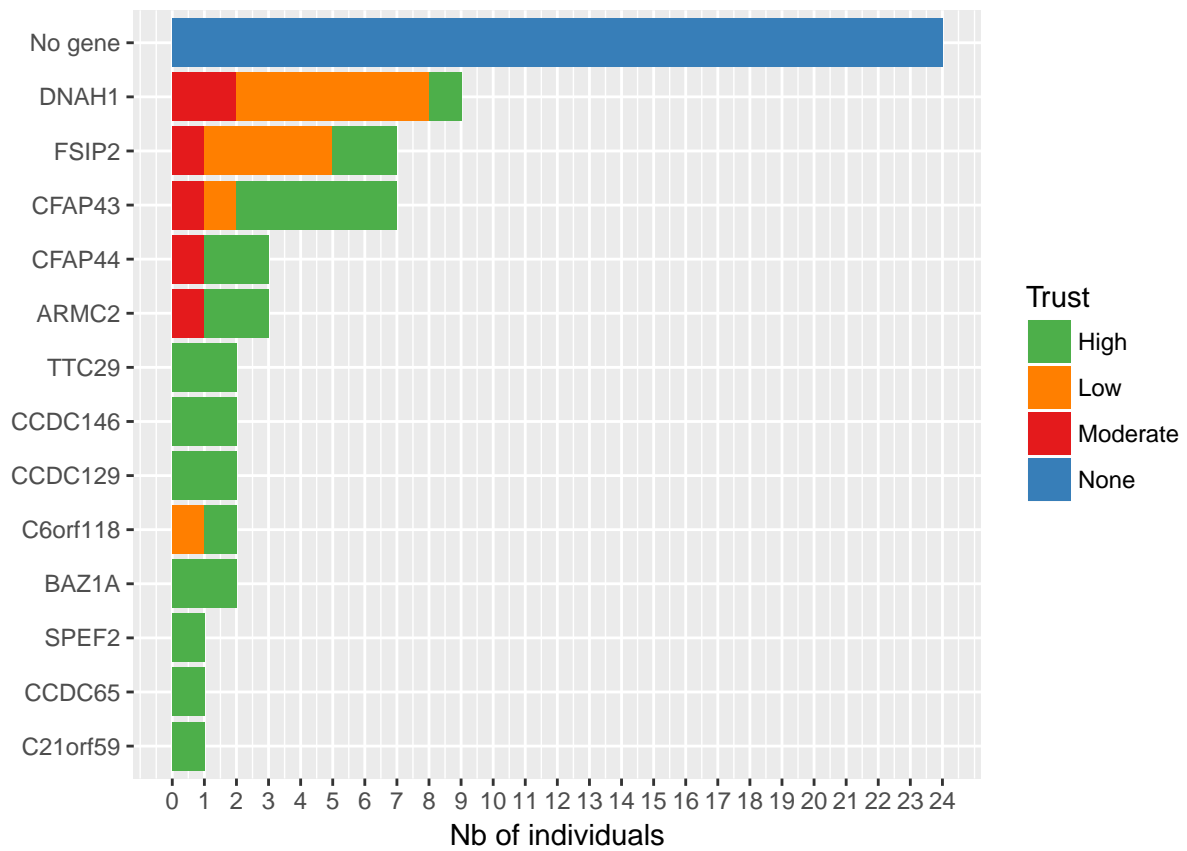
Ainsi, cette analyse révèle l'efficacité de notre pipeline puisqu'elle a permis d'identifier au moins un gène candidat pour p_run_bigmmaf_found % de nos patients. Pour les

autres des analyses individuelles complémentaires sont nécessaires afin d'identifier la cause génétique responsable de leur phénotype.

Une partie de ces différents résultats ont déjà été publiés dans deux articles dont je suis co-auteur :

1. **Whole exome cohort study and analysis of mouse and Trypanosoma models demonstrate the importance of WDR proteins in flagellogenesis and male fertility**, *Nat Genet* (soumis) : Cette article présente nos différents résultats dans la caractérisation des gènes *WDR96* et *WDR52* ainsi que les différentes preuves de leur implication dans le phénotype MMAF.
2. **Whole-exome sequencing of familial cases of multiple morphological abnormalities of the sperm flagella (MMAF) reveals new DNAH1 mutations** : En plus des résultats évoqués précédemment pour la famille MMAF2, cet article inclus ceux de ... patients de cette cohorte présentant des variants sur le gène *DNAH1*

Pour les autres, notre équipe travaille actuellement à la caractérisation des différents gènes afin de comprendre les processus moléculaires



4.3 Conclusion

Au cours de ces différentes études nous avons pu identifier les variants pouvant expliquer les phénotypes de ... des différents patients que nous avons analysé que ce soit au sein d'études familiales ou bien au sein de plus large cohorte composés d'individus non apparentés. Bien que ces résultats soient satisfaisant, il faut noter que pour ... patients, soit ... % d'entre eux aucun candidat n'a pu à ce jour être identifié. Pour ces patients, le WES permet désormais de nouvelles approches permettant d'identifier de larges variants structuraux (insertion ou délétions) pouvant eux aussi être responsable du phénotype qui ne sont pas détectés par les analyses classiques. Néanmoins, il semble clair que des avancées soient encore nécessaires afin d'améliorer l'efficacité de ce genre d'étude notamment en créant de nouveaux filtres permettant ainsi d'épurer les listes de variants facilitant ainsi l'identification des gènes candidats.

Chapitre 5

MutaScript

Conclusion

Chapitre 6

Tables des analyses d'exomes

References

- 1000 Genomes Project Consortium, T. 1. G. P., Auton, A., Brooks, L. D., Durbin, R. M., Garrison, E. P., Kang, H. M., ... Abecasis, G. R. (2015). A global reference for human genetic variation. *Nature*, 526(7571), 68–74. <http://doi.org/10.1038/nature15393>
- Adzhubei, I. A., Schmidt, S., Peshkin, L., Ramensky, V. E., Gerasimova, A., Bork, P., ... Sunyaev, S. R. (2010). A method and server for predicting damaging missense mutations. *Nature Methods*, 7(4), 248–9. <http://doi.org/10.1038/nmeth0410-248>
- Aken, B. L., Achuthan, P., Akanni, W., Amode, M. R., Bernsdorff, F., Bhai, J., ... Flicek, P. (2017). Ensembl 2017. *Nucleic Acids Research*, 45(D1), D635–D642. <http://doi.org/10.1093/nar/gkw1104>
- Amberger, J., Bocchini, C., & Hamosh, A. (2011). A new face and new challenges for Online Mendelian Inheritance in Man (OMIM). *Human Mutation*, 32(5), 564–567. <http://doi.org/10.1002/humu.21466>
- Amdani, S. N., Jones, C., & Coward, K. (2013). Phospholipase C zeta (PLC ζ) : Oocyte activation and clinical links to male factor infertility. *Advances in Biological Regulation*, 53(3), 292–308. <http://doi.org/10.1016/j.jbior.2013.07.005>
- Arnaiz, O., Cohen, J., Tassin, A., & Koll, F. (2015). Remodeling Cildb, a popular database for cilia and links for ciliopathies. *Cilia*, 4, P21. <http://doi.org/10.1186/2046-2530-4-S1-P21>
- Austin-Tse, C., Halbritter, J., Zariwala, M. A., Gilberti, R. M., Gee, H. Y., Hellman, N., ... Hildebrandt, F. (2013). Zebrafish Ciliopathy Screen Plus Human Mutational Analysis Identifies C21orf59 and CCDC65 Defects as Causing Primary Ciliary Dyskinesia. *American Journal of Human Genetics*, 93(4), 672–86. <http://doi.org/10.1016/j.ajhg.2013.08.015>
- Baker, K. E., & Parker, R. (2004). Nonsense-mediated mRNA decay : terminating erroneous gene expression. *Current Opinion in Cell Biology*, 16(3), 293–9. <http://doi.org/10.1016/j.ceb.2004.03.003>
- Ben Khelifa, M., Coutton, C., Zouari, R., Karaouzène, T., Rendu, J., Bidart, M., ... Ray, P. F. (2014). Mutations in DNAH1, which encodes an inner arm heavy chain dynein, lead to male infertility from multiple morphological abnormalities of

the sperm flagella. *American Journal of Human Genetics*, 94(1), 95–104. <http://doi.org/10.1016/j.ajhg.2013.11.017>

Brown, P. R., Miki, K., Harper, D. B., & Eddy, E. M. (2003). A-Kinase Anchoring Protein 4 Binding Proteins in the Fibrous Sheath of the Sperm Flagellum. *Biology of Reproduction*, 68(6), 2241–2248. <http://doi.org/10.1095/biolreprod.102.013466>

Chang, Y.-F., Imam, J. S., & Wilkinson, M. F. (2007). The Nonsense-Mediated Decay RNA Surveillance Pathway. *Annual Review of Biochemistry*, 76(1), 51–74. <http://doi.org/10.1146/annurev.biochem.76.050106.093909>

Chung, M.-I., Kwon, T., Tu, F., Brooks, E. R., Gupta, R., Meyer, M., ... Wallingford, J. B. (2014). Coordinated genomic control of ciliogenesis and cell movement by RFX2. *ELife*, 3, e01439. <http://doi.org/10.7554/eLife.01439>

Conkar, D., Culfa, E., Odabasi, E., Rauniyar, N., Yates, J. R., & Firat-Karalar, E. N. (2017). The centriolar satellite protein CCDC66 interacts with CEP290 and functions in cilium formation and trafficking. *Journal of Cell Science*, 130(8). Retrieved from <http://jcs.biologists.org/content/130/8/1450.long>

Dekomien, G., Vollrath, C., Petrasch-Parwez, E., Boev ?, M. H., Akkad, D. A., Gerding, W. M., & Epplen, J. T. (2010). Progressive retinal atrophy in Schapendoes dogs : mutation of the newly identified CCDC66 gene. *Neurogenetics*, 11(2), 163–174. <http://doi.org/10.1007/s10048-009-0223-z>

DePristo, M. A., Banks, E., Poplin, R., Garimella, K. V., Maguire, J. R., Hartl, C., ... Pritchard, E. (2011). A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nature Genetics*, 43(5), 491–498. <http://doi.org/10.1038/ng.806>

Dowdle, J. A., Mehta, M., Kass, E. M., Vuong, B. Q., Inagaki, A., Egli, D., ... Keeney, S. (2013). Mouse BAZ1A (ACF1) is dispensable for double-strand break repair but is essential for averting improper gene expression during spermatogenesis. *PLoS Genetics*, 9(11), e1003945. <http://doi.org/10.1371/journal.pgen.1003945>

Firat-Karalar, E. N., Sante, J., Elliott, S., & Stearns, T. (2014). Proteomic analysis of mammalian sperm cells identifies new components of the centrosome. *Journal of Cell Science*, 127(Pt 19), 4128–33. <http://doi.org/10.1242/jcs.157008>

Gaedigk, R., Duncan, A. M., Miyazaki, I., Robinson, B. H., & Dosch, H. M. (1994). ICA1 encoding p69, a protein linked to the development of type 1 diabetes, maps to human chromosome 7p22. *Cytogenetics and Cell Genetics*, 66(4), 274–6. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/8162706>

Gerding, W. M., Schreiber, S., Schulte-Middelmann, T., de Castro Marques, A., Atorf, J., Akkad, D. A., ... Petrasch-Parwez, E. (2011). Ccdc66 null mutation causes retinal degeneration and dysfunction. *Human Molecular Genetics*, 20(18), 3620–3631.

<http://doi.org/10.1093/hmg/ddr282>

Hong, E. P., Go, M. J., Kim, H.-L., & Park, J. W. (2017). Risk prediction of pulmonary tuberculosis using genetic and conventional risk factors in adult Korean population. *PloS One*, 12(3), e0174642. <http://doi.org/10.1371/journal.pone.0174642>

Hu, Y., Yu, H., Shaw, G., Renfree, M. B., & Pask, A. J. (2011). Differential roles of TGIF family genes in mammalian reproduction. *BMC Developmental Biology*, 11, 58. <http://doi.org/10.1186/1471-213X-11-58>

Imai, Y., Morita, H., Takeda, N., Miya, F., Hyodo, H., Fujita, D., ... Komuro, I. (2015). A deletion mutation in myosin heavy chain 11 causing familial thoracic aortic dissection in two Japanese pedigrees. *International Journal of Cardiology*, 195, 290–292. <http://doi.org/10.1016/j.ijcard.2015.05.178>

Ivliev, A. E., 't Hoen, P. A. C., Roon-Mom, W. M. C. van, Peters, D. J. M., & Sergeeva, M. G. (2012). Exploring the Transcriptome of Ciliated Cells Using In Silico Dissection of Human Tissues. *PLoS ONE*, 7(4), e35618. <http://doi.org/10.1371/journal.pone.0035618>

Jaffe, K. M., Grimes, D. T., Schottenfeld-Roames, J., Werner, M. E., Ku, T.-S. J., Kim, S. K., ... Burdine, R. D. (2016). c21orf59/kurly Controls Both Cilia Motility and Polarization. *Cell Reports*, 14(8), 1841–9. <http://doi.org/10.1016/j.celrep.2016.01.069>

Khan, A. O., Budde, B. S., Nürnberg, P., Kawalia, A., Lenzner, S., & Bolz, H. J. (2017). Genome-wide linkage and sequence analysis challenge CCDC66 as a human retinal dystrophy candidate gene and support a distinct NMNAT1-related fundus phenotype. *Clinical Genetics*. <http://doi.org/10.1111/cge.13022>

Kumar, P., Henikoff, S., & Ng, P. C. (2009). Predicting the effects of coding non-synonymous variants on protein function using the SIFT algorithm. *Nature Protocols*, 4(7), 1073–1081. <http://doi.org/10.1038/nprot.2009.86>

Lee, B., Park, I., Jin, S., Choi, H., Kwon, J. T., Kim, J., ... Cho, C. (2011). Impaired spermatogenesis and fertility in mice carrying a mutation in the Spink2 gene expressed predominantly in testes. *The Journal of Biological Chemistry*, 286(33), 29108–17. <http://doi.org/10.1074/jbc.M111.244905>

Lek, M., Karczewski, K. J., Minikel, E. V., Samocha, K. E., Banks, E., Fennell, T., ... Exome Aggregation Consortium, D. G. (2016). Analysis of protein-coding genetic variation in 60,706 humans. *Nature*, 536(7616), 285–91. <http://doi.org/10.1038/nature19057>

Lunter, G., & Goodson, M. (2011). Stampy : A statistical algorithm for sensitive and fast mapping of Illumina sequence reads. *Genome Research*, 21(6), 936–939. <http://doi.org/10.1101/gr.111120.110>

Mally, M. I., Cirulli, V., Hayek, A., & Otonkoski, T. (1996). ICA69 is expressed equally in the human endocrine and exocrine pancreas. *Diabetologia*, 39(4), 474–80. Retrieved

from <http://www.ncbi.nlm.nih.gov/pubmed/8777998>

Martin, S., Kardorf, J., Schulte, B., Lampeter, E. F., Gries, F. A., Melchers, I., ... Pfützner, A. (1995). Autoantibodies to the islet antigen ICA69 occur in IDDM and in rheumatoid arthritis. *Diabetologia*, 38(3), 351–5. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/7758883>

Martin, S., Lampasona, V., Dosch, M., & Pietropaolo, M. (1996). Islet cell autoantigen 69 antibodies in IDDM. *Diabetologia*, 39(6), 747. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/8781774>

McLaren, W., Gil, L., Hunt, S. E., Riat, H. S., Ritchie, G. R. S., Thormann, A., ... Cunningham, F. (2016). The Ensembl Variant Effect Predictor. *Genome Biology*, 17(1), 122. <http://doi.org/10.1186/s13059-016-0974-4>

Ng, S. B., Buckingham, K. J., Lee, C., Bigam, A. W., Tabor, H. K., Dent, K. M., ... Bamshad, M. J. (n.d.). Exome sequencing identifies the cause of a Mendelian disorder. <http://doi.org/10.1038/ng.499>

Nielsen, R., Paul, J. S., Albrechtsen, A., & Song, Y. S. (2011). Genotype and SNP calling from next-generation sequencing data. *Nature Reviews. Genetics*, 12(6), 443–51. <http://doi.org/10.1038/nrg2986>

Pelak, K., Shianna, K. V., Ge, D., Maia, J. M., Zhu, M., Smith, J. P., ... Goldstein, D. B. (2010). The characterization of twenty sequenced human genomes. *PLoS Genetics*, 6(9), e1001111. <http://doi.org/10.1371/journal.pgen.1001111>

Sironen, A., Hansen, J., Thomsen, B., Andersson, M., Vilkkilä, J., Toppari, J., & Kotaja, N. (2010). Expression of SPEF2 During Mouse Spermatogenesis and Identification of IFT20 as an Interacting Protein1. *Biology of Reproduction*, 82(3), 580–590. <http://doi.org/10.1095/biolreprod.108.074971>

Sironen, A., Kotaja, N., Mulhern, H., Wyatt, T. A., Sisson, J. H., Pavlik, J. A., ... Lee, L. (2011). Loss of SPEF2 function in mice results in spermatogenesis defects and primary ciliary dyskinesia. *Biology of Reproduction*, 85(4), 690–701. <http://doi.org/10.1095/biolreprod.111.091132>

Stassi, G., Schloot, N., & Pietropaolo, M. (1997). Islet cell autoantigen 69 kDa (ICA69) is preferentially expressed in the human islets of Langerhans than exocrine pancreas. *Diabetologia*, 40(1), 120–2. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/9028728>

Su, Z., Łabaj, P. P., Li, S. S., Thierry-Mieg, J., Thierry-Mieg, D., Shi, W., ... Shi, L. (2014). A comprehensive assessment of RNA-seq accuracy, reproducibility and information content by the Sequencing Quality Control Consortium. *Nature Biotechnology*, 32(9), 903–14. <http://doi.org/10.1038/nbt.2957>

Winer, S., Astsaturov, I., Cheung, R., Tsui, H., Song, A., Gaedigk, R., ... Dosch, H.-M. (2002). Primary Sjögren's syndrome and deficiency of ICA69. *The Lancet*, 360(9339),

1063–1069. [http://doi.org/10.1016/S0140-6736\(02\)11144-5](http://doi.org/10.1016/S0140-6736(02)11144-5)

Zheng, Y., Zhang, J., Wang, L., Zhou, Z., Xu, M., Li, J., & Sha, J.-H. (2006). Cloning and characterization of a novel sperm tail protein, NYD-SP28. *International Journal of Molecular Medicine*, 18(6), 1119–25. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/17089017>