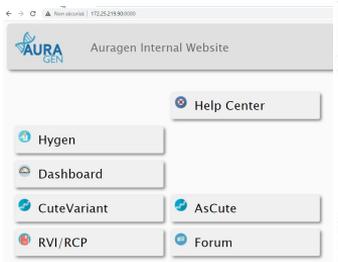
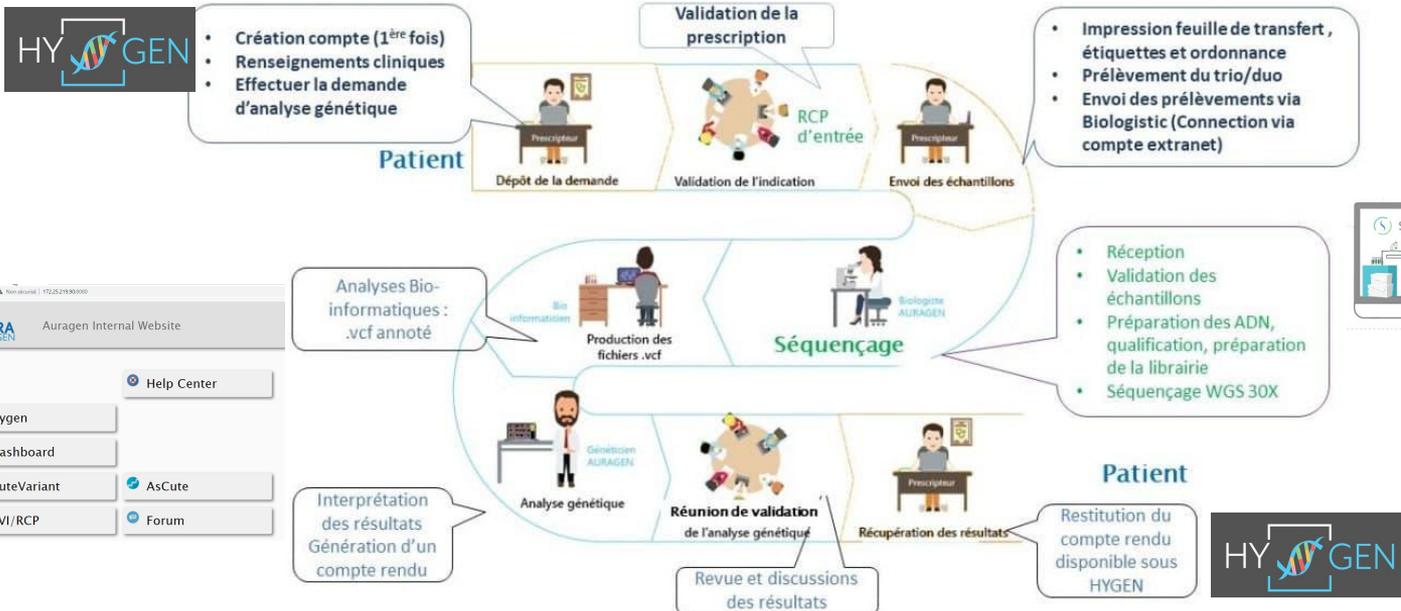


Plateforme Auragen Infra, routine & dev bioinfo

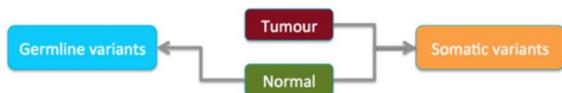
15 mai 2024

Bioinfodiag

Virginie Bernard, Auragen



Cancer

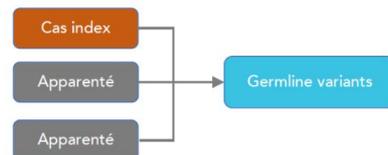


WGS tumeur 60-90x
normal 30-55x

WES tumeur 300x
normal 150x

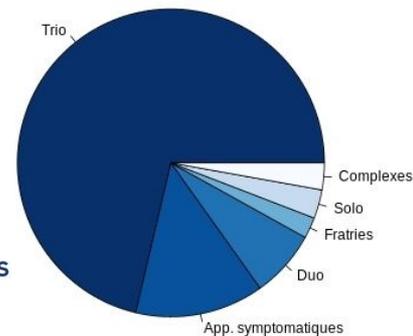
RNAseq tumeur >64 Mp
(fusions)

Maladies Rares & Oncogénétique



WGS 30-55x

Un ou plusieurs échantillons
au sein d'une famille





Organisation équipes bioinfo Auragen CLB (Lyon) -& INRIA, CHUGA (Grenoble)

Alain Viari - responsable SI

Jean François Scariot - ingénieur - infra

Maëlle Martinet Gerphagnon - ingénieur info

Anthony Ferrari - responsable bioinfo KC

Anne Sophie Sertier - bioinfo dev

Elise Dugat ingénieurs - bioinfo prod

Julien Thevenon & Virginie Bernard -
responsables bioinfo MR

**Quentin Charret & Clément Lionnet -
ingénieurs bioinfo prod & dev**

+ internes en génétiques

Alexandre, Laury, Bertrand,
Benjamin, Alexis ...



Identification et spécification d'un besoin
→ cahier des charges

1. Revue de littérature
2. Identification d'un gold standard et/ou de jeux de données internes pour l'évaluation
3. Comparaison de méthodes et sélection de la plus appropriée, customisations
4. Création et intégration d'un pipeline en suivant les recommandations de l'INCa
5. Validation du pipeline sur jeux de données internes et gold standard
6. Industrialisation et mise en production versionnée

GENOME IN A BOTTLE

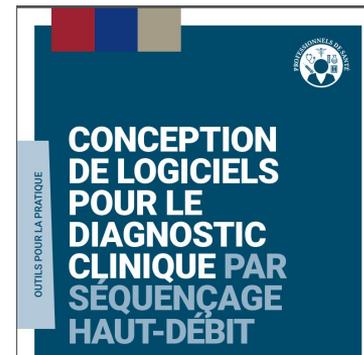
Our mission is to provide the authoritative characterization of human genomes.



nextflow



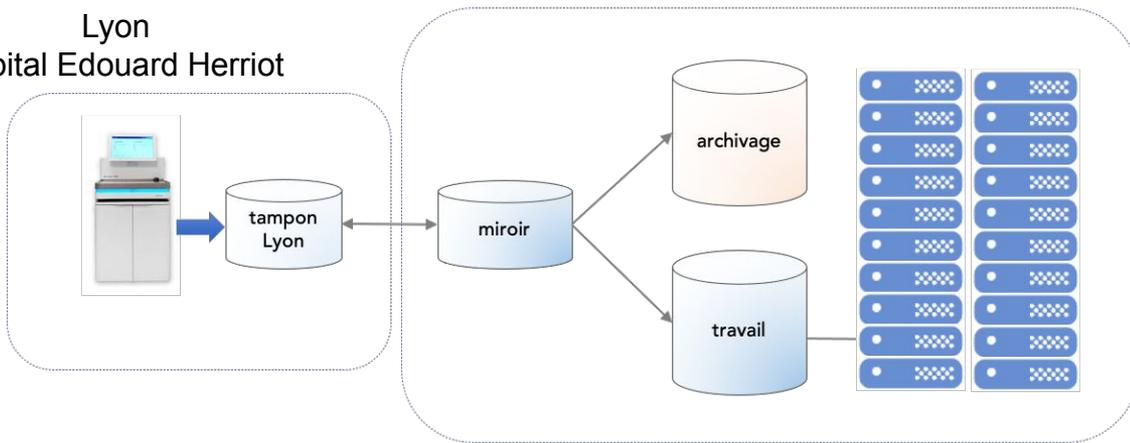
- ✓ Traçabilité
- ✓ Reproductibilité
- ✓ Fiabilité



Quelques chiffres

DataCenter – Eolas - Grenoble

Lyon
Hopital Edouard Herriot



- 20 noeuds de calcul
- 1500 CPUs
- 13 To RAM
- 2.5 Po stockage

Configuration HPC

Atos

Hébergement HDS

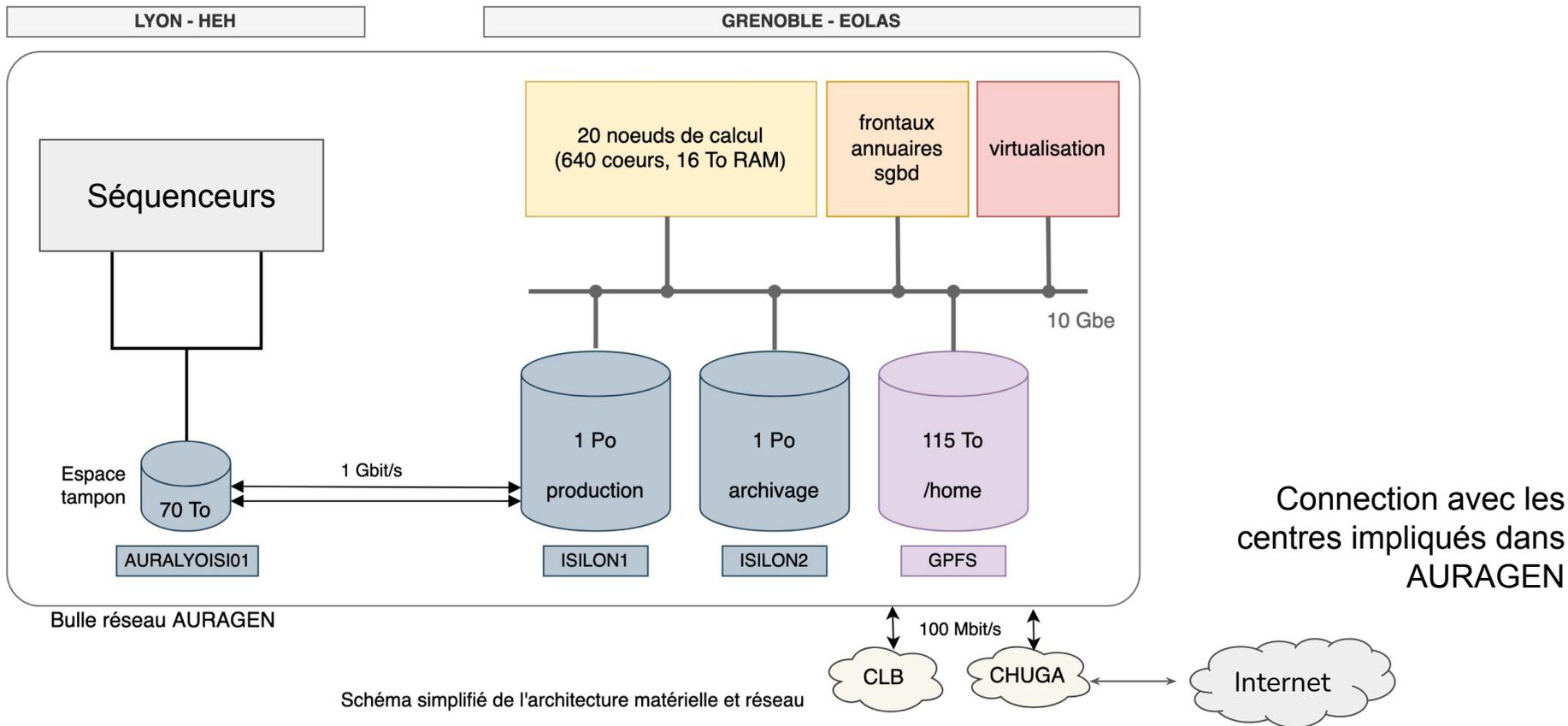
eolas
GROUPES BUSINESS & DECISION

Maintenance

AZCOM
RESADIA



Schéma simplifié de l'infrastructure



System Usage

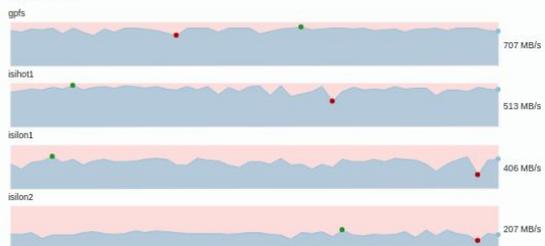
Disk Usage



Slurm usage



Disk IO Speed



Whoard Tasks

State

Filter:

task	context	state	update time
VepGenotype	MR-2400099	Running	2024-05-14 23:04:30
CoronaQC	MR-2400305	Done-Success	2024-05-14 23:04:02
Mantis	MR-2304778	Running	2024-05-14 23:03:10
GenotypeGivf	MR-2400099	Done-Success	2024-05-14 22:55:55
VepGenotype	MR-2304778	Done-Success	2024-05-14 22:55:47
VepGenotype	MR-2400052	Running	2024-05-14 22:52:19
WgsMatchParQC	MR-2400328	Done-Success	2024-05-14 22:50:49
VepGenotype	MR-2305169	Running	2024-05-14 22:38:44
GenotypeGivf	MR-2400052	Done-Success	2024-05-14 22:37:17
GenotypeGivf	MR-2305169	Done-Success	2024-05-14 22:27:48

Log Level top is most recent

level	message
log	[24:05:15-01:04:30] [job-task] notify launched task: VepGenotype context: MR-2400099
log	[24:05:15-01:04:29] [job-batch] VepGenotype-MR-2400099 submit: Submitted batch job 13454553 status: 0



Fichier brut produit lors du séquençage → conduit à un fichier FASTQ

Sample Sheet → informations librairies qui permet entre autre d'associer un identifiant de prélèvement à ses lectures

Métadonnées → informations cliniques sur le prélèvement, le dossier

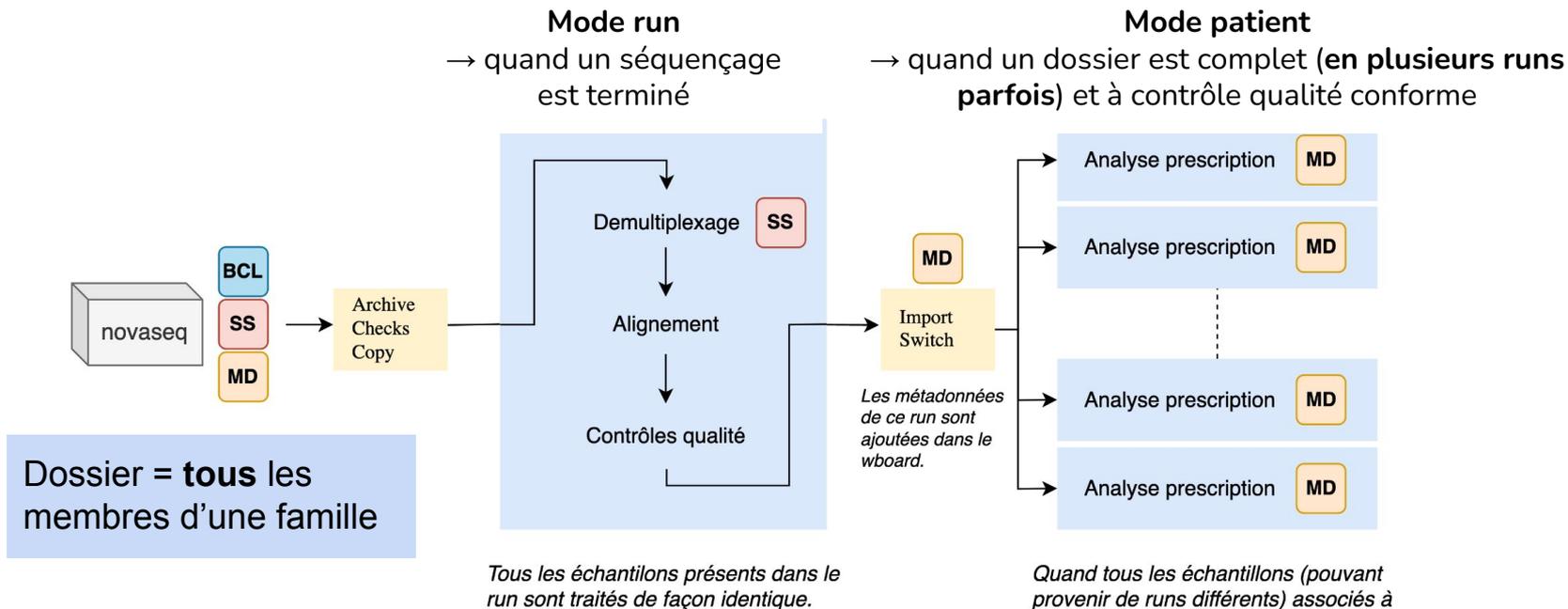
SS : Samplesheet émis par le SGL

MD : Métadonnées associées à la prescription. Transmises par le SGL et l'outil de prescription

Schéma simplifié des processus

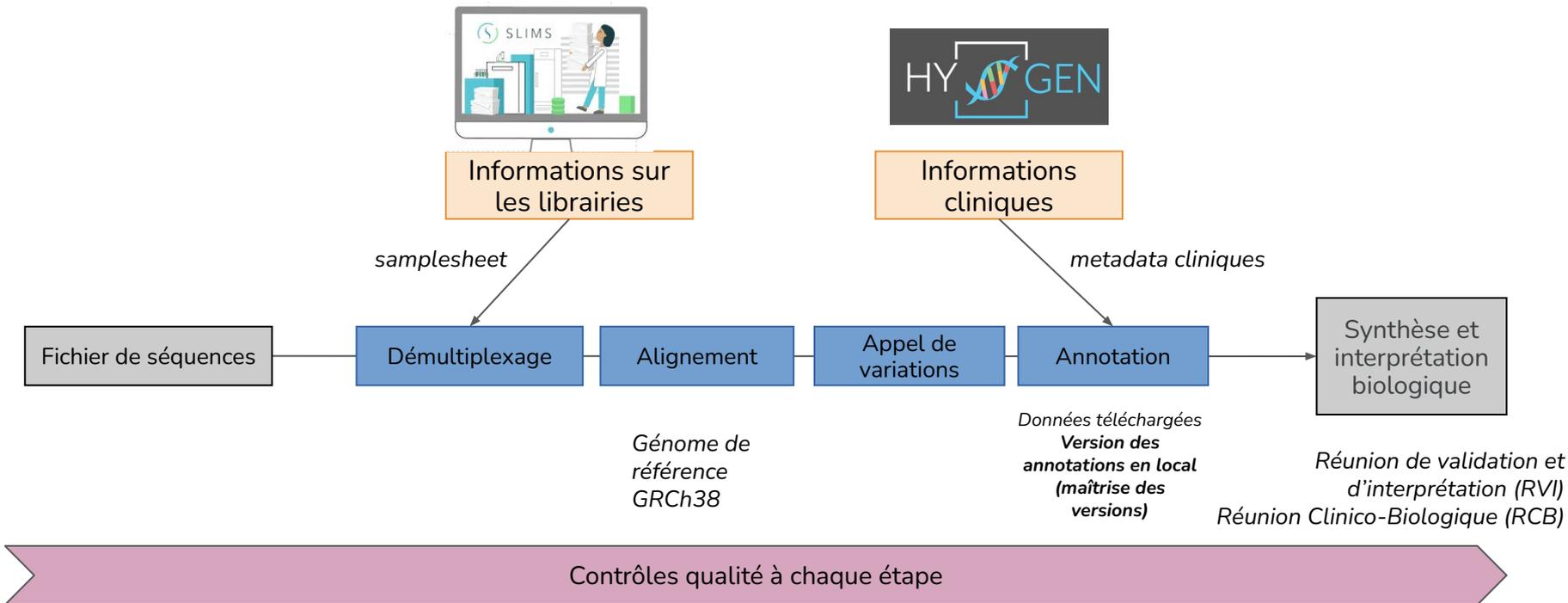
Séquençage des échantillons dans le flux de leur réception.

Traitement bioinformatique une fois les dossiers complets avec QC conformes



SS : Samplesheet émis par le SGL

MD : Métadonnées associées à la prescription. Transmises par le SGL et l'outil de prescription



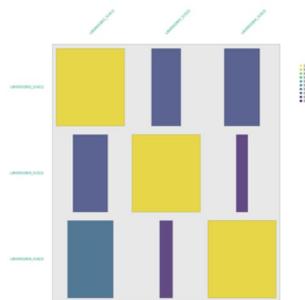
Les contrôles d'identité vigilance requis "conformes" pour une interprétation

Quels contrôles de l'identité vigilance sont mis en place en **post séquençage**?

Les informations patient sont elles conformes à la prescription?

Le genre observé est il bien conforme à la prescription ?
→ ratio profondeur chrY chrX

Organe	IDIS	Exemple	Index	Type	Sexe	Sexual	Impact	Genre
SN-2000976	NR0012	LR00000902_014	1	PR	F	✓	✓	✓
SN-2000976	NR0012	LR00000902_015	2	PR	F	✓	✓	✓
SN-2000976	NR0012	LR00000904_016	3	PR	M	✓	✓	✓



Genre conforme ✓

Similarité entre prélèvements conforme à la prescription?
→ 100000 SNP chr1 à 22

Similarité entre prélèvements conforme ✓

Est ce que l'échantillon est contaminé?
La contamination fausse-t-elle les résultats ?

→ 2589 positions ayant 3 allèles possibles dans la population
Présence effective de plus de 2 allèles ?
Si contamination, taux impactant les résultats ?

Pas de contamination = conforme
Contamination sans impact sur les variations = conforme ✓

Ce qui est présenté pour interprétation a été validé par tous les contrôles mis en place

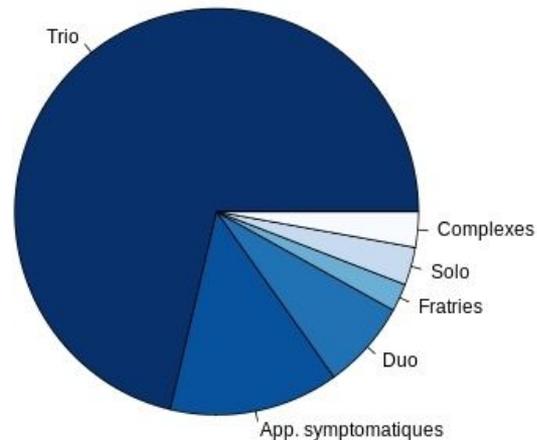
Pipelines adaptés aux structures familiales

→ couverture progressive de chaque structure (en cours)

Rapport de synthèse automatisé interactif

~ 100/120 dossiers par semaine

Interface de tri de variants SeqOne



MR-2101612 < Retour

MR-2101612
Version 5.0 beta

IMRO021
Malformations cérébrales

Informations cliniques

Conclusion de la RCP d'amont
FAVORABLE

Hypothèse diagnostique
Non renseigné

Symptômes renseignés
D012758: Neurodevelopmental delay
O001339: Lisencéphalie
O200134: Encéphalopathie épileptique
O002069: Crises convulsives généralisées tonico-cloniques
O001302: Pachygyrie
O000456: Strabisme

Centre prescripteur
CHU Dijon

Variations retenues

Vous avez sélectionné les 1 variations suivantes :

Variations ponctuelles et petites insertions/délétions

Inter.	Gène/Panel	Allèle & VAF / Signification familiale	HOVSg / HOVSg: Changement protéique	Fréquence ClinGen v3 / Fréquence cohorte	Scores de prédiction	ClinVar / OMIM
P	PAFAM1B1 View	0.58	chr17:g.3967136C>T c.337C>T + s.511 p.Arg131Ter	--	CADD 58.00 GERP++ 2.8 PPH2 -- SpliceAI --	Pathogène 2 maladies OMIM

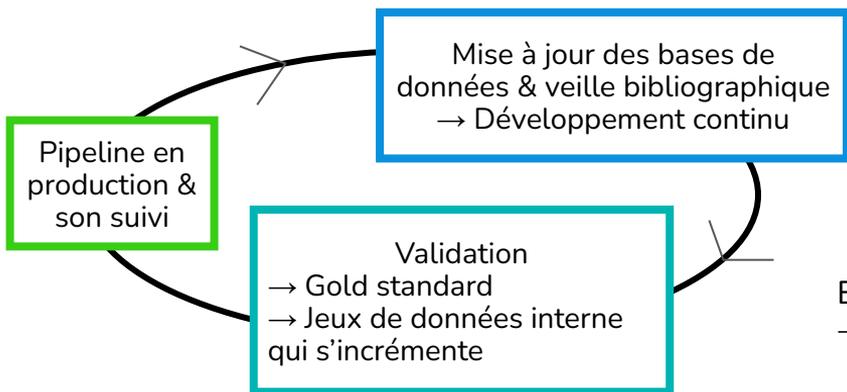
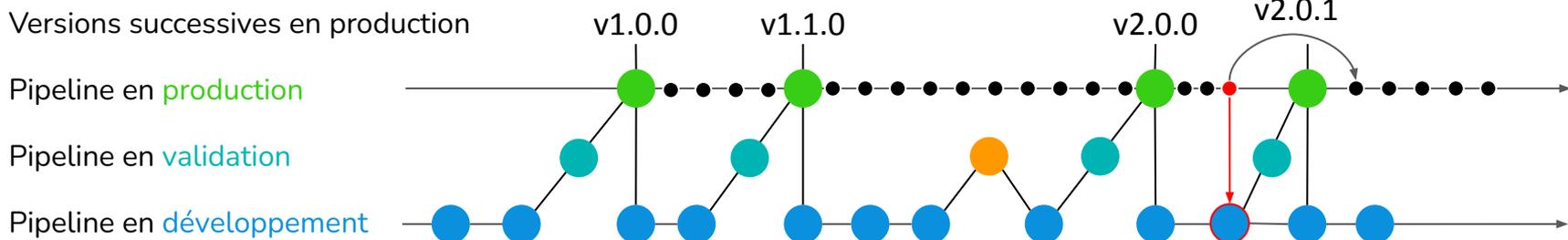
Conclusion du dossier

VCF auragen annoté envoyé (SNV & petites indels)

SeqOne "Auragen"

Amélioration continue des rapports

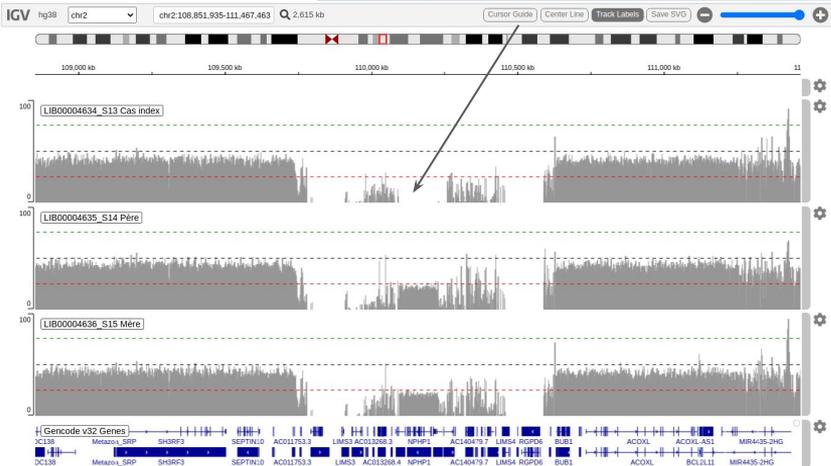
- Analyse d'un dossier en succès
- Analyse d'un dossier en erreur



Gestion des anomalies informatique / bioinformatique **critiques**
 → délai à respecter imposé par les contraintes de la routine

Evolution des pipelines
 → amélioration continue

Délétion homozygote, héritée des parents hétérozygotes chr2 : dossier MR-2000696 → région de faible mappabilité



ABSENT dans les données brutes:

CNVnator recherche une rupture de profondeur

→ en amont pas de lecture (région de faible mappabilité)

→ del homo pas de lecture

→ en aval pas de lecture (région de faible mappabilité)

= pas de rupture

Vu par interpréteurs via les contrôles qualité ciblés

Cible *codant* à moins de 10 lectures

Les 6 gènes de la filière ayant une couverture < 10 lectures pour une partie des régions codantes +/-50 bases.

Gène	Pourcentage de la cible <i>codant</i> des gènes à 20 lectures	Pourcentage de la cible <i>codant</i> des gènes à 10 lectures
HGSNAT	97.43	99.27
INPP5E	98.56	98.87
NPHP1	0.00	0.00
OPNTILW	66.19	71.68
OPNTIMW	6.08	35.28
RPGR	92.01	99.07

Pour ces gènes, nous listons ici tous les défauts impactant la séquence *codante*. Concernant les régions introniques, nous listons les défauts à moins de 50 bases d'un exon, avec une profondeur entre 0 à 8 lectures sur un intervalle de taille supérieur à 5 bases sont détaillés. Une absence de tableau est possible.

Intervalles de la cible "*codant*" à moins de 10 lectures et non récurrents

Gène	Nombre de bases < 10 lectures	Nombre de lectures sur l'intervalle (max)	Distance à l'exon (bases)	Intervalle du défaut
NPHP1	137266	[0,9]	0	chr2:110092811-110230076

→ bilans de couverture des gènes d'intérêt

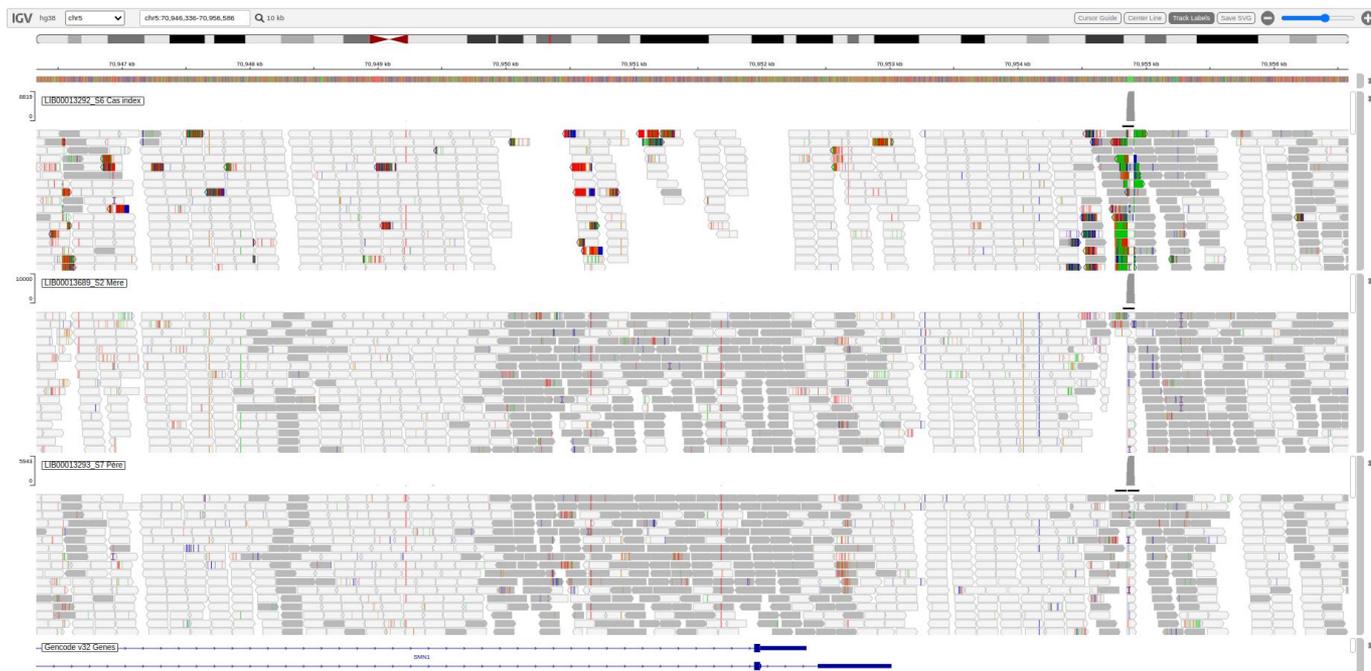
Premier "contrôle positif" pour mettre en place la nouvelle approche de détection des délétions homozygotes (cf. slides suivantes)

→ Mise en place d'un **appel des "délétions homo"** dans les régions à faible mappabilité
 Reads blancs : reads de faible mappabilité **SMN1**

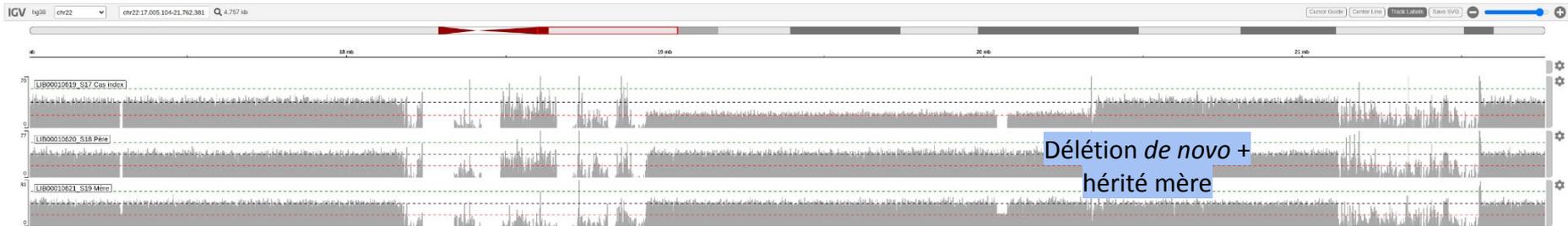
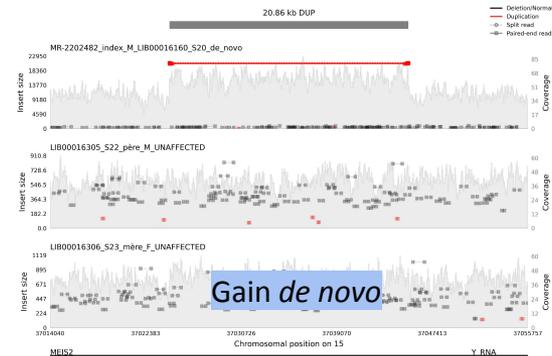
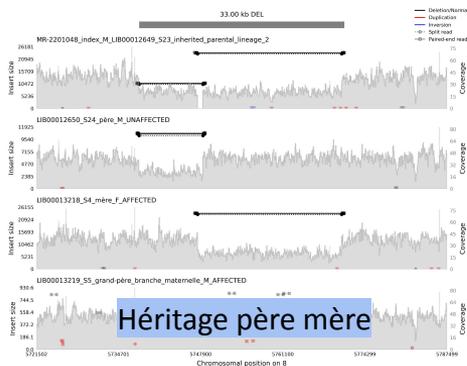
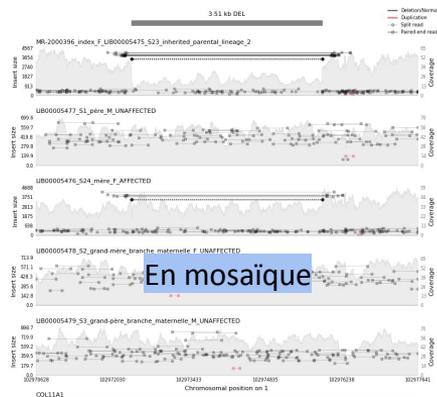
CI

Père

Mère



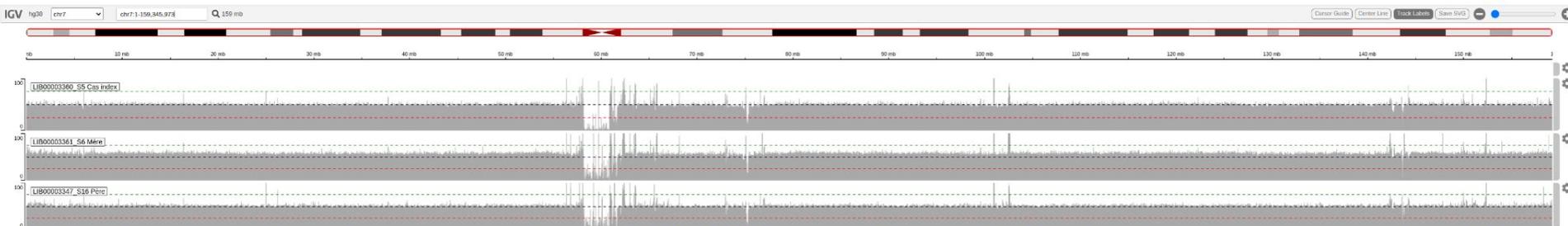
CI
Père
Mère



Régions d'homozygotie

Les régions d'homozygotie d'une taille supérieure à 2Mb sont représentées dans le tableau ci-dessous.

Intervalle	Taille (en Mb)
chr7:6840356-29640320	22.80
chr7:35946191-56339116	20.39
chr7:65831419-72962554	7.13
chr7:77058056-96745610	19.69



Régions d'homozygotie

Les régions d'homozygotie d'une taille supérieure à 2Mb sont représentées dans le tableau ci-dessous.

Intervalle	Taille (en Mb)
chr7:6840356-29640320	22.80
chr7:35946191-56339116	20.39
chr7:65831419-72962554	7.13
chr7:77058056-96745610	19.69

B-allèle frequency

Exploitation de ~1 million de polymorphismes fréquents sur le génome

Critères requis : 20 lectures

BAF entre 0 (ref) et 1 (homozygote), centrée sur 0.5 (hétérozygote)

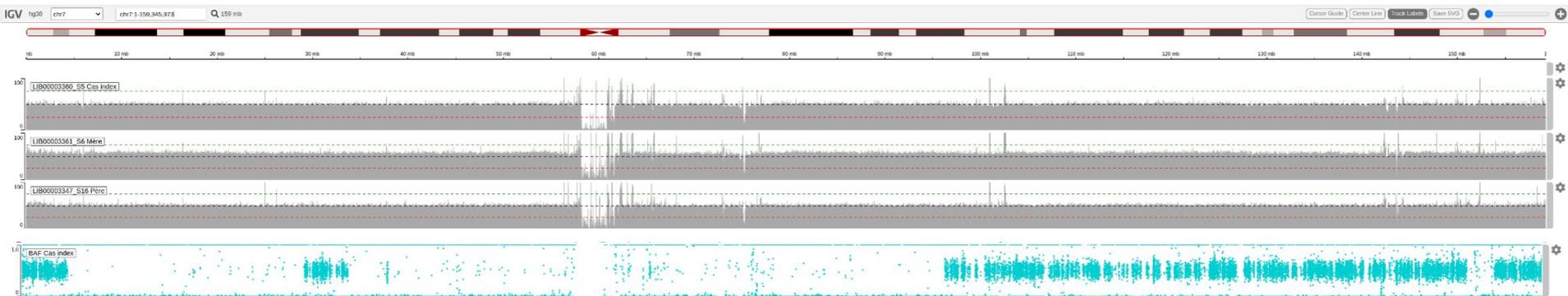
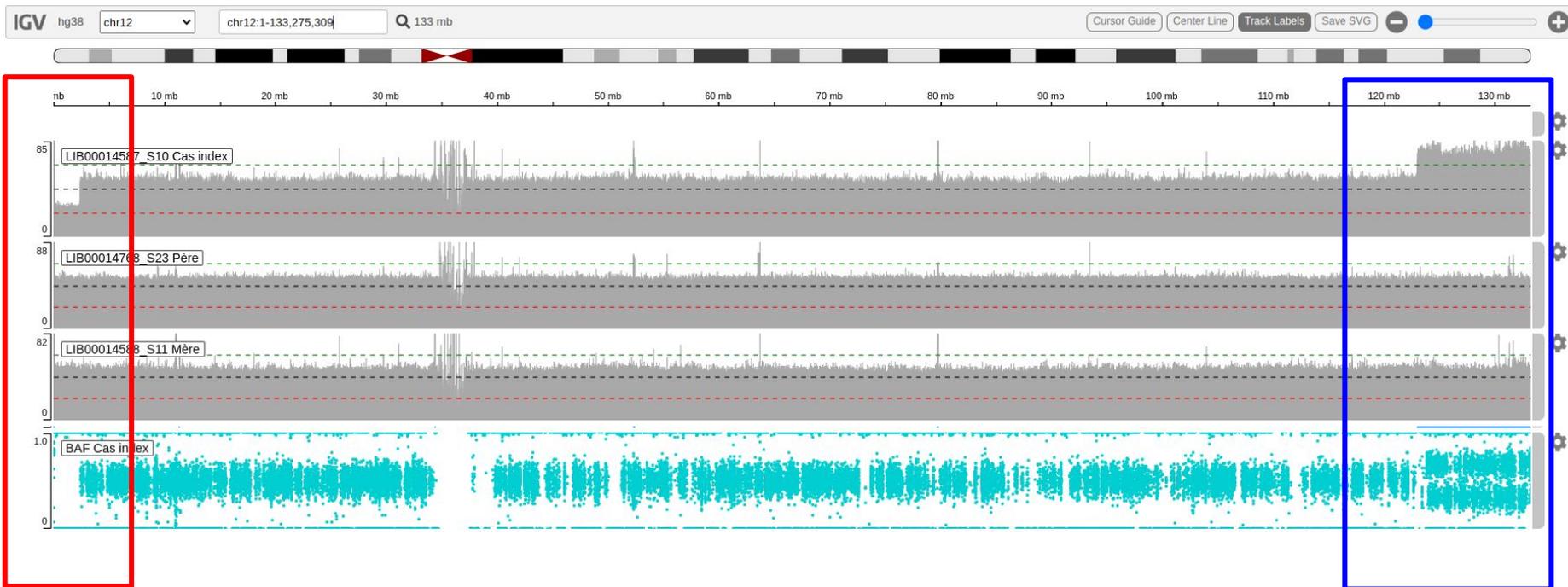


Illustration SNP array like pour le cas index - CNV



BAF des gains et pertes “classiques”

Illustration SNP array like pour le cas index - mosaïque

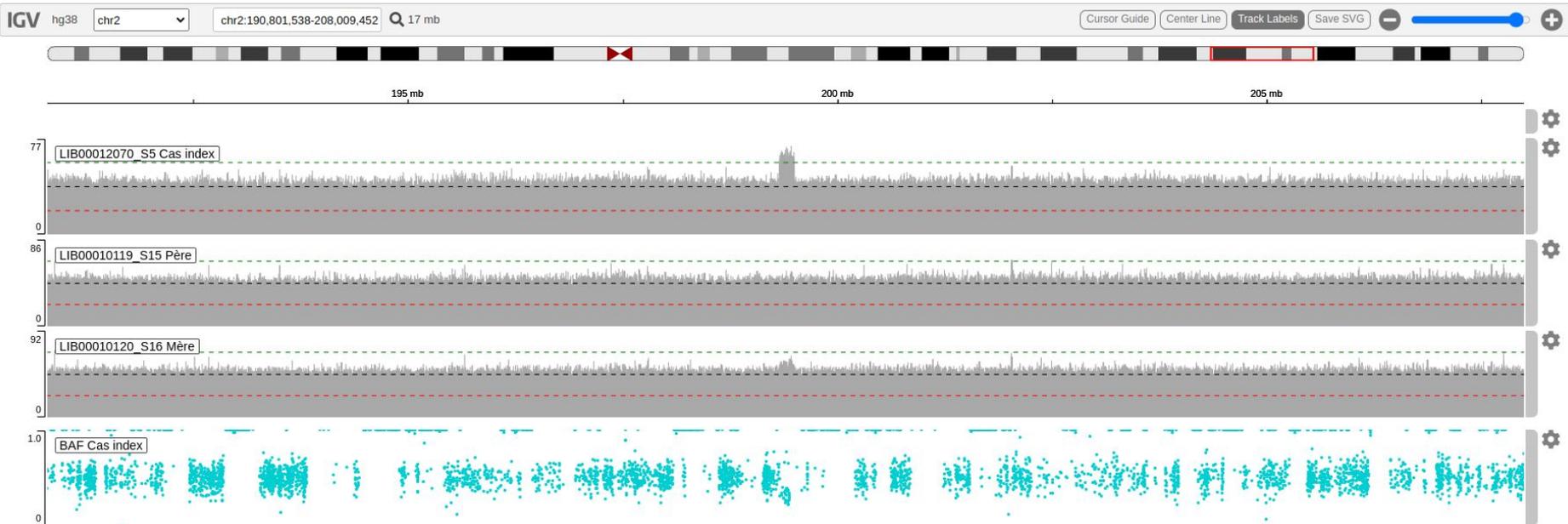


Illustration SNP array like Pour aller plus loin...

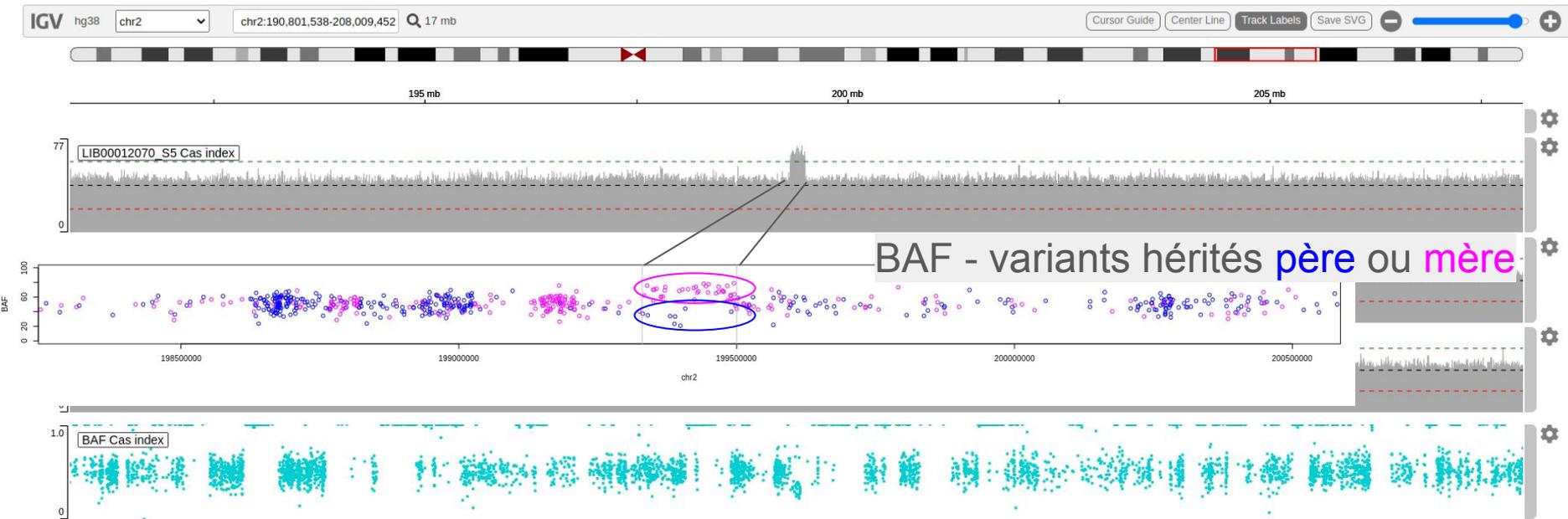


Illustration SNP array like Pour aller plus loin...

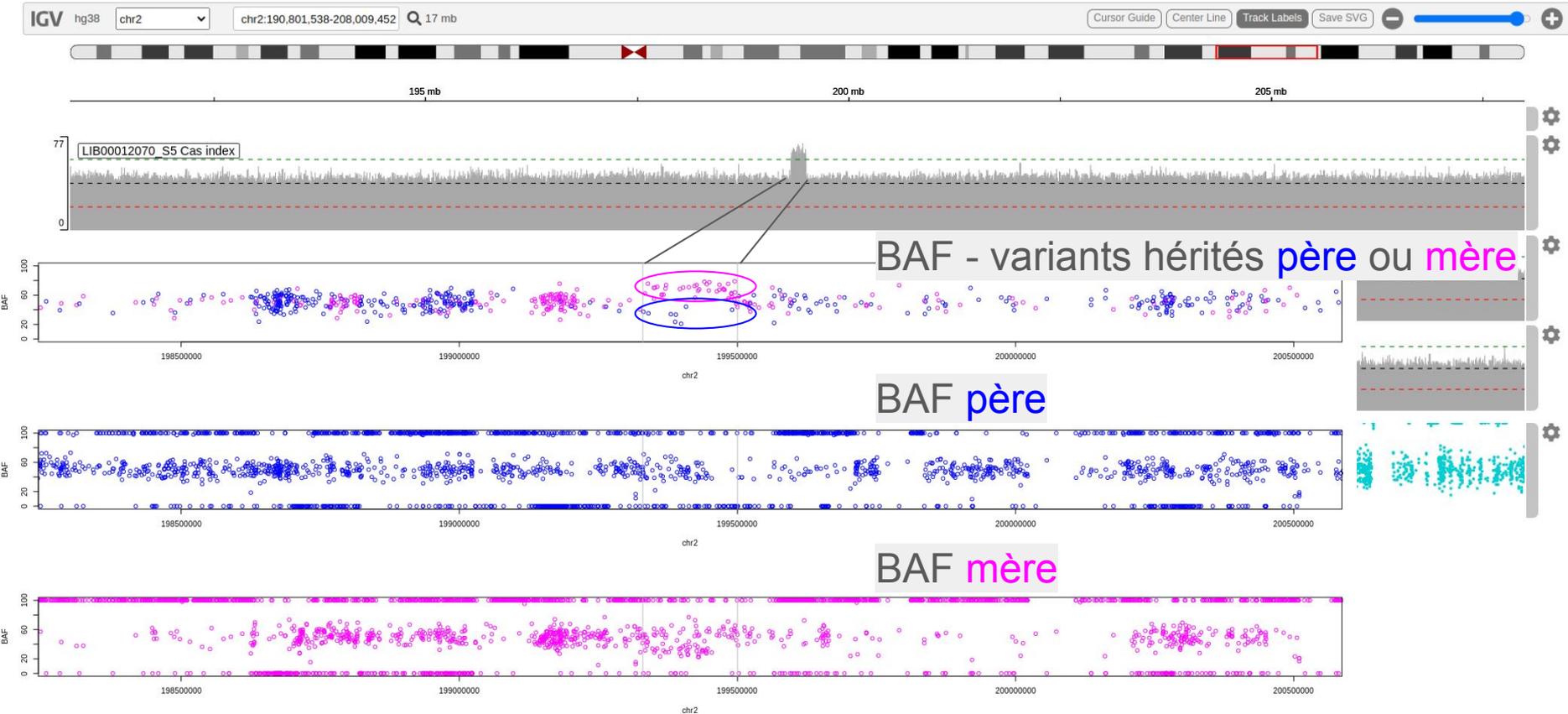
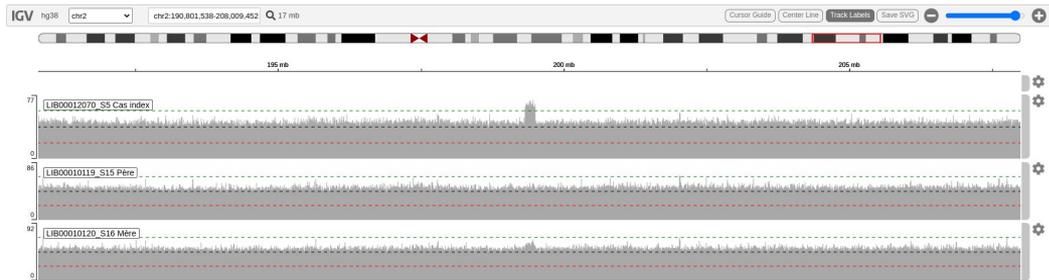
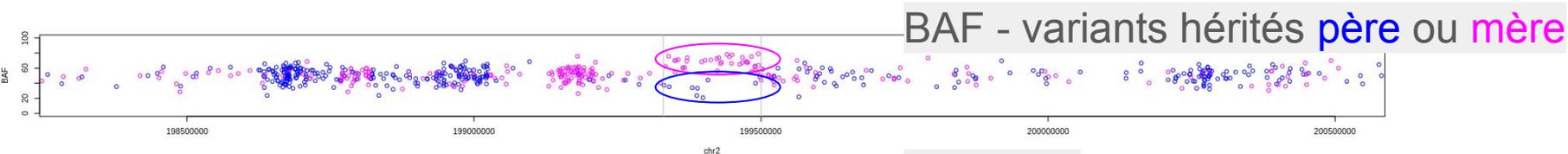


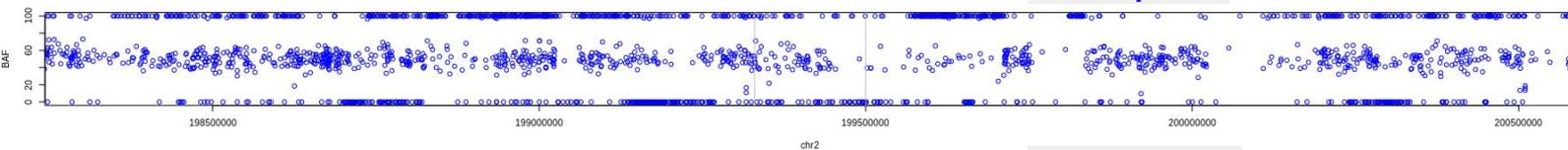
Illustration SNP array like Pour aller plus loin...



Visible sur le BAM ✓
Renforcé par profils BAF ✓
Confirmé par MLPA ✓



BAF père



BAF mère

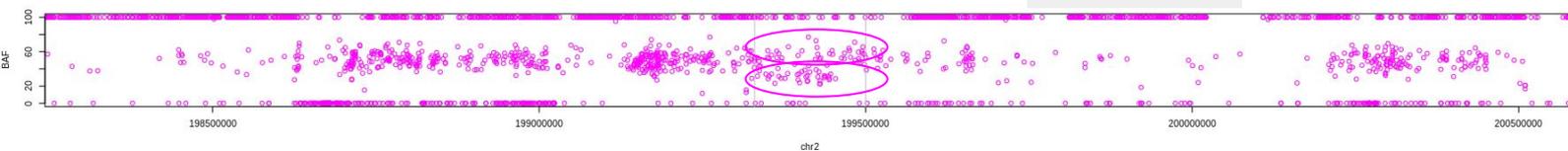
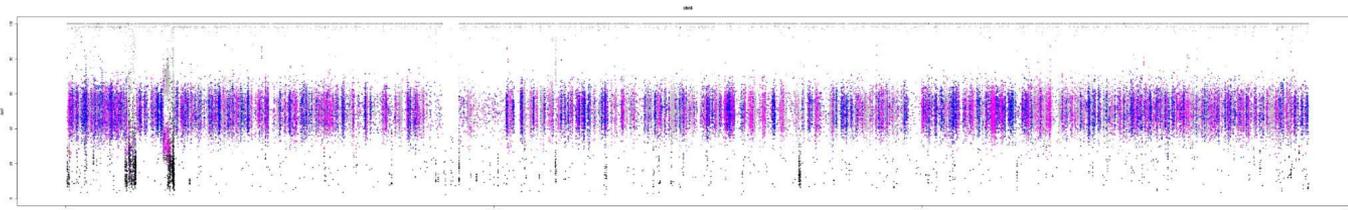
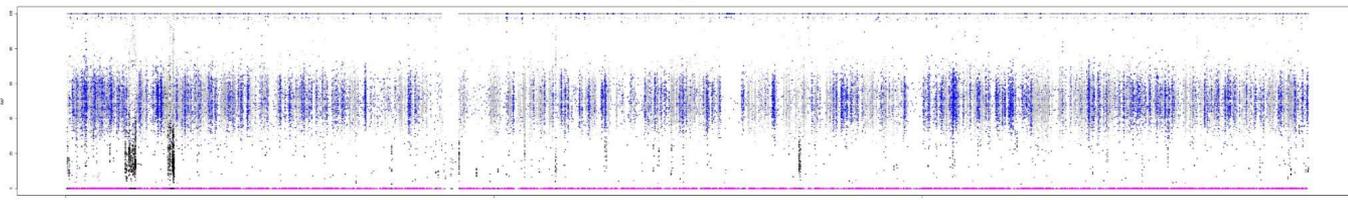


Illustration SNP array like & héritage père mère - profil standard

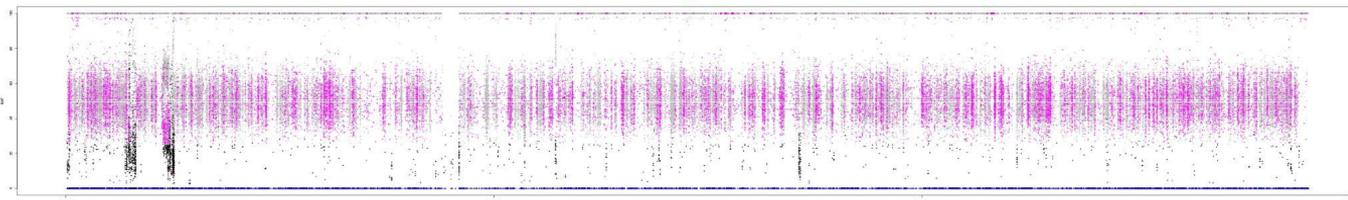
CI



Père



Mère

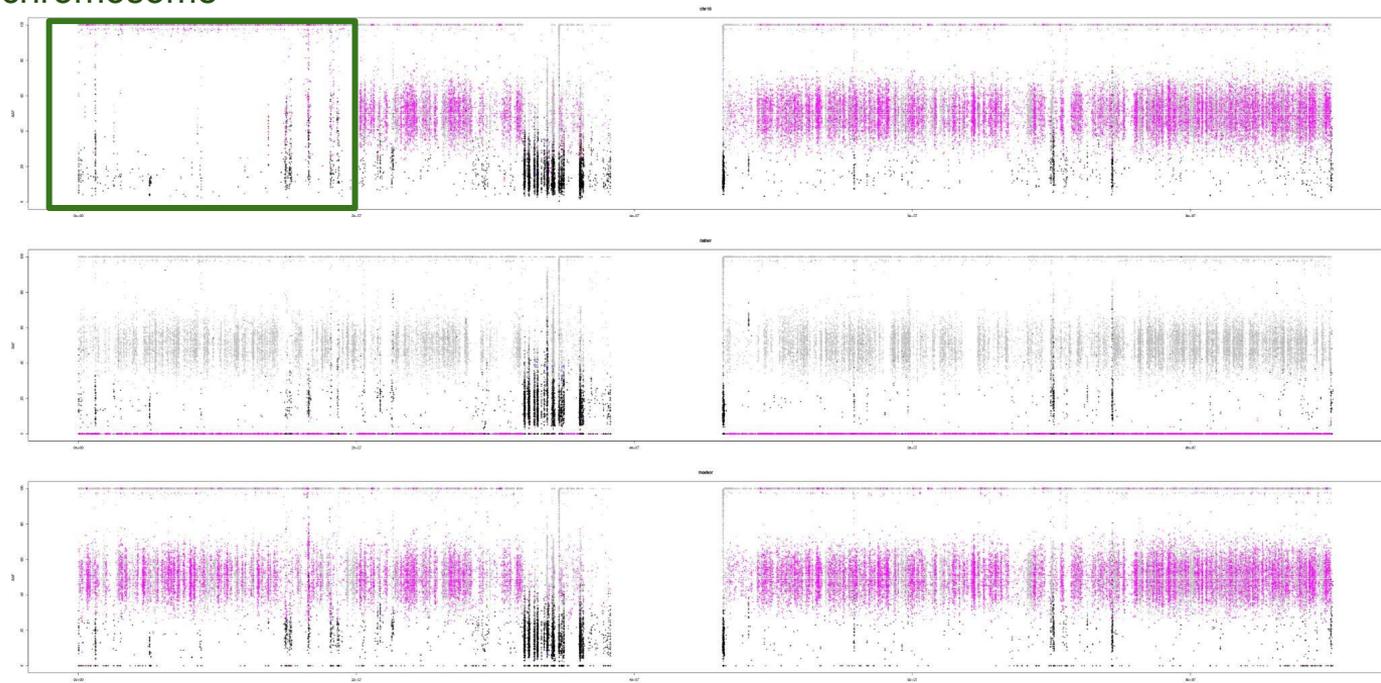


- SNP hérité du père

- SNP hérité de la mère

Illustration SNP array like & héritage père mère - profil UPD

Région d'isodisomie : Origine des
SNP : un seul chromosome
maternel

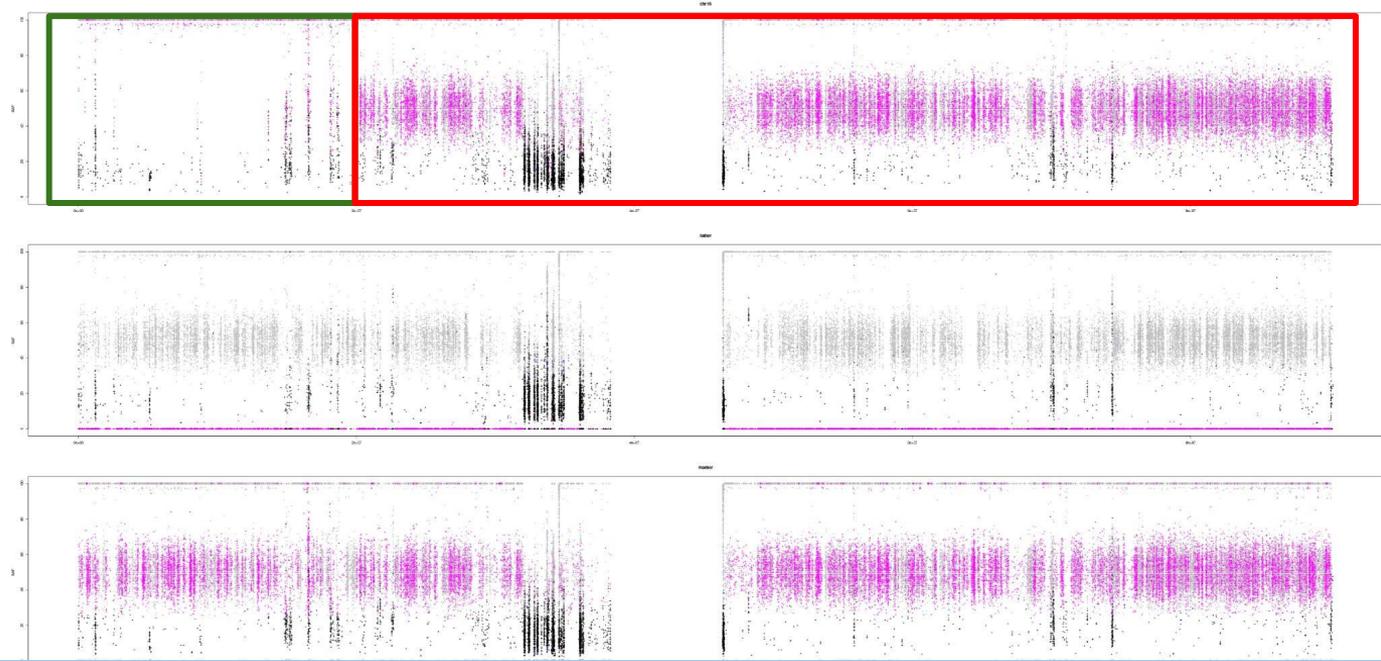
A small icon consisting of a square with diagonal yellow and black stripes, followed by a horizontal line.

WORK IN PROGRESS

Illustration SNP array like & héritage père mère - profil UPD

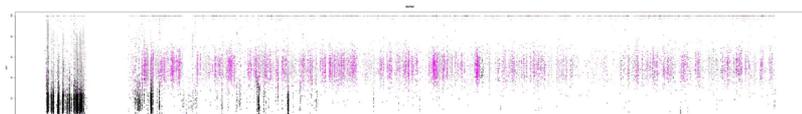
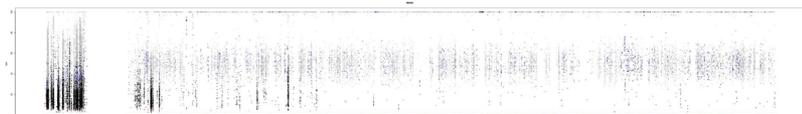
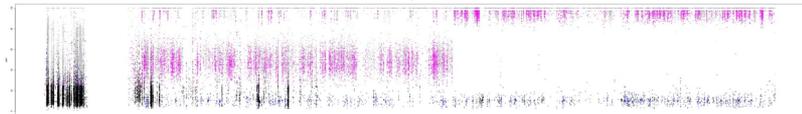
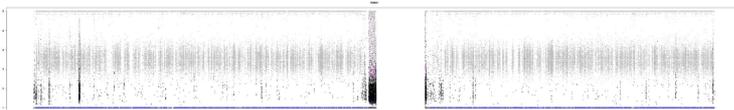
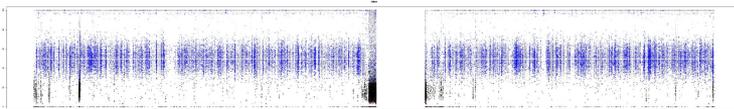
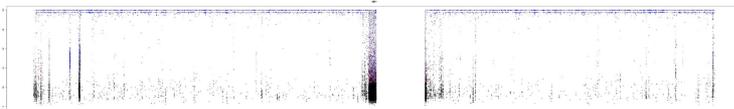
Région d'isodisomie : Origine des SNP : un seul chromosome maternel

Région d'hétérodisomie : Origine des SNP : les deux chromosomes **maternels**



=> Disomie uniparentale de tout le chr16 : UPD(16) Avec isodisomie 16pter et hétérodisomie sur le reste du chr

Illustration SNP array like Quoi prioriser ?



WORK IN PROGRESS

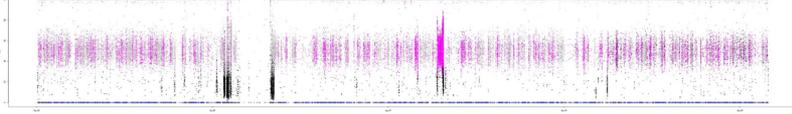
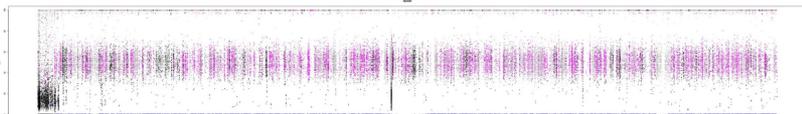
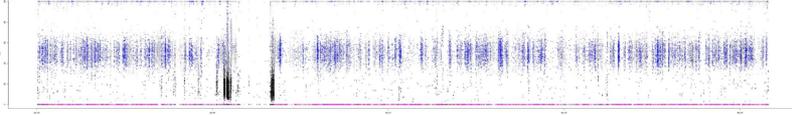
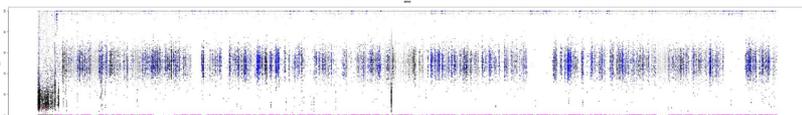
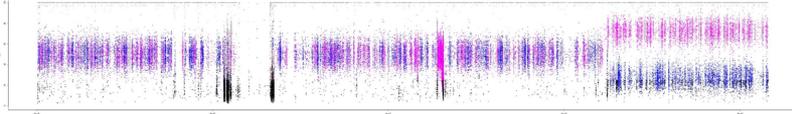
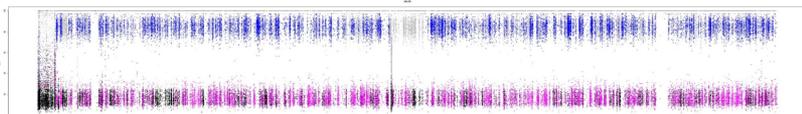
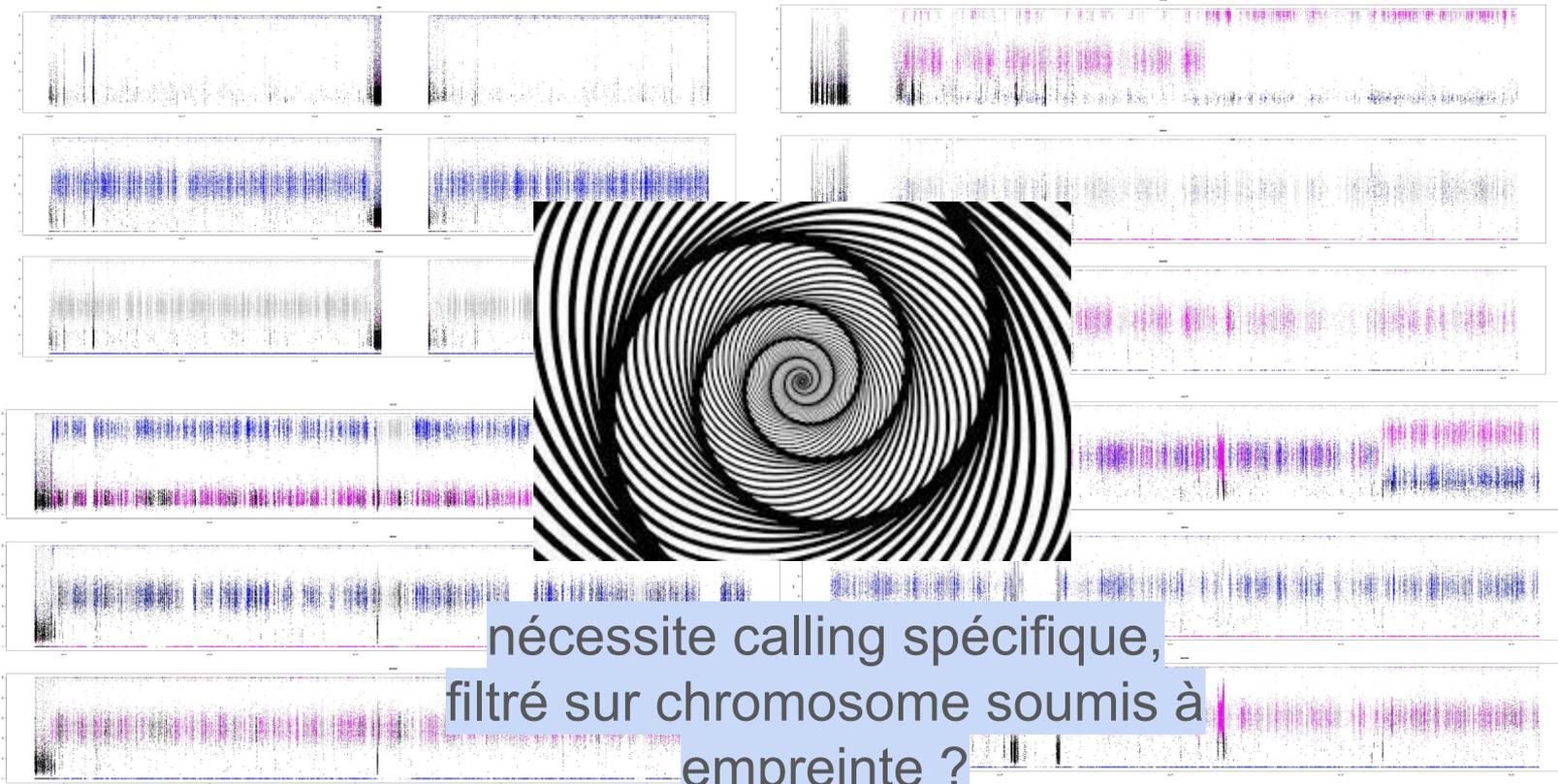


Illustration SNP array like Quoi prioriser ?



WORK IN PROGRESS

nécessite calling spécifique,
filtré sur chromosome soumis à
empreinte ?

Mise en place wet & dry Mise en place dry & bio
Importance MAJEURE pour permettre les améliorations

Livraisons pour interprétation dans le flux
→ automatisation de la prod

Confiance!

La suite ?

Auragen : todo liste des améliorations encore grande :)
⇒ Objectif diag!

Collecteur et Analyseur de Données
→ Projets soumis pour les améliorations de demain



WORK IN PROGRESS





Merci!