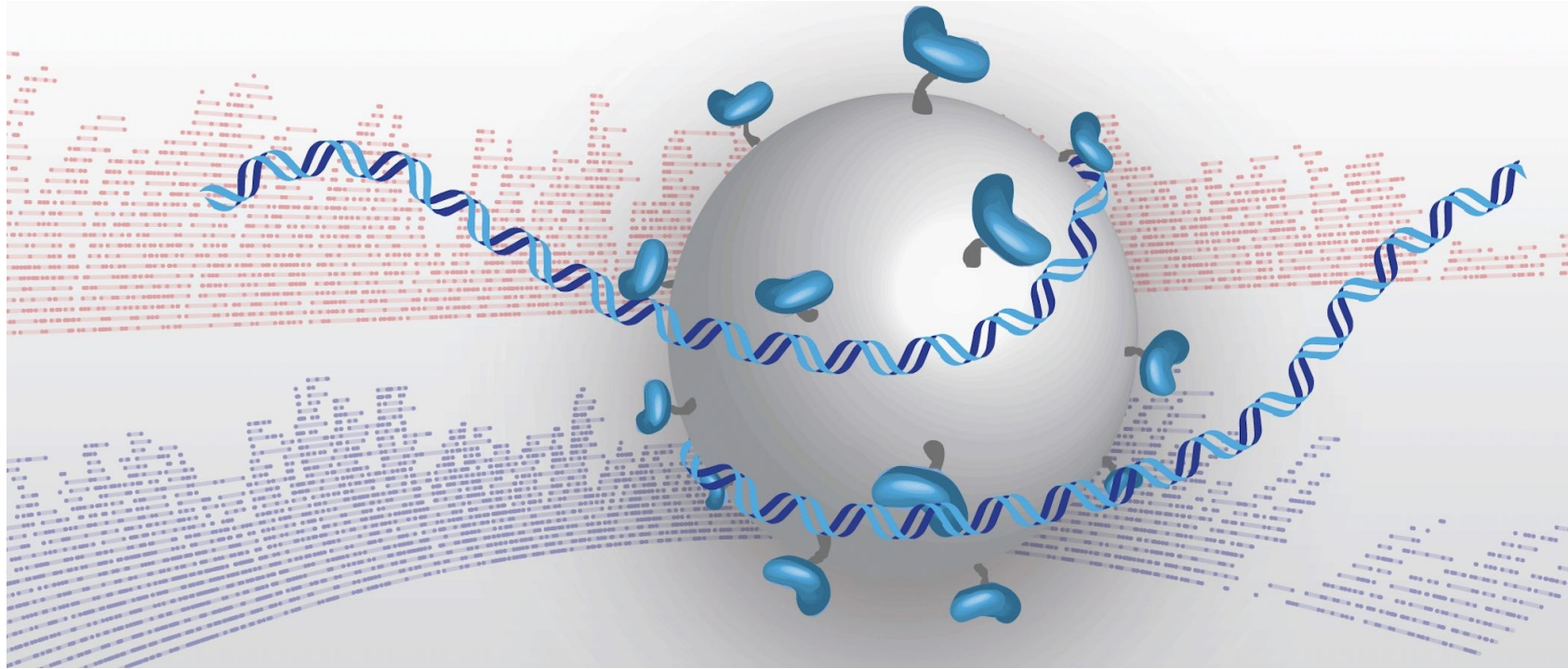


# Haplotagging: état de l'art

## Phasage de génomes à bas coût pour les études populationnelles



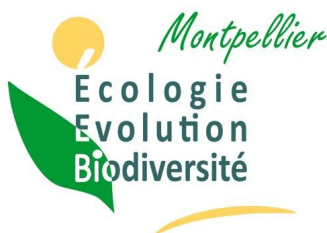


# Projet **DevOCGen** : Développement et applications de nouveaux outils pour la gestion et la conservation des populations naturelles à partir de données génomiques

## Porteurs

Simon BOITARD  
Raphaël LEBLOIS

Plateforme **GenSeq**  
**UAR MEEB**



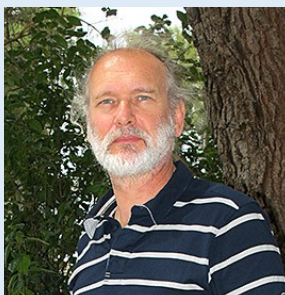
**Un projet multi-partenaires Occitanie**





# Projet **DevOCGen** - Plateforme **GenSeq**

**Axe méthodologique :** Développement de l'haplotagging  
*Librairies à « lectures liées » permettant de reconstituer  
l'information de la phase haplotypique*



**Érick Desmarais**



**Frédérique Cerqueira**



**Anais Bordes**



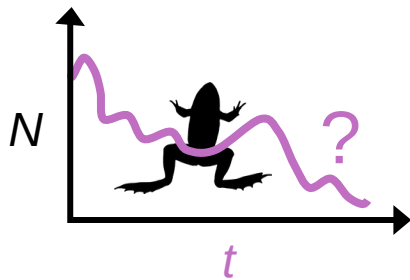
**Cathy Liautard-Haag**



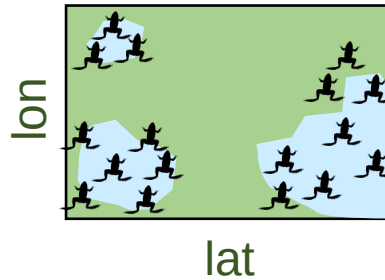
**Valentina Neglia**



# Plusieurs dimensions à renseigner : *le temps, l'espace, le génome*



**TEMPS**



**ESPACE**

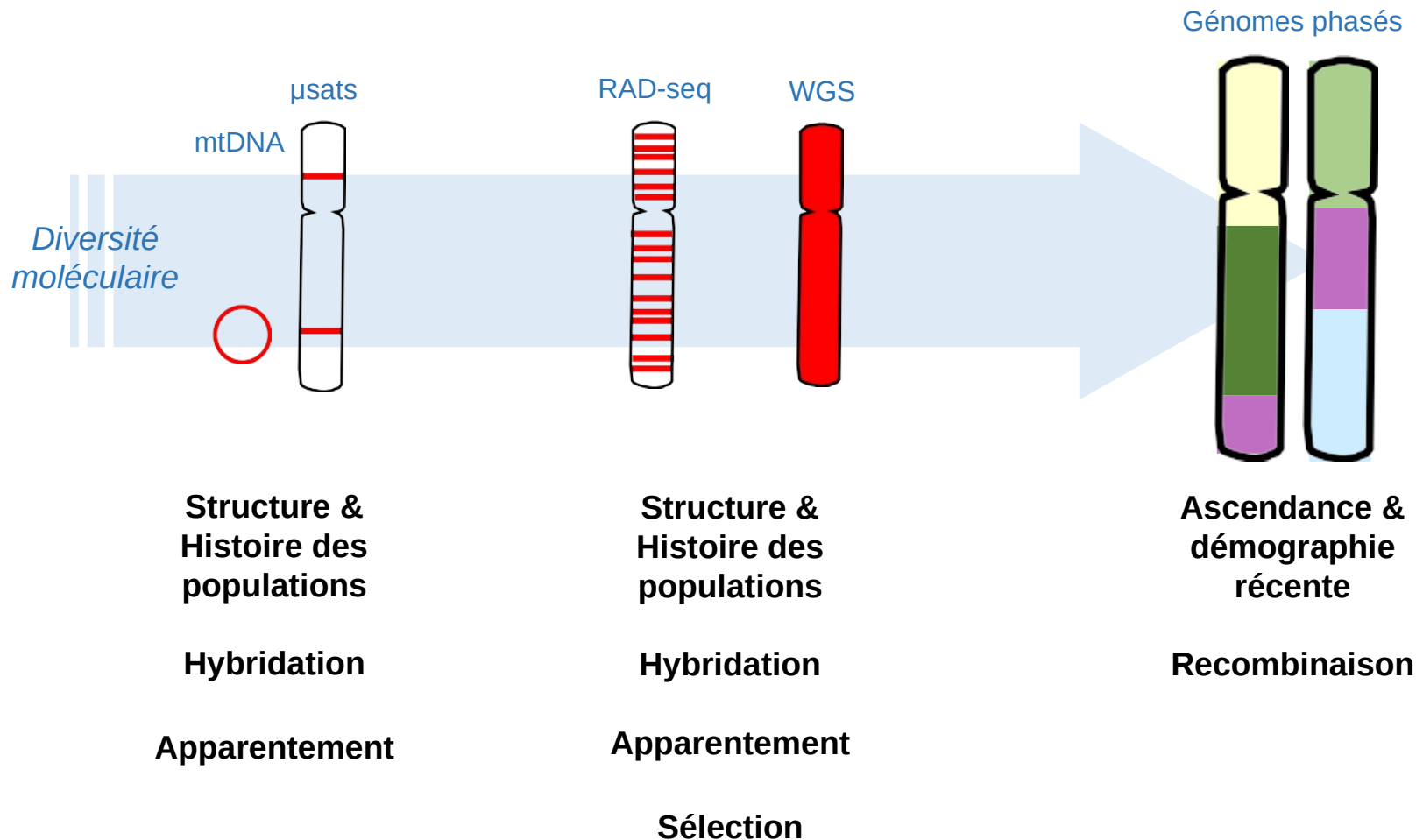


**GENOME**

## Défis :

- suivis temporels rares, partiels, ou inexistants
- respecter la démarche « **éviter-réduire**-compenser »
- maximiser l'information extraite de chaque échantillon

# Quelle information peut-on espérer extraire des données génomiques ?



# Utilités de l'information haplotypique pour l'étude de la biodiversité

## ***POUR FAIRE QUOI?***

- Étude de la **recombinaison**
- Détection de la **sélection**
- Reconstruction de l'**histoire évolutive**
- Inférence de l'**ascendance** et de la **démographie récente**
- Reconstruction du **Graphe de Recombinaison Ancestrale (ARG)**
- Étude des **variants structuraux**
- **Assemblage** de génomes et de pangénomes

# Utilités de l'information haplotypique pour l'étude de la biodiversité

## COMMENT PHASER LES GENOMES ?

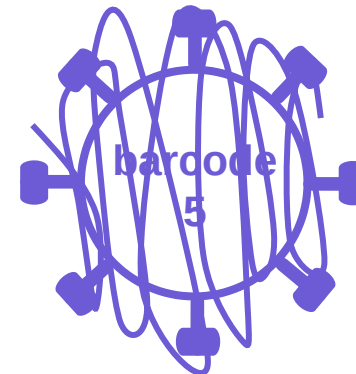
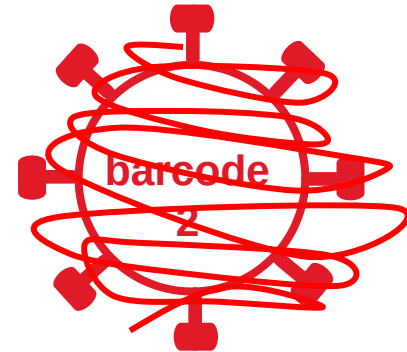
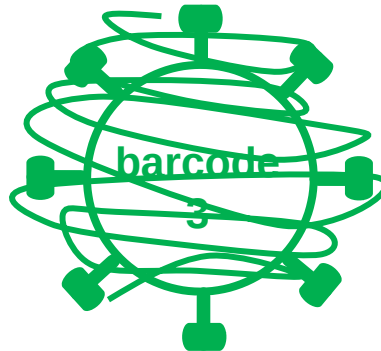
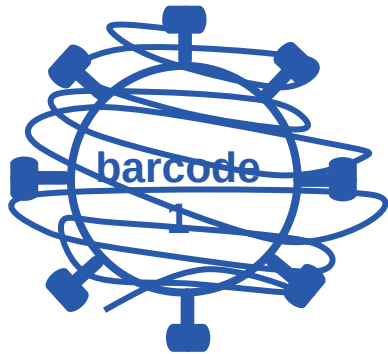
- **Séquençage short-read (*Illumina*):**
  - + économique, taux d'erreur faible
  - perte de l'information haplotypique au-delà de ~1kb
- **Séquençage long-read (*PacBio*, *ONT*):**
  - + phase physique préservée (limitée par qualité de l'ADN)
  - coûts mal adaptés aux études populationnelles
- **Séquençage linked-read (Haplotagging):**
  - + combine les deux avantages
  - coûteux et débit limités

Linked-Reads



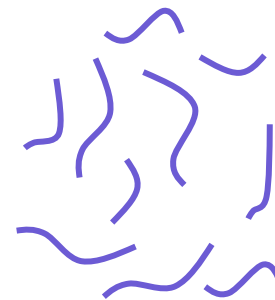
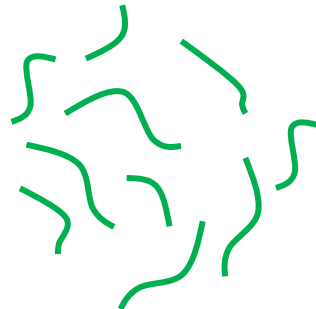
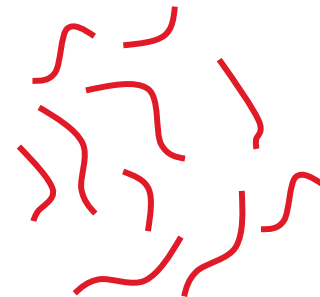
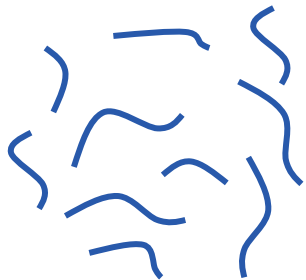
# L'haplotagging

**LINKED-READ QU'EST-CE QUE C'EST?**

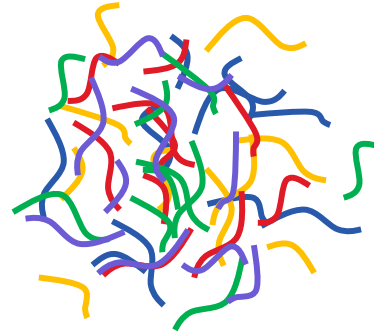




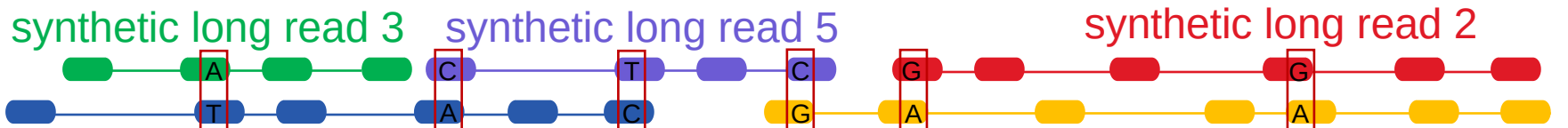
# Un barcode/index par molécule



# Reconstruction bioinformatique des reads longs synthétiques



séquençage short-read  
et mapping



génom

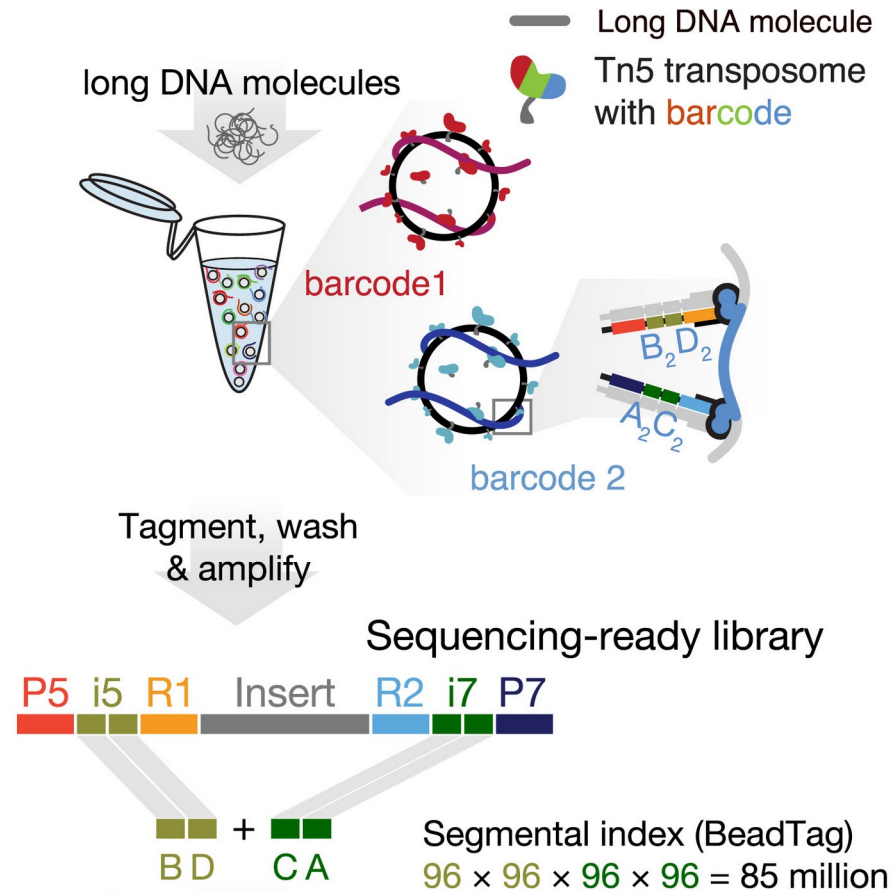


Phasage physique

haplotype maternel  
ACTCGG

haplotype paternel  
TACGAA

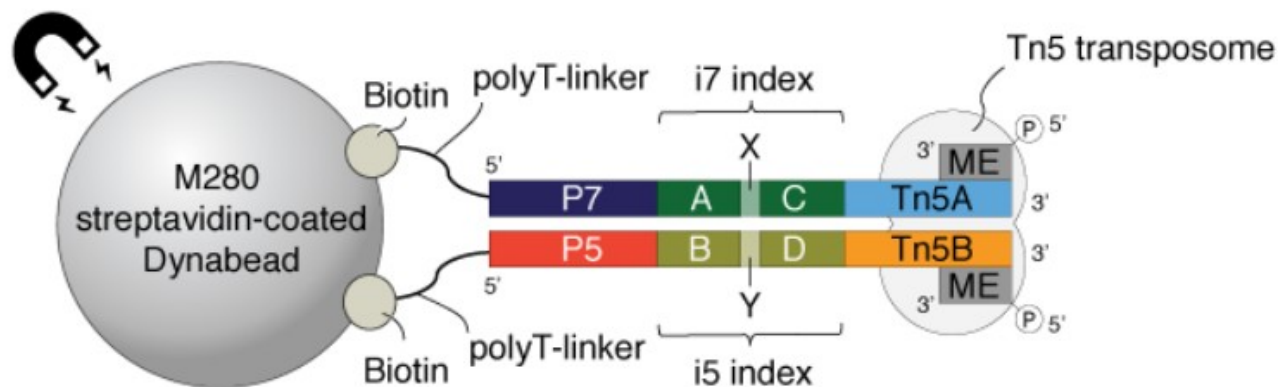
# Récap



# L'haplotagging

## COMMENT CA MARCHE?

- **Combine fragmentation et marquage** de l'ADN (transfert d'adaptateurs indexés) en une seule réaction rapide de tagmentation
- Utilise la **transposase Tn5 liée à des billes**
- **Chaque bille contient** des complexes Tn5-adaptateur avec **une combinaison unique d'index** (modules A-C sur i7, B-D sur i5)



# Un jeu de $96^4$ billes différentes !



# L'haplotagging

## **AVANTAGES**

- Reconstitution de **reads longs synthétiques**
- **Phasage de blocks mégabasiqes** à l'échelle individuelle
- **Réaction enzymatique rapide** avec PCR
- **Bas coût**, moins cher que les librairies Illumina commerciales
- Applicable chez une **large gamme d'organismes**



# Les étapes

1. Extraction/quali/quantité d'ADN de haut poids moléculaire
2. Construction des libraires d'haplotagging et séquençage
3. **Analyses bio-informatiques** pour le phasage des génomes

# Matériel = HMW DNA

## QUALITÉ DE L'ADN EST CRUCIALE POUR LA RÉUSSITE DU PHASAGE

### Échantillonnage :

- Optimiser les conditions échantillonnage et de conservation
- Flash freezing Azote liquide, éthanol 96% ou RNAlater à -20/-80°C

### Extraction :

- HMW MagBead (Zymo research)
- Monarch genomic DNA (NEB)
- Nanobind Tissue (PacBio)
- "Phénol-chloroforme

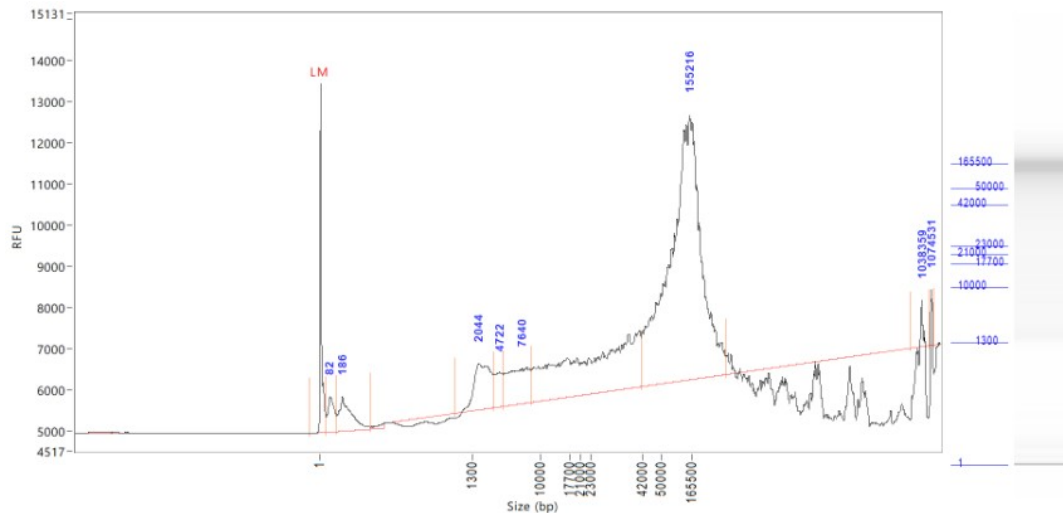
### Élimination des fragments courts - sélection de taille :

- SRE kit (PacBio) pour éliminer le molécules <10 kb (moyen)

# QUALITE - QUANTITE ADN

## Vérification sur Femto Pulse (Agilent):

- Maximiser les fragments >10kpb



**Agilent Femto Pulse**  
**DNA profile**

≥ 50% of DNA ≥ 30 kb

Très très peu d'ADN nécessaire: <1ng ADN/indiv

Préparation des ADN une grosse partie du temps et du budget

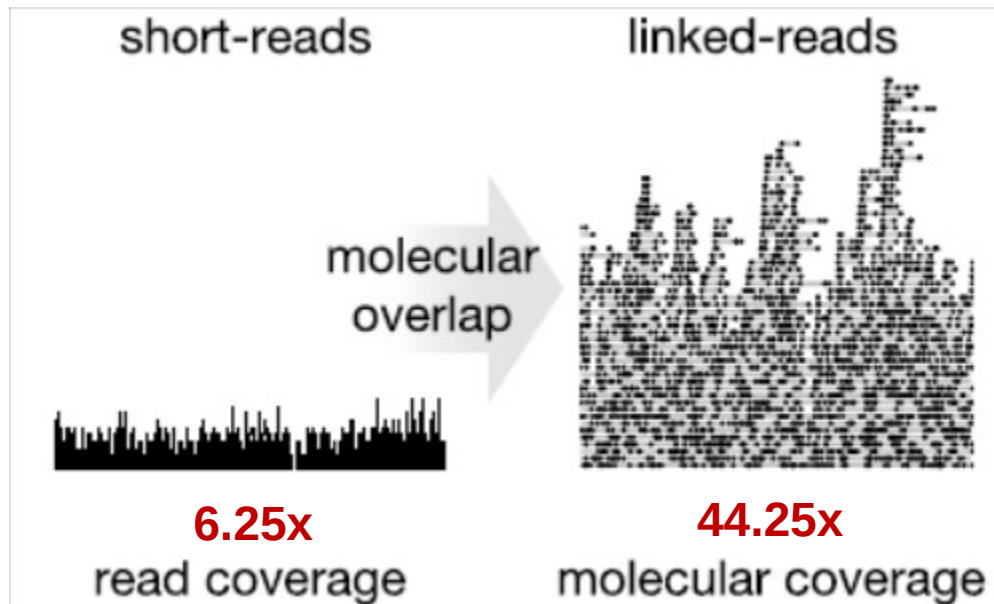
# Notion de couverture moléculaire

Pour un génome de 1Gb - ligne S4

En utilisant des molécules d'ADN de 50kb en entrée :

*Les molécules haplotaggées devraient couvrir  $50^3 \times 885000 = 44.2\text{Gb}$*

**>> Couverture moléculaire théorique de 44.2x par individu**



Objectif de couverture moléculaire de 50X

Avec couverture de reads de 5x et des ADN de taille moyenne de 50kb:

En moyenne 10 reads par molécule

# Design de construction des librairies

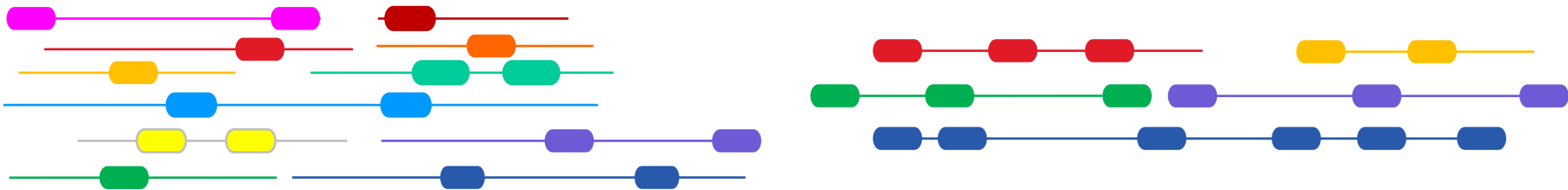
Choix FlowCell : xx million reads



Taille génome: Couverture moyenne par position génome (5x)



Nombre et taille moyenne des molécules ADN génomique



**IMPORTANCE DE LA TAILLE DES MOLÉCULES**  
**IMPORTANCE DE LA QUANTITÉ D'ADN UTILISÉ /**  
**DIMINUER LE NOMBRE DE BILLES SÉQUENCÉES**

# Design de construction des librairies

## Input Parameters

Number of individuals multiplexed:

Choose a flowcell type:

Choose read length:

Haploid genome size (Mb):

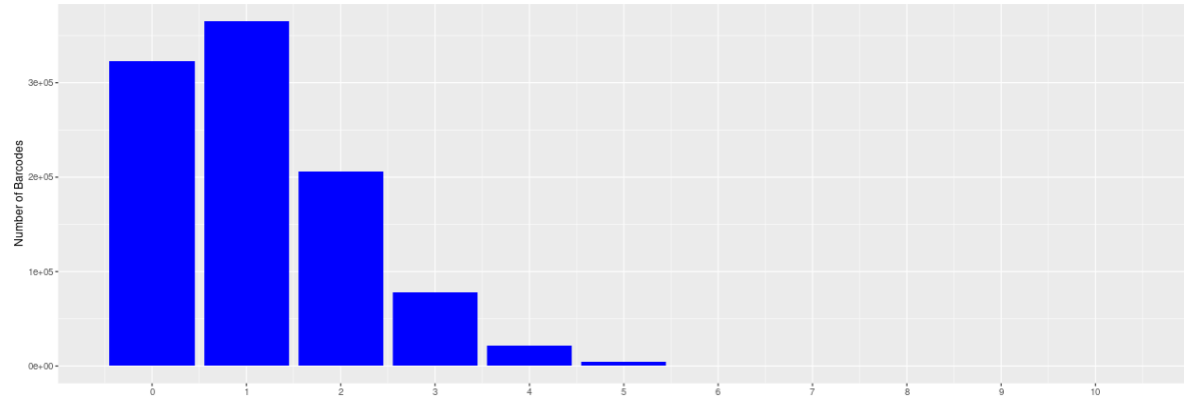
Median size of HMW DNA molecules (bp):

Number of beads used:

Number of HMW DNA molecules used:

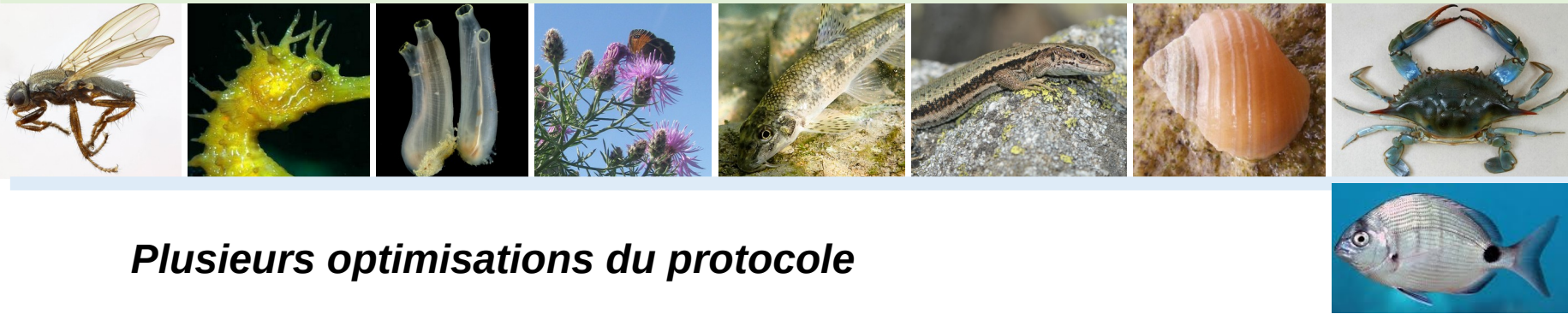
Fraction of the tagmentation volume used:

## Nombre de copies de chaque barcode :

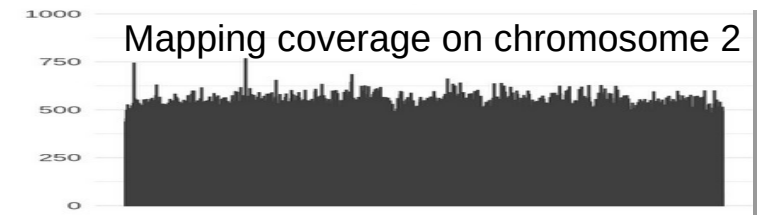
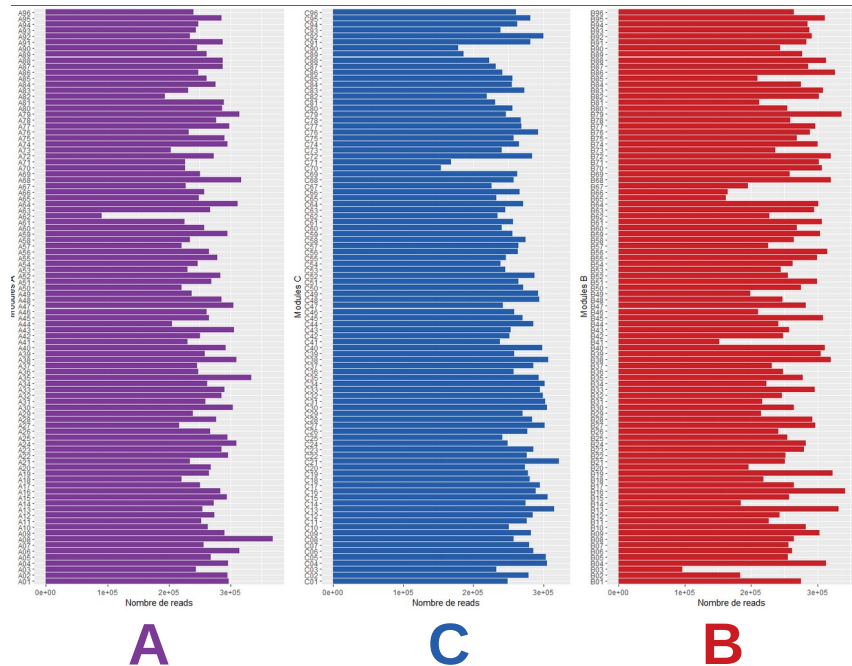




# Espèces déjà analysées en haplotagging à *GenSeq*



*Plusieurs optimisations du protocole*



# Les coûts

## 1. Extraction d'ADN de haut poids moléculaire (+ *Femto*)

Extraction : ~10€/échant

Femto: ~10€/échant

SRE: ~5€/échant

## 2. Construction des libraires d'haplotagging

*~700€ par plaque de 96 échantillon*

## 3. Séquençage : ~ 3200€ pour un génome de 1Gb

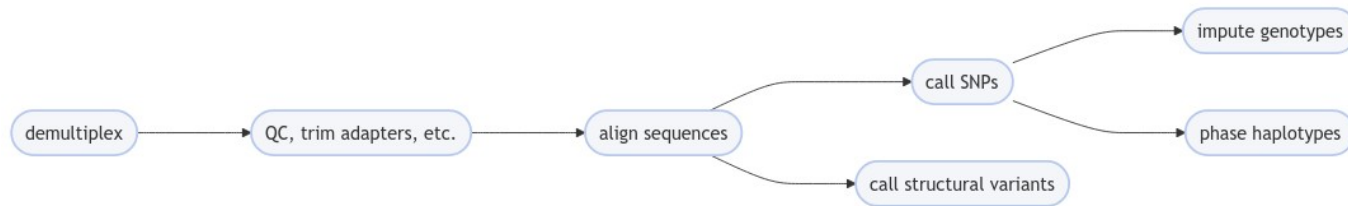
## 4. Bioinformatique (stockage + analyse)

# Pipeline bio-informatique

## Harpy: a pipeline for processing haplotagging linked-read data

Pavel V. Dimens<sup>1,\*</sup>, Ryan P. Franckowiak<sup>1</sup>, Azwad Iqbal<sup>1</sup>, Jennifer K. Grenier<sup>2</sup>, Paul R. Munn<sup>2</sup>, Nina Overgaard Therkildsen<sup>1</sup>

*Bioinformatics Advances*, 2025, **00**, vbaf133  
<https://doi.org/10.1093/bioadv/vbaf133>  
Advance Access Publication Date: 5 June 2025  
**Application Note**



HARPY v3.0

Therkildsen Lab Cornell GIH Source Submit Issue

**Pavel Dimens**

<https://github.com/pdimens/HARPY>

**GETTING STARTED**

- Install
- Input Format
- Linked-Read Data
- Common Options
- Guides
- Resources
- Troubleshooting

**WORKFLOWS**

- Align
- Assembly
- Convert deprecated
- Deconvolve
- Downsample deprecated
- Demultiplex
- Impute
- Metasassembly
- Template
- Phase
- Validate
- QC
- Simulate
- SNP

Powered by RETYPE

### Home



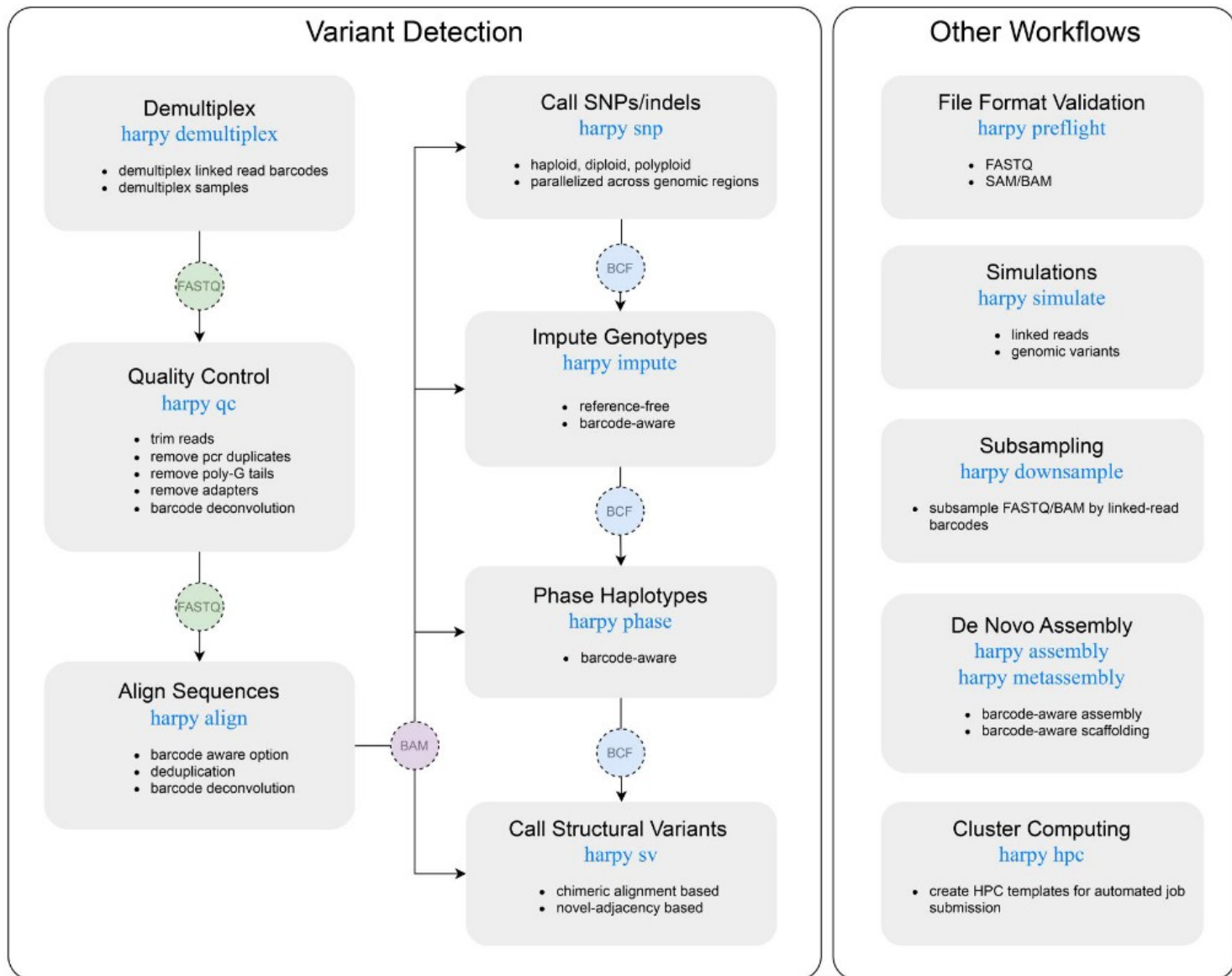
Harpy is a [haplotagging data](#) processing pipeline for Linux-based systems-- at least it was prior to the release of version 2. Now, it can process linked-read data from haplotagging, TELLseq, stLFR, and even regular non-linked WGS data. It uses all the magic of [Snakemake](#) under the hood to handle the workflow decision-making, but as a user, you just interact with it like a normal command-line program. Harpy employs both well known and niche programs to take raw linked-read sequences and process them to become called SNP genotypes (or haplotypes) or large structural variants (inversions, deletions, duplications). Most of the settings are pre-configured and the settings you can modify are done at the command line. Some parts of this documentation will refer to haplotagging specifically as we either forgot to update parts of the documentation or require you (the user) to do a data conversion for some parts of Harpy to work with non-haplotagging linked-read data. As always, feel free to drop an [issue](#) or open a [Discussion](#) on GitHub.

### Harpy Commands

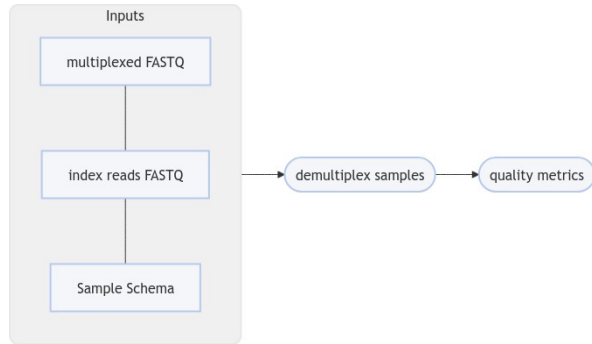
Harpy is modular, meaning you can use different parts of it independent from each other. Need to only align reads? Great! Only want to call variants? Awesome! All modules are called by `harpy <workflow>`. For example, use `harpy align` to align reads.

Command	Description
<code>align</code>	Align sample sequences to a reference genome
<code>assembly</code>	Create a genome assembly from linked-reads

# Pipeline Harpy



# Rapports Harpy : démultiplexage



Samples

96



Min % Valid

75.773



Max % Valid

98.521



Min BX Absent

0



Max BX Absent

0

## General Per-Sample Haplotag Barcode Statistics

Below is a table listing all the samples Harpy processed and their associated haplotag barcode statistics as determined by the reads in the **forward read only**. If for some reason **TotalBarcodes** equals 0, then there may be an issue with the format of your FASTQ headers.

CSV

Sample	TotalReads	TotalBarcodes	ValidBarcodes	ValidPercent	InvalidBarcodes	InvalidPercent
All	All	All	All	All	All	All
Lsaxa_TRA_VPA_010	1856694	1856694	1807087	97.328	49607	2.672
Lsaxa_TRA_VPA_011	8699420	8699420	8518240	97.917	181180	2.083
Lsaxa_TRA_VPA_021	13392423	13392423	13130035	98.041	262388	1.959
Lsaxa_TRA_VPA_034	6180942	6180942	6044892	97.799	136050	2.201

Overall barcode quality per sample



# Rapports Harpy : alignement



Contigs

63

Intervals

50 kb

Average Depth

1.46

Stdev Depth

0.3

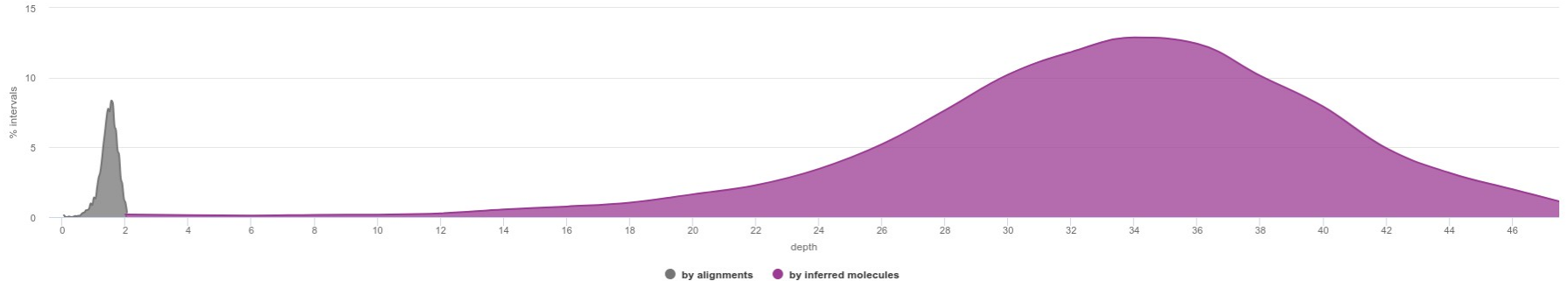
Mol. Average Depth

32.29

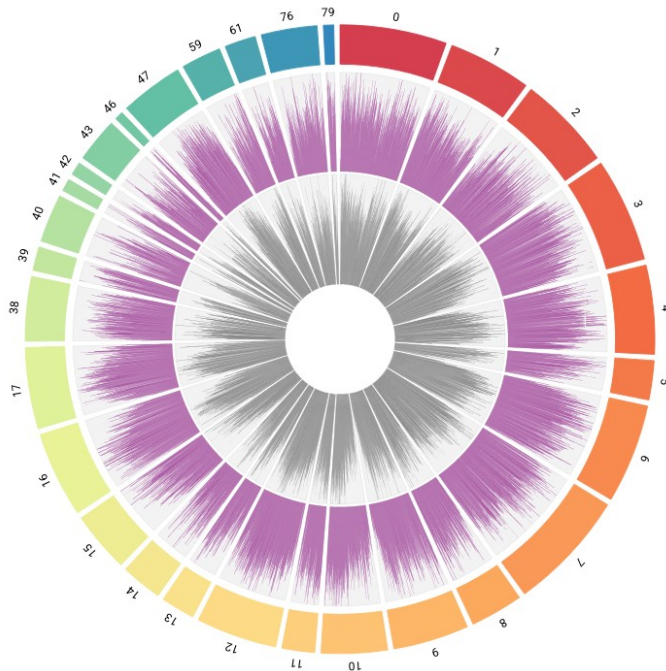
Mol. Stdev Depth

7.08

couverture reads vs molécules synthétiques :

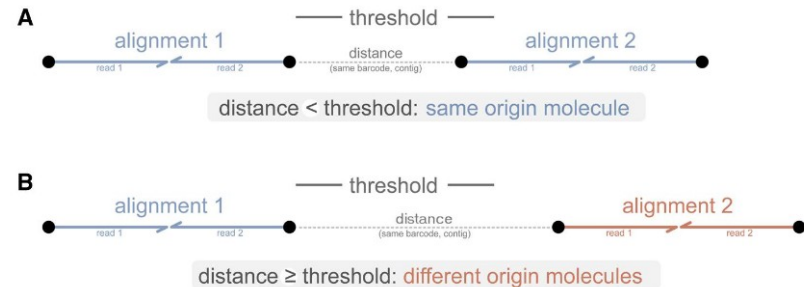


inferred molecules calculate "effective" or "molecule" coverage



Etape de déconvolution :

contig





# Rapports Harpy : alignement



Contigs  
2,751

Unique  
Barcodes  
633,220

Molecule  
Threshold  
200,000

Unique  
Molecules  
1,477,736

Valid BX  
Records  
2,229,336

Invalid BX  
Records  
130,732

Singletons  
(%)  
76.27

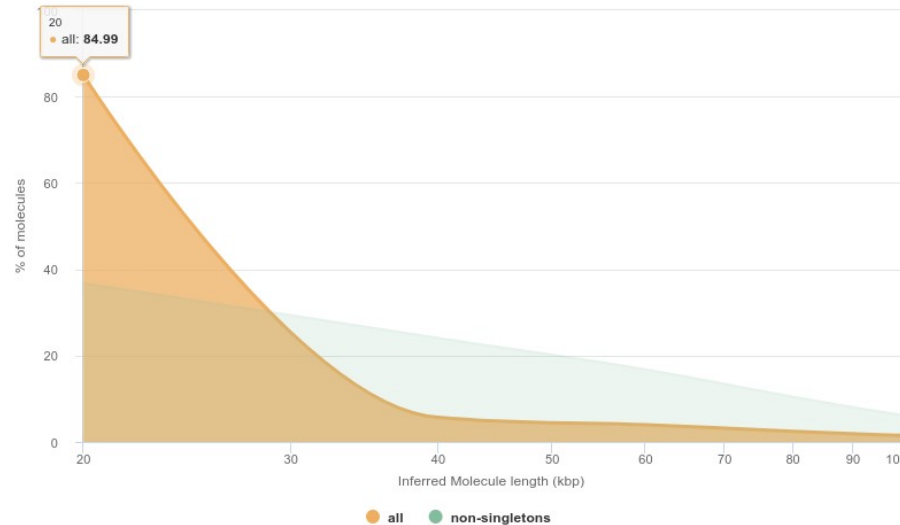
Avg  
molecules  
per contig  
538

CSV

contig	valid BX records	invalid BX records	molecules	N50	N75	N90
7	140434	8161	91487	61403	37267	19609
0	131344	7845	84806	61214	37206	19892
3	129261	7754	83637	61014	37051	19967
6	120673	6981	78249	61018	37419	20024
2	116393	6882	75984	60772	36652	19529
4	106623	6190	69440	60383	36279	19315



## Inferred Molecule Length

lengths reported as kilobases (kbp)



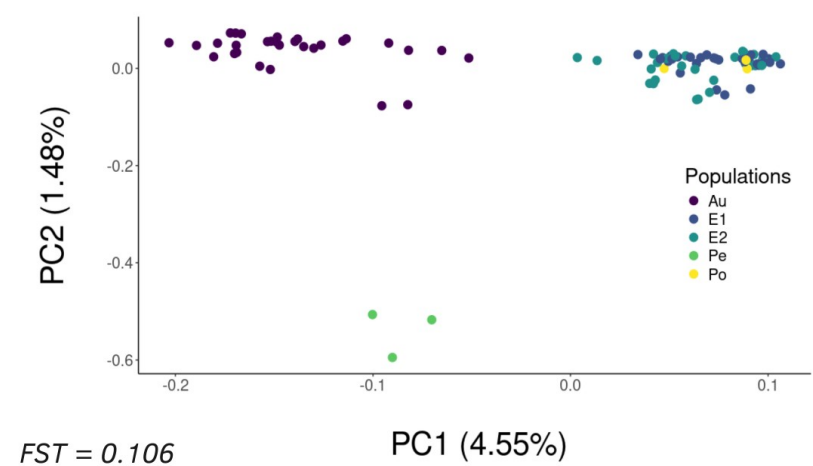
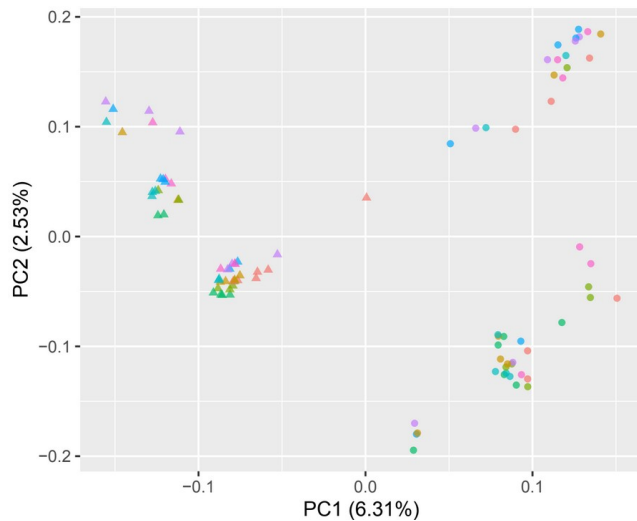
unique molecules: 633220 (350667 non-singletons)

## A close-up photograph of a butterfly with orange and black wings perched on a vibrant purple thistle flower. The flower has many thin, spiky petals. The background is a clear, bright blue sky. The image is framed by a white border.

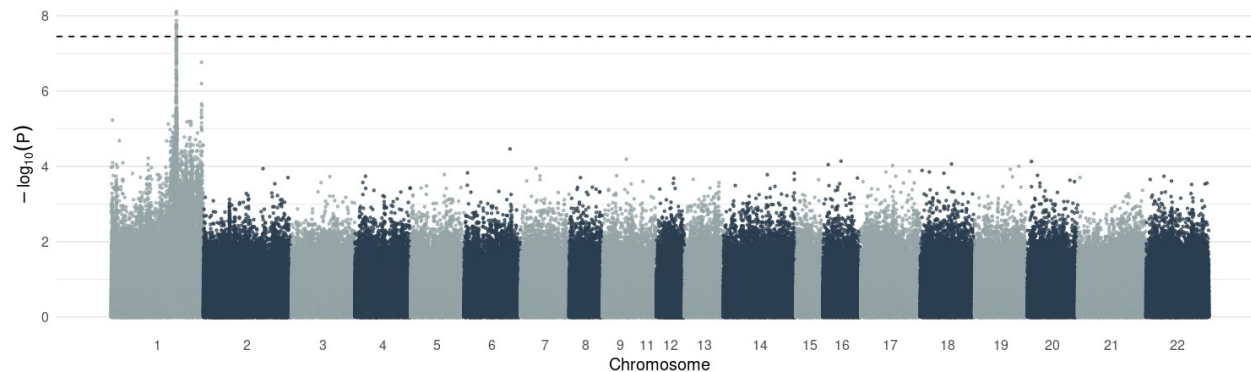


## Phased haplotype block length

# Analyses de structure génétique



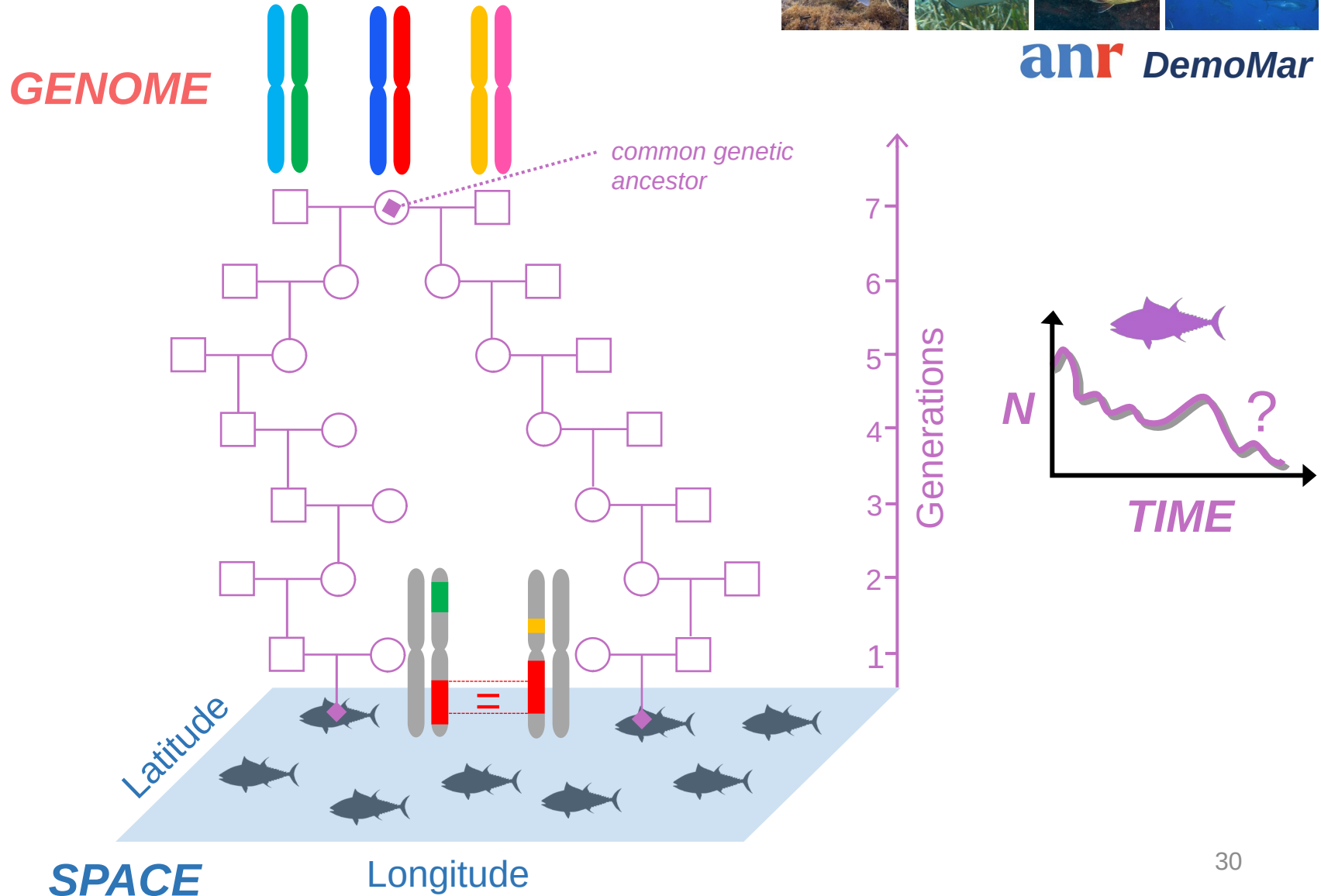
GWAS For Sex, mode dominant for male, PCA as covariates



# Analyse des co-ascendances récentes

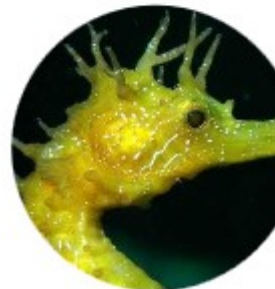


anr DemoMar



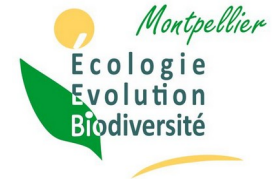
# CONCLUSIONS : exemples d'applications de l'haplotagging

- ✓ Détection des segments IBD pour étudier la démographie récente (centaurée)
- ✓ Détection de signaux d'introggression adaptative (cione)
- ✓ Détection d'inversions chromosomiques et de leurs points de cassure (hippocampes)



# Remerciements

L'équipe *GenSeq* (UAR MEEB)



Partenaires du projet *DevOCGen*

Groupe de discussion haplotagging

**MPI** : Cecile Molinier, Lauric Reynes, Frank Chan, Marek Kucka

**Roscoff** : Claire Daguin, Thomas Broquet, Alexis Simon

**Rennes** : Claire Mérot, Claire Lemaitre, Fabrice Legeai

**Paris** : Pierre de Villemereuil, Mélanie Januario, Pascaline Chifflet-Belle, Elise Gay

Développement et soutien bioinfo

**Harpy** : Pavel Dimens

**Plateforme MBB** : Khalid Belkhir, Benjamin Penaud, Iago Bonnici

