



HAL
open science

Pipeline for Haplotype Frequencies Estimation from Pooled Targeted Sequencing in Maize

Minguella Raphaël, Aurélie Canaguier, Delphine Madur, Aurélie Bérard, Isabelle Le Clainche, Agustin O. Galaretto, Damien Hinsinger, Stephane Nicolas, Patricia Faivre Rampant

► **To cite this version:**

Minguella Raphaël, Aurélie Canaguier, Delphine Madur, Aurélie Bérard, Isabelle Le Clainche, et al.. Pipeline for Haplotype Frequencies Estimation from Pooled Targeted Sequencing in Maize. AG France Génomique, Jun 2024, Evry courcouronnes, France. 2013, 10.1093/molbev/mst016 . hal-04670307

HAL Id: hal-04670307

<https://hal.inrae.fr/hal-04670307v1>

Submitted on 21 Aug 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

Minguella Raphaël¹, Canaguier Aurélie¹, Madur Delphine², Berard Aurélie¹, Le Clainche Isabelle¹, Galaretto Agustin-Oscar², Hinsinger Damien D.¹, Nicolas Stéphane D.², Faivre Rampant Patricia¹

¹ INRAE, EPGV, Evry, France ² INRAE, GQE, Gif-sur-Yvettes, France

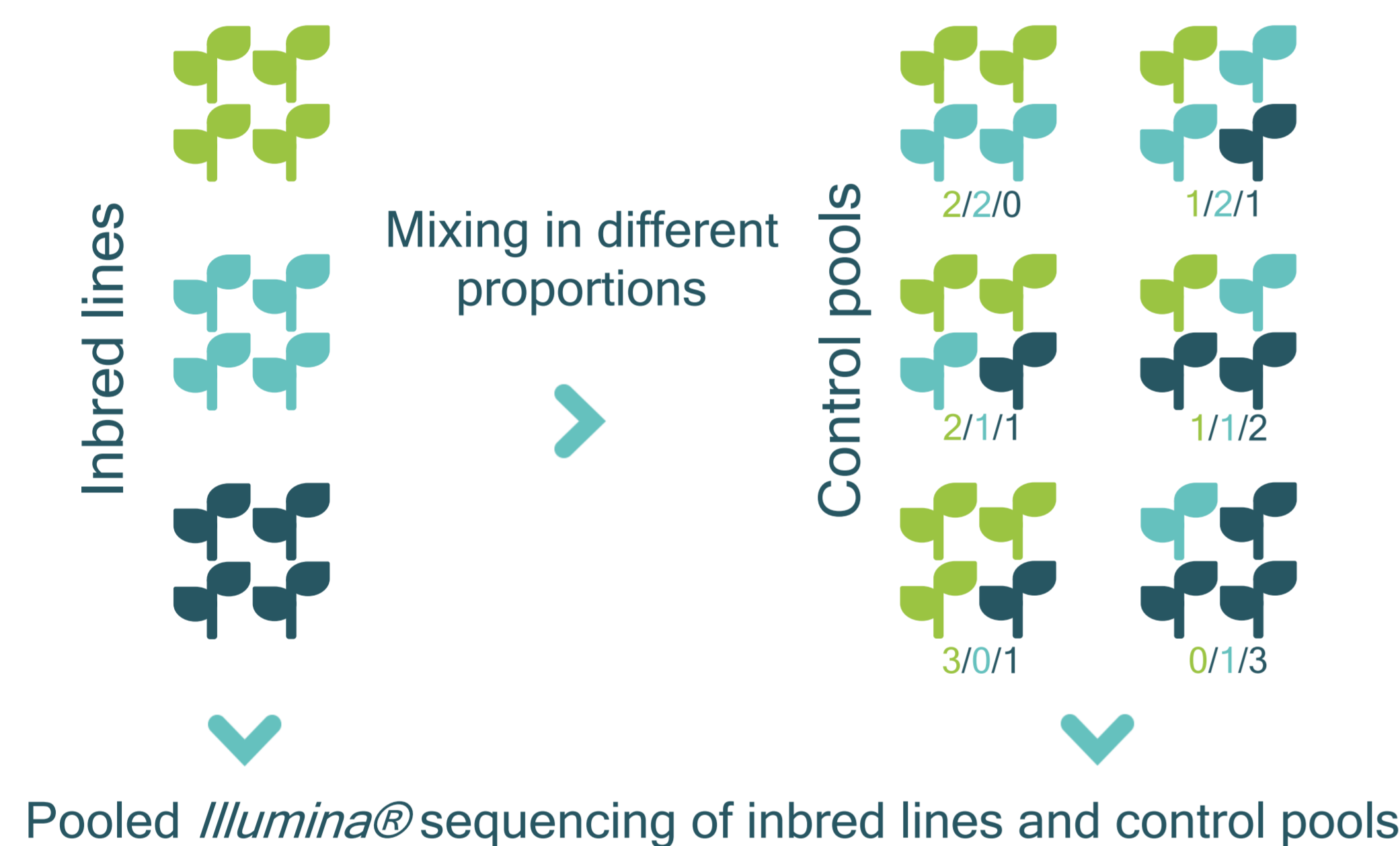
1. Background

Haplotypes are useful markers in population genetics due to a tighter link to populations history than SNP and are therefore considered more informative for populations structuration analyses. However, capturing populations **diversity** can be challenging because it requires to genotype many individuals which can be very expensive. An usual solution is to genotype populations in **pool**, meaning that several individuals of a same population are mixed and their DNA extraction is done in pool. Unfortunately, information about haplotypes is lost during the process because the DNA is fragmented and mixed. Several approaches have been proposed to rebuild haplotypes with reads overlap. Here implemented one of these algorithm in a pipeline to estimates short haplotypes frequencies from targeted **genotyping by sequencing** (tGBS) technology, which **reduces costs** by sequencing only desired genomic regions.

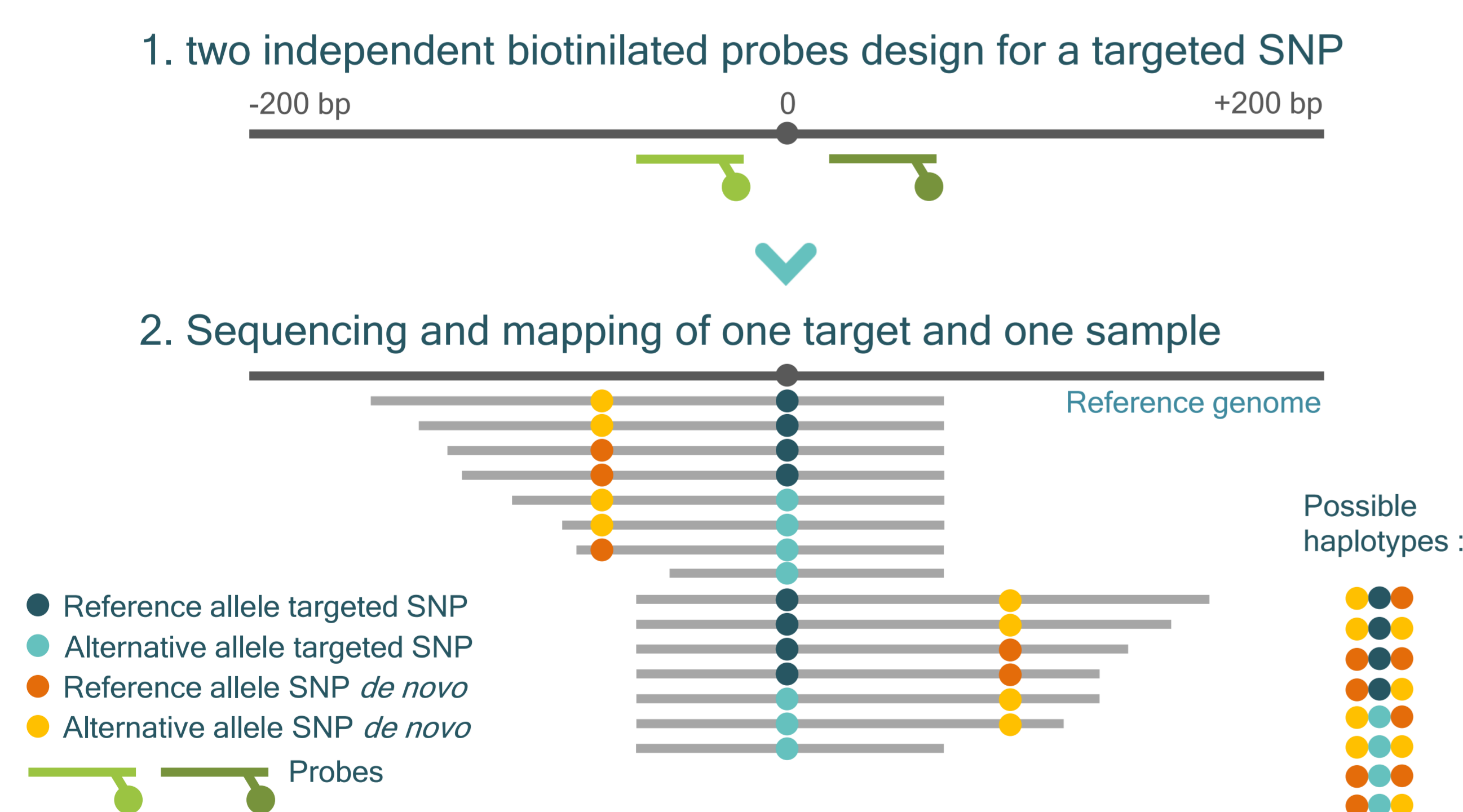
We assessed the accuracy of our pipeline using control pools with known haplotypes frequencies and we show that our pipeline gives correct haplotypic frequencies estimations when frequencies are higher than 5%.

2. Control pools design

To **evaluate** haplotypic frequencies estimation quality we designed **control pools** that are **mixes** of DNA from homozygous inbred lines in known proportions (or known F1 hybrids). Therefore, we can calculate **expected haplotype frequencies** for each control pool from the haplotyping of inbred lines that is assumed to be correct, and then, compare it to the **observed frequencies** in each control pools. We used 30 control pools : 5 F1 hybrids and 25 mixes of 3 inbred lines including different genetic groups.

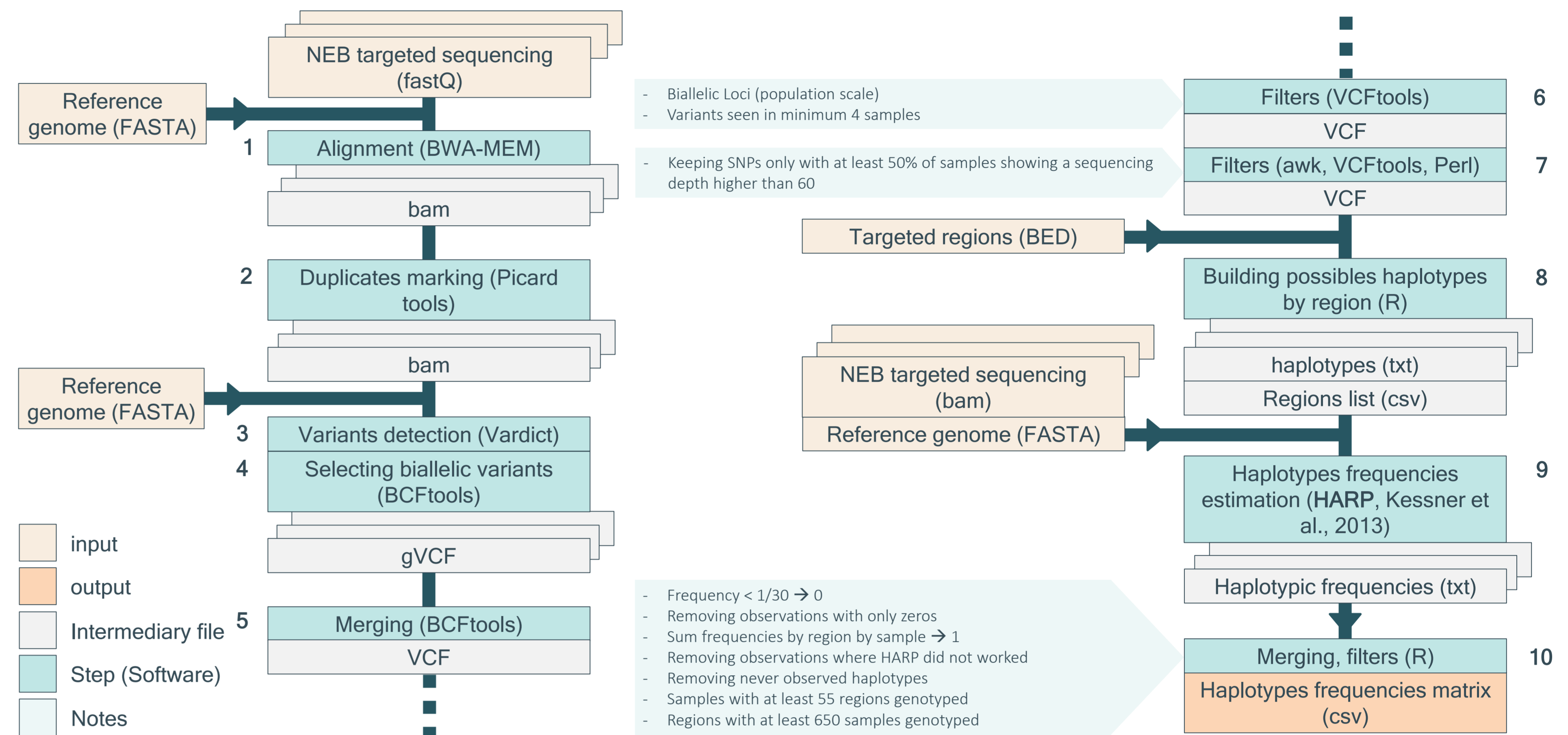


3. tGBS : NEBnext® direct genotyping solution



➤ Pipeline for Haplotype Frequencies Estimation from Pooled Targeted Sequencing in Maize

4. Pipeline



5. Performances

To estimate the quality of haplotypes detection, we assessed the qualitative detection of this haplotype (right) for each haplotypes of each region and each pool. Then, we calculated for each expected frequency the proportion of each category that we plotted (right).

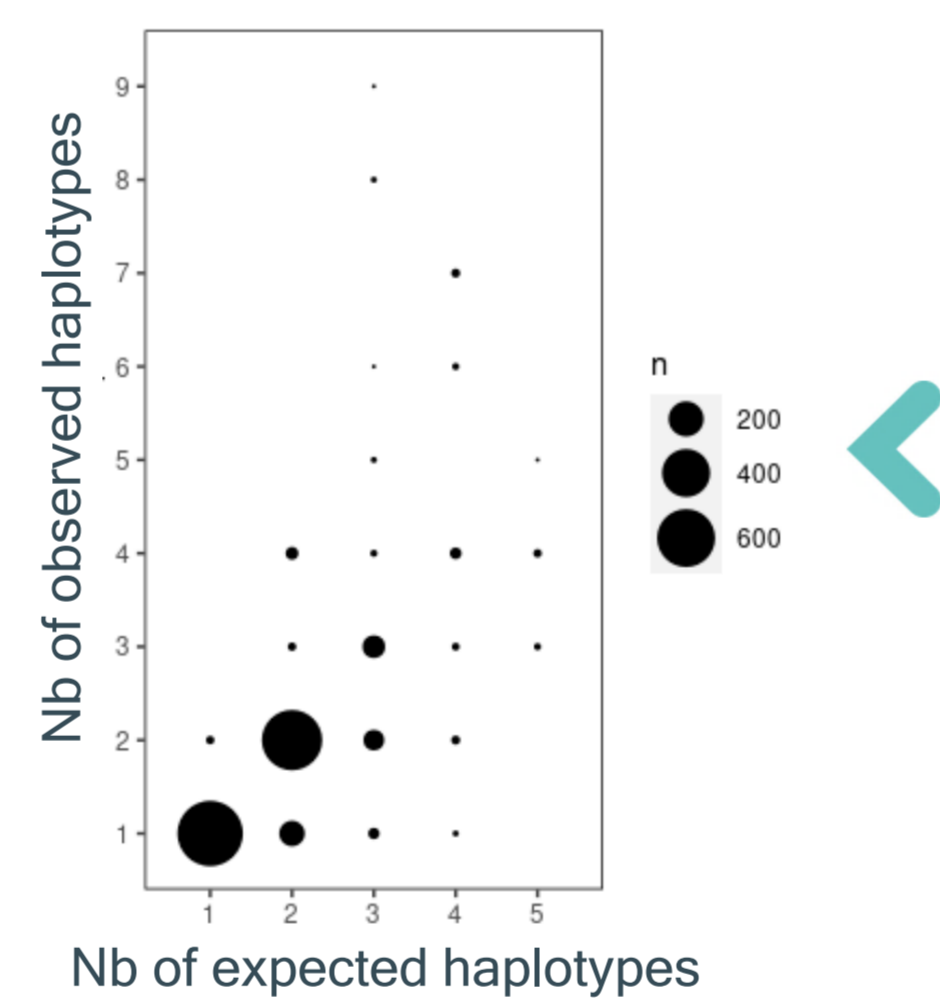
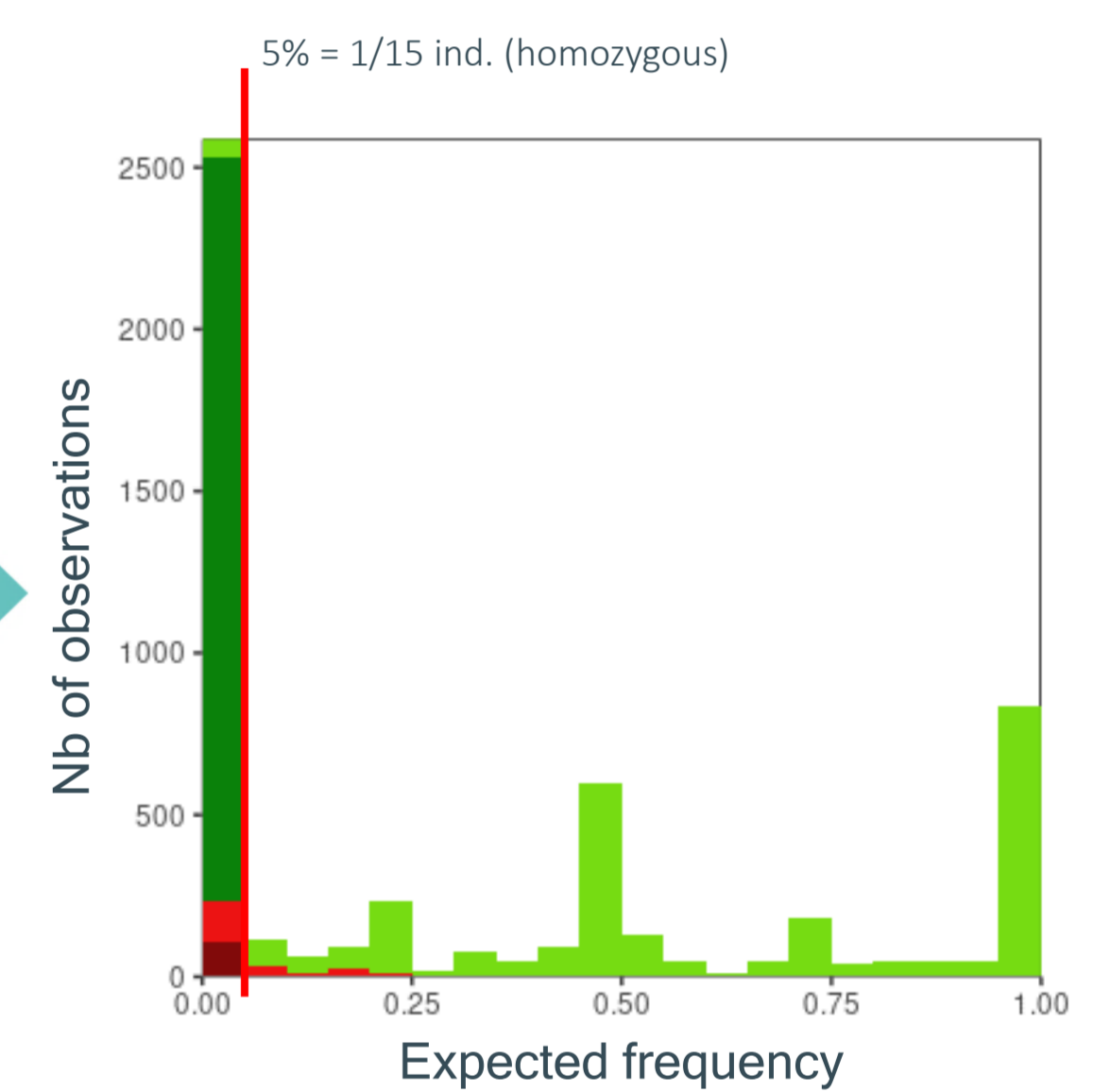
When haplotype frequency >5% :

- Very few false positives (=unexpected but observed haplotypes)
- Very few false negatives (=undetected haplotypes)

Haplotypes presence or absence

expected	observed	detected
expected	observed	Unexpected and not observed
expected	observed	undetected
expected	observed	Unexpected but observed

Good haplotypes detection if found in at least 1/15 ind. in the pool



To estimate the ability of our pipeline to properly detect all haplotypes in a pool, we calculated the number of observed and expected haplotypes for each region of each control pool :

- Correct number of detected haplotypes (1-2 haplotypes in a pool)
- When >2 haplotypes in the pool : a few undetected haplotypes
- >3 haplotypes observed due to residual heterozygosity in inbred lines
- Few cases with detected hap. > expected hap. (HARP errors)

The pipeline retrieves the expected number of haplotypes

Haplotypes frequencies

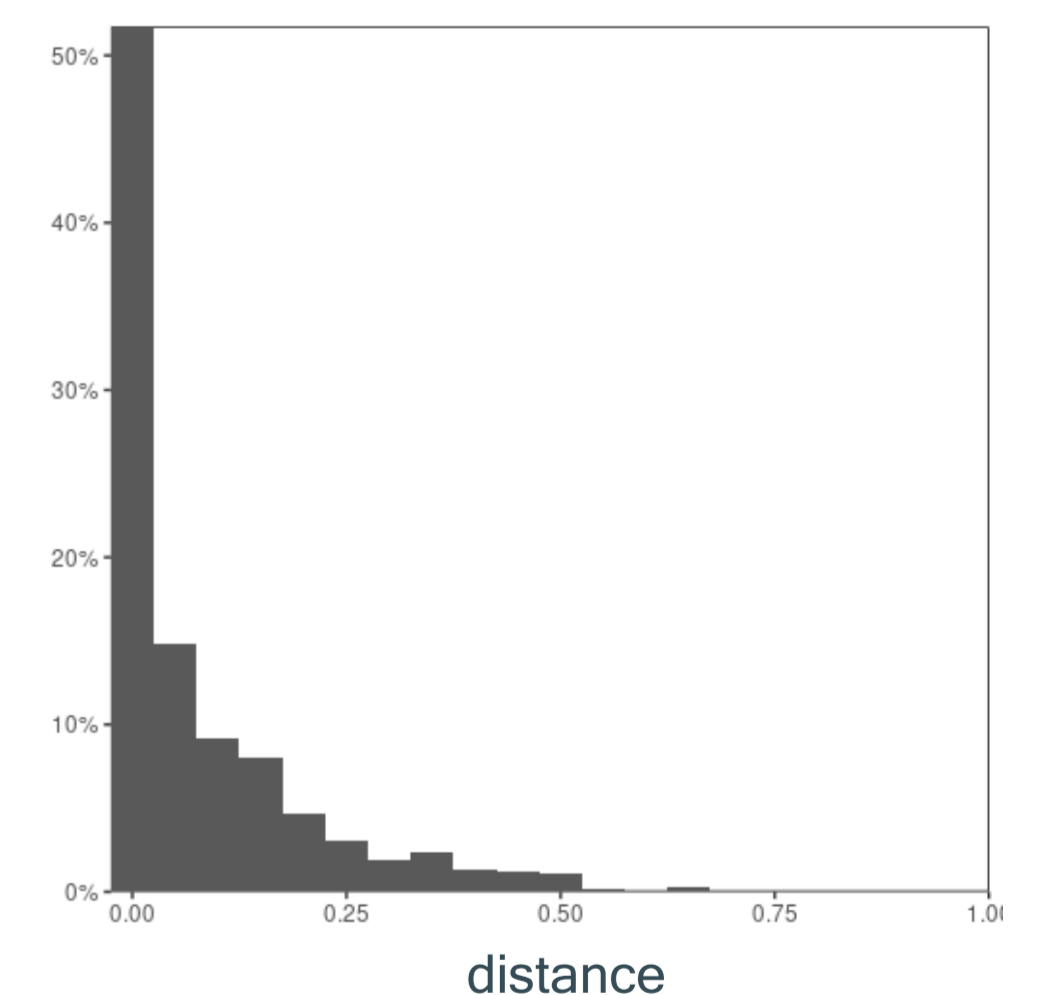
expected	observed	freq
expected	observed	1
expected	observed	0

Observed frequencies close from expected ones, no very distant observation

To estimate the quality of frequencies estimation, we calculated for each region and each pool the sum of absolute differences between expected and observed haplotypic frequencies :

$$distance = \frac{1}{2} \sum_{haplos} |f_{obs} - f_{exp}|$$

- >50% of freq. correctly estimated (distance < 0.05)
- >80% of freq. estimated with a distance < 0.2



6. Conclusions & perspectives

- Haplotypes are detected if there frequency is more than 5% (=1/15 individual in the pool) and frequencies are correctly estimated in pools
- Cheap short haplotypes sequencing approach, but improvement could be useful during probes design
- First results on actual data show that maize landraces are more diversified than maize inbred lines
- Maize landraces harbor haplotypes that are not in inbred lines

Center
Île-de-France - Versailles-Saclay

Acknowledgment



Bibliography

NEBnext® direct genotyping solution:
Emerman AB, Bowman SK, Barry A, Henig N, Patel KM, Gardner AF, Hendrickson CLJCPiMB: NEBNext Direct: A Novel, Rapid, Hybridization-Based Approach for the Capture and Library Conversion of Genomic Regions of Interest. 2017;119(1):7.30. 31-37.30. 24.

HARP :
Darren Kessner, Thomas L. Turner, John Novembre, Maximum Likelihood Estimation of Frequencies of Known Haplotypes from Pooled Sequence Data, Molecular Biology and Evolution, Volume 30, Issue 5, May 2013, Pages 1145-1158, <https://doi.org/10.1093/molbev/mst016>

➤ INRAE - Unité EPGV US1279
CEA - Institut de biologie François Jacob
Site d'Evry - Bat G1
2 rue Gaston Crémieux
91057 - Evry Cedex
France
support-epgv@inrae.fr

<https://www6.versailles-grignon.inrae.fr/epgv>