



**HAL**  
open science

# Estimating the allele frequency spectrum and detecting selective sweeps from pooled next generation sequencing samples

Simon S. Boitard, Christian Schlötterer, Robert Koffler, Ram Vinay Pandey,  
Andreas Futschik

## ► To cite this version:

Simon S. Boitard, Christian Schlötterer, Robert Koffler, Ram Vinay Pandey, Andreas Futschik. Estimating the allele frequency spectrum and detecting selective sweeps from pooled next generation sequencing samples. 20. Annual Meeting of the Society for Molecular Biology and Evolution (SMBE), Jun 2012, Dublin, Ireland. 2012. hal-02807025

**HAL Id: hal-02807025**

<https://hal.inrae.fr/hal-02807025>

Submitted on 6 Jun 2020

**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.

# Estimating the allele frequency spectrum and detecting selective sweeps from pooled next generation sequencing samples.



**Simon Boitard**<sup>1</sup>, **Christian Schlötterer**<sup>2</sup>, **Robert Kofler**<sup>2</sup>, **Ram Vinay Pandey**<sup>2</sup>, **Andreas Futschik**<sup>3</sup>. (1) Laboratoire de Génétique Cellulaire, INRA, Toulouse, France. (2) Institut für Populationsgenetik, Vetmeduni Vienna, Vienna, Austria. (3) Institute of Statistics and Operations Research, University of Vienna, Vienna, Austria.

## SUMMARY

Due to its cost effectiveness, next generation sequencing of pools of individuals (Pool-Seq) is becoming a popular strategy for genome-wide estimation of allele frequencies in population samples. Since the allele frequency spectrum (AFS) provides information about past episodes of selection, Pool-Seq is also a promising design for genomic scans for selection. Here, we introduce a statistical method for estimating the AFS and detecting selective sweeps in a Pool-Seq sample. This method accounts for the uncertainty concerning the true allele frequencies in the pool, which might typically be higher for sites with low coverage or bad quality scores. We apply our method to simulated data at low coverage ( $0.5\times$  per chromosome), and to real data from *Drosophila melanogaster*. A python program implementing our method, denoted **Pool-hmm**, is freely available at <https://qgp.jouy.inra.fr/>.

## MODELLING THE POOL-SEQ DATA

Consider a pool of  $n$  chromosomes from the same population. At position  $i$  on the genome, let  $Y_i$  be the number of derived alleles in the pool ( $0 \leq Y_i \leq n$ ) and  $r_i$  be the number of reads covering the position. Observed read data at position  $i$  is written  $Z_i = (Z_{i,1}, \dots, Z_{i,r_i})$ , where  $Z_{i,j}$  is equal to 1 (0) if the allele at read  $j$  is derived (ancestral). If  $p_{i,j}$  is the probability for a sequencing error at read  $j$  (deduced from quality score),

$$\mathbb{P}(Z_i | Y_i) = \prod_{j:Z_{i,j}=1} \left( (1-p_{i,j}) \frac{Y_i}{n} + p_{i,j} \left(1 - \frac{Y_i}{n}\right) \right) * \prod_{j:Z_{i,j}=0} \left( (1-p_{i,j}) \left(1 - \frac{Y_i}{n}\right) + p_{i,j} \frac{Y_i}{n} \right)$$

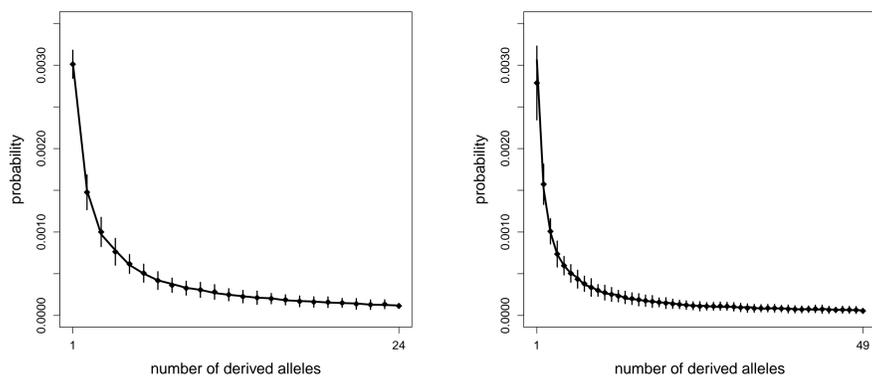
## DETECTION OF SELECTIVE SWEEPS - METHODS

Detection of selective sweeps is based on a Hidden Markov Model. Each position  $i$  is assumed to have a hidden state  $X_i$ , with 3 possible values: "Selection", "Intermediate" or "Neutral". These 3 values are associated to different AFS. The "Neutral" AFS is estimated from whole genome data. The "Intermediate" and "Selection" AFS are derived from the "Neutral" AFS as in (Nielsen et al., 2005). To account for the uncertainty concerning the true allele frequencies in the pool, emission probabilities are

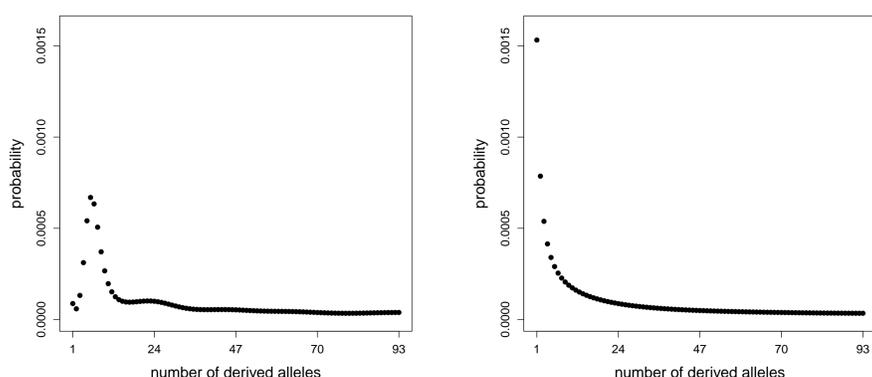
$$\mathbb{P}(Z_i | X_i) = \sum_{Y_i=0}^n \mathbb{P}(Z_i | Y_i) \mathbb{P}(Y_i | X_i)$$

Hidden states (and consequently selective sweeps) are predicted from the observed Pool-Seq data using the Viterbi algorithm.

## ESTIMATION OF THE AFS



**Simulation results.** 100 pools of  $n = 25$  (left) and  $n = 50$  (right) chromosomes of length  $L = 100\text{kb}$  were simulated under the coalescent with  $\theta = 0.003$  and  $\rho = 0.003$ . Solid lines (diamonds) show the average AFS estimated from the complete sequences (Pool-Seq data at  $100\times$  coverage). **If quality scores are reliable, our estimation of the AFS is very accurate.**

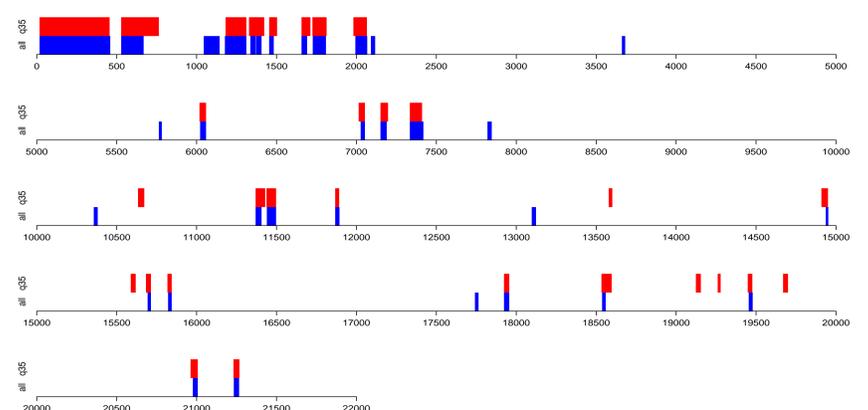


**Folded AFS on the X chromosome of *D. melanogaster*.** Estimated from a pool of 97 flies sequenced at  $100\times$  coverage, using all base calls (left) or only those with PHRED score greater than 35 (right). **The AFS is biased because quality scores overestimate the sequencing error probabilities.**

## DETECTION OF SELECTIVE SWEEPS - RESULTS

Sample size	$n = 25$	$n = 50$	$n = 100$	$n = 200$
Pool-Seq data	0.91	0.90	0.91	0.90
Sequence data	0.89	0.88	0.87	0.87

**Simulation results.** Samples were simulated with selection intensity  $s = 500$ , other parameters are as in the left box. **Pool-Seq data and complete sequence data provide the same detection power.**



**Selective sweeps on the X chromosome of *D. melanogaster*,** detected using all base calls (blue) or base calls with PHRED score greater than 35 (red). **Sweep detection is not affected by the bias of the "Neutral" AFS.**

## REFERENCES

- Boitard, S., Schlötterer, C., Nolte, V., Pandey, R., Futschik, A., 2012. Detecting selective sweeps from pooled next generation sequencing samples. *Mol. Biol. Evol.*, doi: 10.1093/molbev/mss090.
- Nielsen, R., Williamson, L., Kim, Y., Hubisz, M., Clark, A., Bustamante, C., 2005. Genomic scans for selective sweeps using SNP data. *Genome Research* 15, 1566–1575.