



**HAL**  
open science

## Performance evaluation of adaptive introgression classification methods

Jules Romieu, Ghislain Camarata, Pierre-André Crochet, Miguel de Navascués, Raphaël Leblois, François Rousset

► **To cite this version:**

Jules Romieu, Ghislain Camarata, Pierre-André Crochet, Miguel de Navascués, Raphaël Leblois, et al.. Performance evaluation of adaptive introgression classification methods. MCEB 2024 : Mathematical and computational Evolutionary biology 2024, Jun 2024, St Martin de Londres, France. hal-04713611

**HAL Id: hal-04713611**

**<https://hal.inrae.fr/hal-04713611v1>**

Submitted on 29 Sep 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.

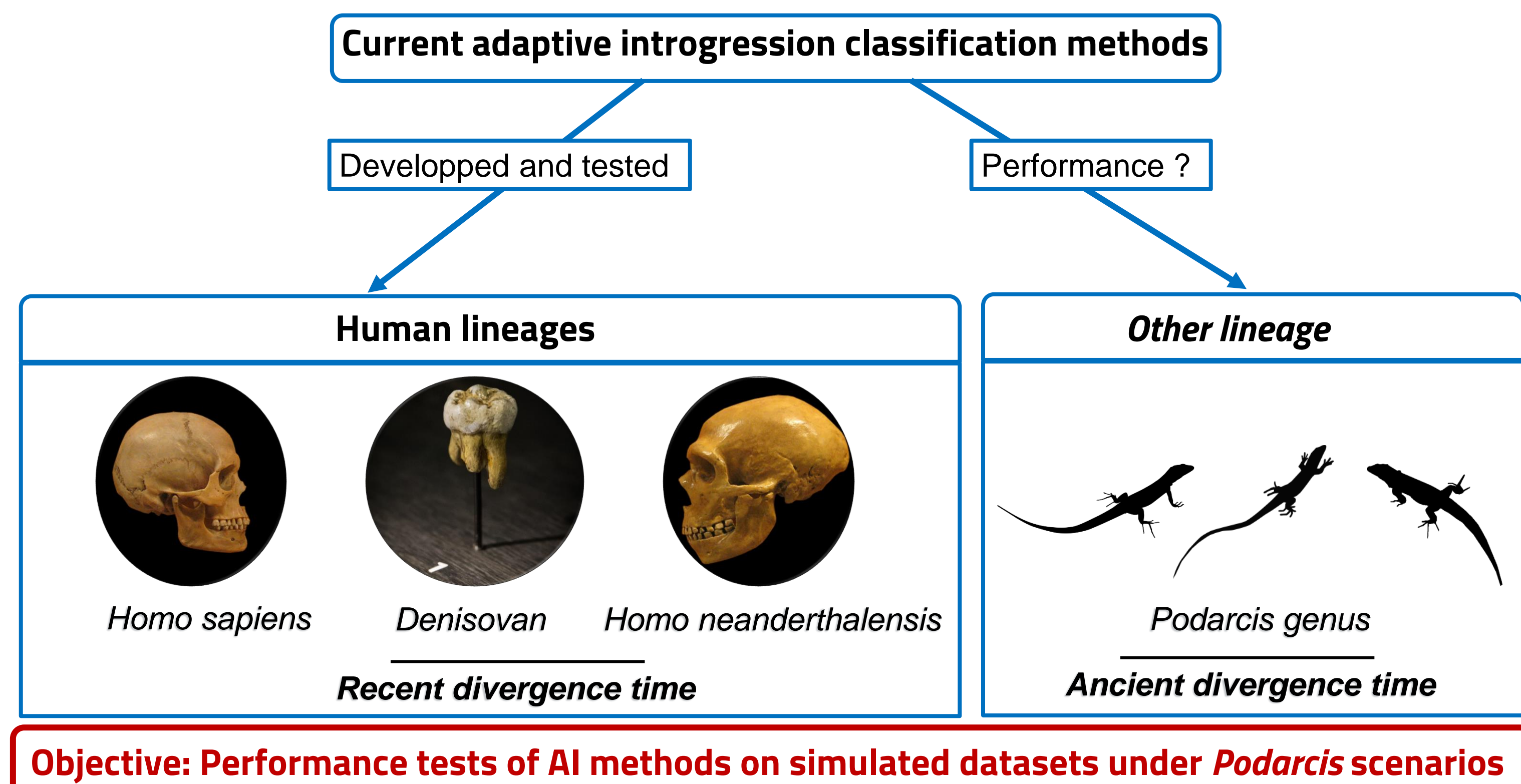
L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Performance evaluation of adaptive introgression classification methods

Jules Romieu<sup>1,2</sup>, Ghislain Camarata<sup>1,2</sup>, Pierre-André Crochet<sup>3</sup>, Miguel de Navascués<sup>2</sup>, Raphaël Leblois<sup>2</sup>, and François Rousset<sup>1</sup>

<sup>1</sup>ISEM, Univ Montpellier, CNRS, IRD, Montpellier, France, <sup>2</sup>CBGP, INRAE, CIRAD, IRD, Institut Agro, Univ Montpellier, Montpellier, France, <sup>3</sup>CEFE, CNRS, Univ Montpellier, EPHE, IRD, Montpellier, France

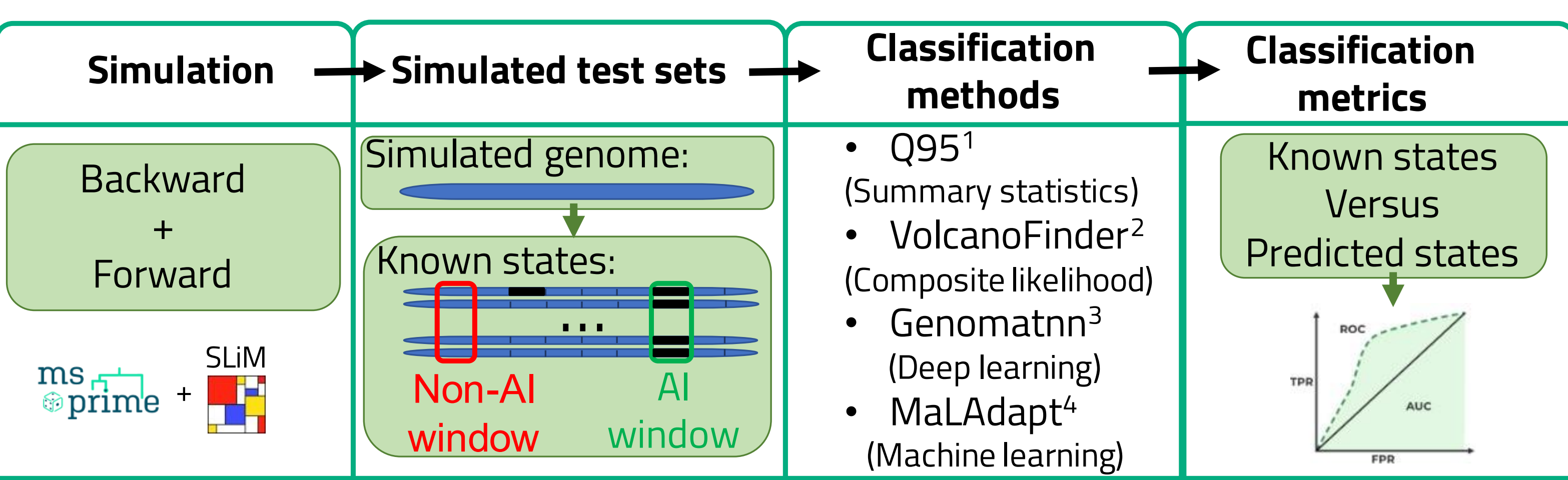
## 1 Background and objective



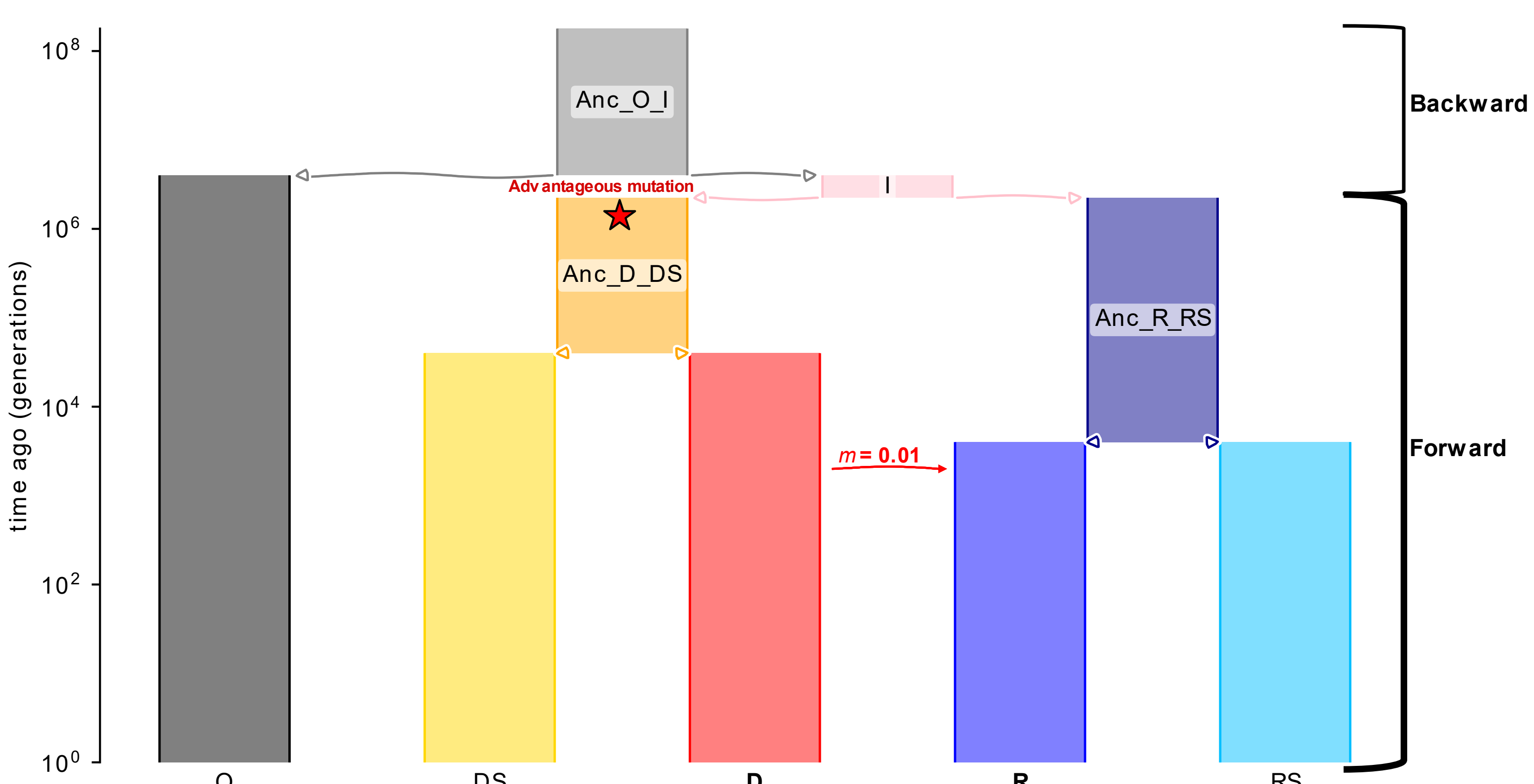
## 2 Summary statistic and methods tested

Q95 (1%, 100%) (Racimo <i>et al.</i> 2017)	VolcanoFinder (Setter <i>et al.</i> 2020)	Genomatnn (Gower <i>et al.</i> 2021)	MaLAdapt (Zhang <i>et al.</i> 2023)
<ul style="list-style-type: none"> <li>Summary statistic</li> <li>50kb windows</li> <li>High frequency of alleles uniquely shared between the donor and the recipient populations</li> <li>Uses samples from 3 populations</li> </ul>	<ul style="list-style-type: none"> <li>Composite likelihood-based genomic scan</li> <li>One test site every kb</li> <li>Excess intermediate frequency polymorphism in the flanking region of the genome</li> <li>Uses only the recipient population sample</li> </ul>	<ul style="list-style-type: none"> <li>Deep learning (Convolutional Neural Network)</li> <li>100kb windows</li> <li>Black box (train on genotype matrix)</li> <li>Uses samples from 3 populations</li> </ul>	<ul style="list-style-type: none"> <li>Machine learning (Extra Trees Classifier)</li> <li>50kb windows</li> <li>Train on summary statistics (alleles shared between donor and recipient, linkage disequilibrium, genetic diversity, etc)</li> <li>Uses samples from 3 populations</li> </ul>

## 3 Simulation test approach



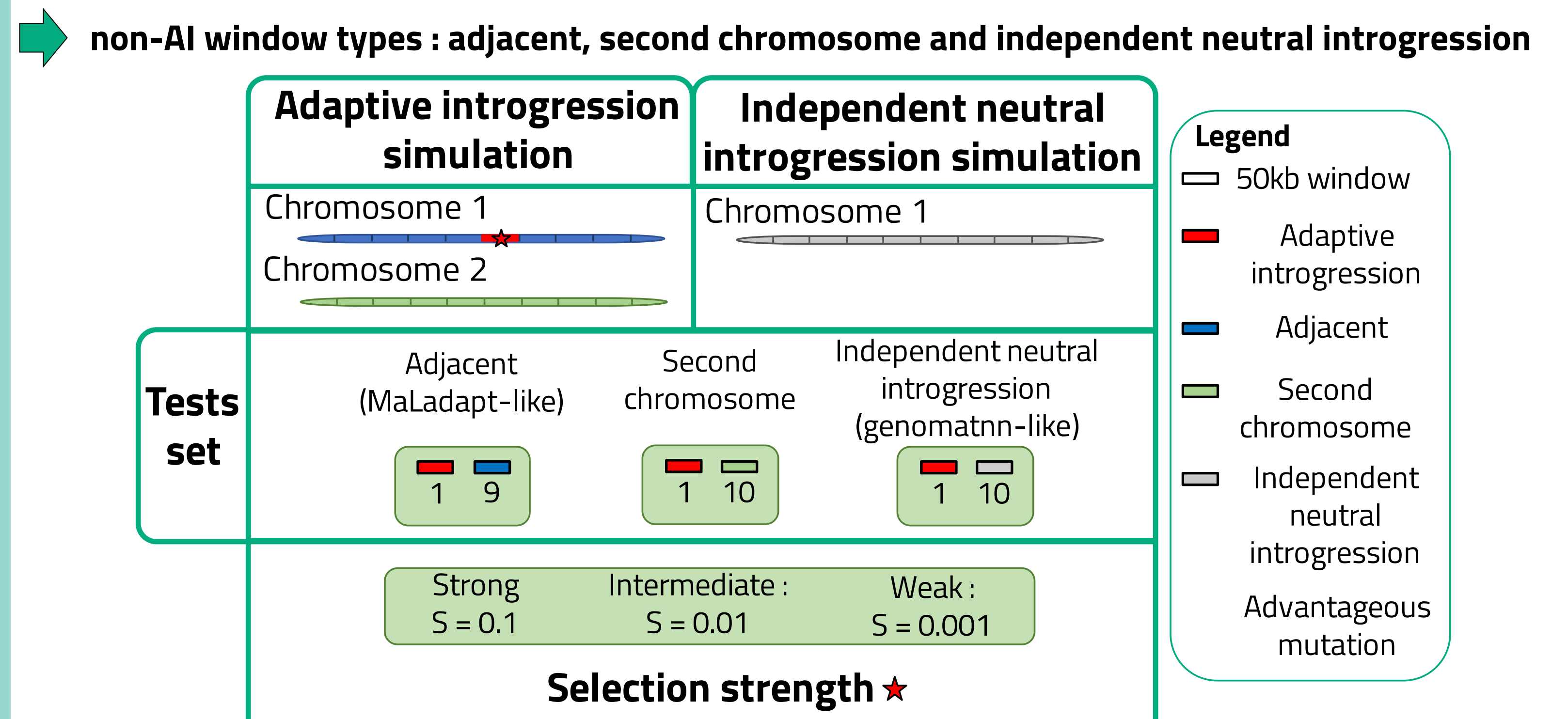
## 4 Demographic model



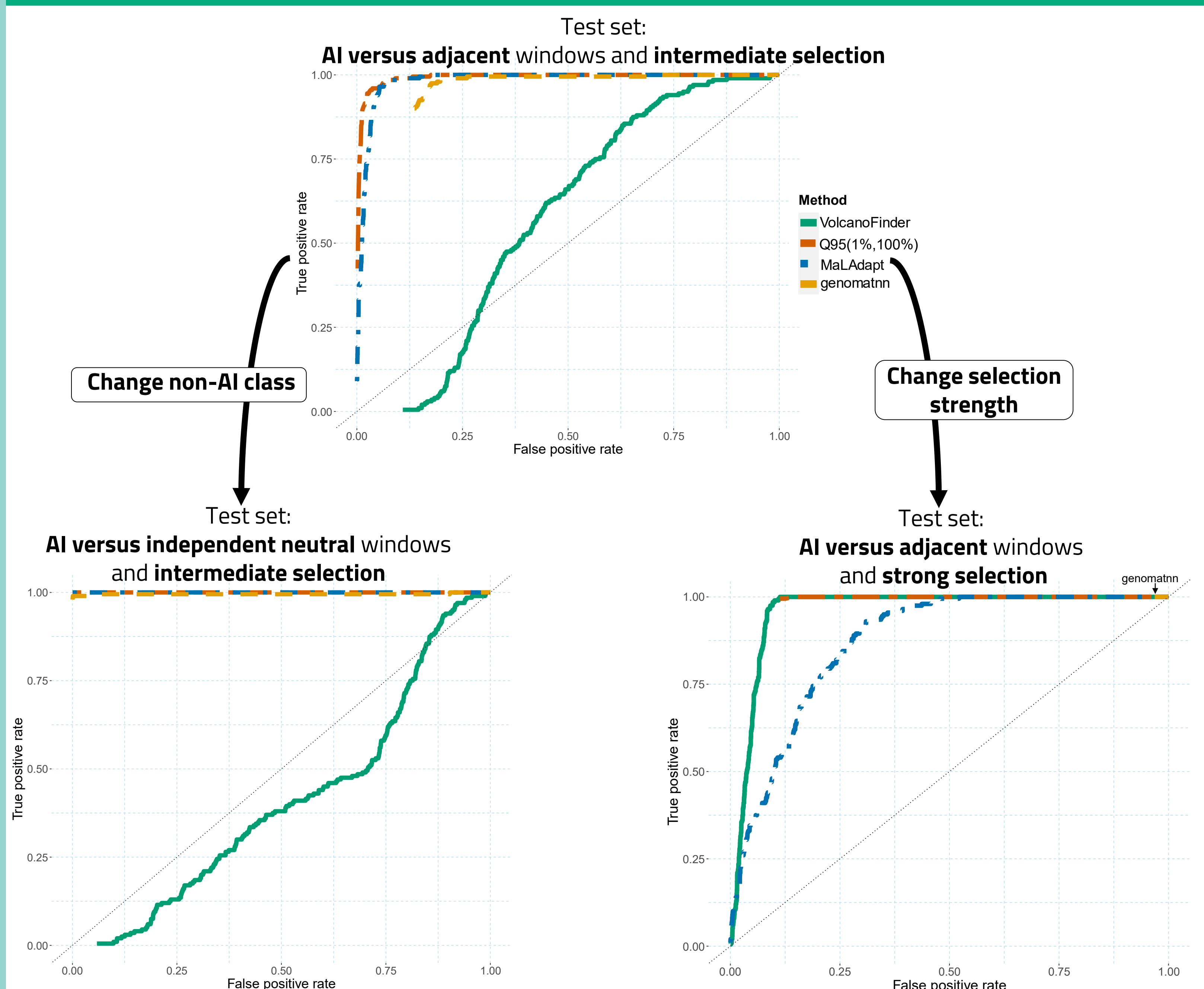
**Figure 1 – Demographic model used to generate simulated dataset.** Anc\_O\_I = Outgroup and Ingroup ancestral population, O = Outgroup population, I = Ingroup population, Anc\_D\_DS = Donor and non-introgressing ancestral population, Anc\_R\_RS = Recipient and non-introgressed ancestral population, DS = sister to the donor population (non-introgressing population), D = Donor population (introgressing population), R = Recipient population (introgressed population), RS = Recipient sister population (non-recipient population), Population size (N) = 10,000.

## 5 Importance of genetic architecture

### Hitchhiking effect on method's performance?



## 6 Results: ROC curves



## 7 Discussion and conclusion

- Best method in our tests: Q95 summary statistic
- Better performance with strong selection but leads to an increase in hitchhiking effect
  - Take adjacent windows into account
- Simulation-based inference methods: Factors that have negative impacts on performance:
  - Test sets with non-AI windows types non-used in train sets (ex :genomatnn)
  - Misspecification of demographic model used in train sets
- Importance to take into account different types of neutral windows in trained sets
- Classification methods problems:
  - Define a threshold to discriminate AI/non-AI windows
  - Necessity to use FDR control methods
- Solution: Develop methods to estimate genomic-level introgression

## Bibliography

<sup>1</sup>Racimo F, Marnetto D, Huerta-Sánchez E (2017). Signatures of Archaic Adaptive Introgression in Present-Day Human Populations. *Molecular Biology and Evolution* 34, 296–317.

<sup>2</sup>Setter D, Mousset S, Cheng X, Nielsen R, DeGiorgio M, Hermisson J (2020). VolcanoFinder: Genomic scans for adaptive introgression. *PLoS Genetics* 16. e1008867.

<sup>3</sup>Gower G, Picazo PI, Fumagalli M, Racimo F (2021). Detecting adaptive introgression in human evolution using convolutional neural networks. *eLife* 10. e64669.

<sup>4</sup>Zhang X, Kim B, Singh A, Sankararaman S, Durvasula A, Lohmueller KE (2023). MaLAdapt Reveals Novel Targets of Adaptive Introgression From Neanderthals and Denisovans in Worldwide Human Populations. *Molecular Biology and Evolution* 40, msad001.