



**HAL**  
open science

# PSAP-Genomic-Regions: A Method Leveraging Population Data to Prioritize Coding and Non-Coding Variants in Whole Genome Sequencing for Rare Disease Diagnosis

Marie-sophie Ogloblinsky, Ozvan Bocher, Chaker Aloui, Anne-louise Leutenegger, Ozan Ozisik, Anaïs Baudot, Elisabeth Tournier-Lasserre, Helen Castillo-Madeen, Daniel Lewinsohn, Donald F Conrad, et al.

## ► To cite this version:

Marie-sophie Ogloblinsky, Ozvan Bocher, Chaker Aloui, Anne-louise Leutenegger, Ozan Ozisik, et al.. PSAP-Genomic-Regions: A Method Leveraging Population Data to Prioritize Coding and Non-Coding Variants in Whole Genome Sequencing for Rare Disease Diagnosis. *Genetic Epidemiology*, 2024, Online ahead of print. 10.1002/gepi.22593 . hal-04746483

**HAL Id: hal-04746483**

**<https://hal.science/hal-04746483v1>**

Submitted on 21 Oct 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.

1 PSAP-genomic-regions: a method leveraging population  
2 data to prioritize coding and non-coding variants in whole  
3 genome sequencing for rare disease diagnosis

4 Marie-Sophie C. Ogloblinsky<sup>1,\*</sup>, Ozvan Bocher<sup>1,2</sup>, Chaker Aloui<sup>3</sup>, Anne-Louise Leutenegger<sup>3</sup>, Ozan Ozisik<sup>4</sup>,  
5 Anaïs Baudot<sup>4</sup>, Elisabeth Tournier-Lasserre<sup>3,5</sup>, Helen Castillo-Madeen<sup>6</sup>, Daniel Lewinsohn<sup>6</sup>, Donald F.  
6 Conrad<sup>6</sup>, Emmanuelle Génin<sup>1,7,¶</sup>, Gaëlle Marenne<sup>1,\*</sup>,¶

7  
8 <sup>1</sup>Univ Brest, Inserm, EFS, UMR 1078, GGB, Brest, France

9 <sup>2</sup>Institute of Translational Genomics, Helmholtz Zentrum München, Munich, Germany

10 <sup>3</sup>Université Paris Cité, Inserm, NeuroDiderot, Unité Mixte de Recherche 1141, Paris, France

11 <sup>4</sup>Aix Marseille Univ, INSERM, Marseille Medical Genetics (MMG), Marseille, France

12 <sup>5</sup>Assistance publique-Hôpitaux de Paris, Service de Génétique Moléculaire Neurovasculaire, Hôpital Saint-  
13 Louis, Paris, France

14 <sup>6</sup>Division of Genetics, Oregon National Primate Research Center, Oregon Health & Science University,  
15 Portland, Oregon, United States of America

16 <sup>7</sup>Centre Hospitalier Régional Universitaire de Brest, Brest, France

17

18 \*Corresponding authors:

19 Email: [marie-sophie.ogloblinsky@inserm.fr](mailto:marie-sophie.ogloblinsky@inserm.fr) (M-S.O.); [gaelle.marenne@inserm.fr](mailto:gaelle.marenne@inserm.fr) (G.M.)

20

21 ¶These authors contributed equally to this work

## 22 Abstract

23 The introduction of next generation sequencing technologies in the clinics has improved rare  
24 disease diagnosis. Nonetheless, for very heterogeneous or very rare diseases, more than half of cases still  
25 lack molecular diagnosis. Novel strategies are needed to prioritize variants within a single individual. The  
26 PSAP method was developed to meet this aim but only for coding variants in exome data. Here, we  
27 propose an extension of the PSAP method to the non-coding genome called PSAP-genomic-regions. In this  
28 extension, instead of considering genes as testing units (PSAP-genes strategy), we use genomic regions  
29 defined over the whole genome that pinpoint potential functional constraints.

30 We conceived an evaluation protocol for our method using artificially-generated disease exomes  
31 and genomes, by inserting coding and non-coding pathogenic ClinVar variants in large datasets of exomes  
32 and genomes from the general population.

33 PSAP-genomic-regions significantly improves the ranking of these variants compared to using a  
34 pathogenicity score alone. Using PSAP-genomic-regions, more than fifty percent of non-coding ClinVar  
35 variants were among the top 10 variants of the genome. On real sequencing data from 6 patients with  
36 Cerebral Small Vessel Disease and 9 patients with male infertility, all causal variants were ranked in the  
37 top 100 variants with PSAP-genomic-regions.

38 By revisiting the testing units used in the PSAP method to include non-coding variants, we have  
39 developed PSAP-genomic-regions, an efficient whole-genome prioritization tool which offers promising  
40 results for the diagnosis of unresolved rare diseases.

41

42 **Keywords:** rare diseases, non-coding variants, whole-genome sequencing, variant prioritization

43

## 44 Introduction

45 Each rare disease affects, by definition, a small number of individuals. However, as a whole, rare  
46 diseases affect about 350 million people world-wide (1). Approximately 80% of rare diseases have a  
47 genetic origin that mostly follows a Mendelian mode of inheritance (2–4). The advent of Next Generation  
48 Sequencing (NGS) and the development of variant pathogenicity prediction tools have allowed, in recent  
49 years, the identification of many genes involved in rare Mendelian diseases. Nonetheless, despite  
50 extensive efforts, the molecular diagnosis is still unknown for more than 50% of rare diseases cases (5–7).  
51 This can mainly be explained by the fact that many rare diseases are characterized by an extreme genetic  
52 heterogeneity, which results in only one individual carrying a specific pathogenic causal variant. This issue  
53 is referred to as the “n-of-one” problem (8).

54 With the advent of high throughput sequencing technologies in clinics, molecular diagnosis is now  
55 often sought through whole exome or whole genome sequencing (WES and WGS respectively). However,  
56 due to the large number of rare variants in each individual genome, causal variants are sought among  
57 very rare and highly pathogenic variants in genes relevant to the current known disease mechanism. The  
58 limited knowledge about gene functions and disease mechanisms can make this strategy unfruitful. To  
59 address the issue of variant prioritization at the level of an individual, the Population Sampling Method  
60 (PSAP) (8) was developed. PSAP computes, for each gene, a null distribution, which is the probability to  
61 observe in the general population a genotype with a CADD pathogenicity score (9) greater than or equal  
62 to the highest one to the highest one observed in the patient for this gene. This initial version of the PSAP  
63 method, which we will refer to as PSAP-genes, has been successfully applied to identify variants of interest  
64 in diverse phenotypes, including male infertility (10–12), recurrent pregnancy loss (13) and ciliary  
65 dyskinesia (14).

66 A current hindrance to the application and generalization of PSAP-genes as a tool for diagnosis is  
67 its restriction to the coding parts of the genome. Indeed, the majority of variants reside in non-coding  
68 parts of the genome (15). Non-coding variants may contribute to explain part of the etiology of rare  
69 diseases (16), as suggested by the large number of GWAS hits located in non-coding regions of the genome  
70 (17). The involvement of non-coding pathogenic variants in rare diseases is further corroborated by the  
71 fact that non-coding regions are heavily involved in the regulation of gene expression. Several prediction  
72 tools have been developed to this end (18–20), but most of them lack a variant-based score for both  
73 coding and non-coding regions. In addition, to be performant, they often require multiple annotations like  
74 Human Phenotype Ontology (HPO) terms (21) to characterize the symptoms or disease of a patient . Thus,  
75 they rely on previous knowledge and rarely go beyond candidate genes.

76 To move beyond the gene as a natural unit of testing for the PSAP method, we need to use  
77 predetermined regions across the whole genome. These regions also need to be defined using functional  
78 information to be used as a cohesive unit for the construction of PSAP null distributions. This challenge of  
79 defining regions along the whole genome has been tackled by Bocher et al. in the context of rare-variant  
80 association testing (22): they describe CADD regions, which are characterized by a lack of observed  
81 variants with high functionally-Adjusted CADD Scores (ACS) in the gnomAD database (23). CADD regions  
82 are expected to reflect functional constraints. CADD regions present the key advantage of providing pre-  
83 defined and functionally-informed regions which can be used to construct PSAP null distributions.

84 We have made available a new implementation of the PSAP method using Snakemake (24)  
85 workflows, called Easy-PSAP (<https://github.com/msogloblinsky/Easy-PSAP>), which features null  
86 distributions constructed with up-to-date allele frequency data and pathogenicity scores. Here, we  
87 introduce PSAP-genomic-regions, an extension of the PSAP method to the non-coding genome by using  
88 the pre-defined CADD regions as testing unit instead of genes. This is an innovative strategy to prioritize  
89 variants at the scale of an individual genome. PSAP-genomic-regions is now available in Easy-PSAP. We

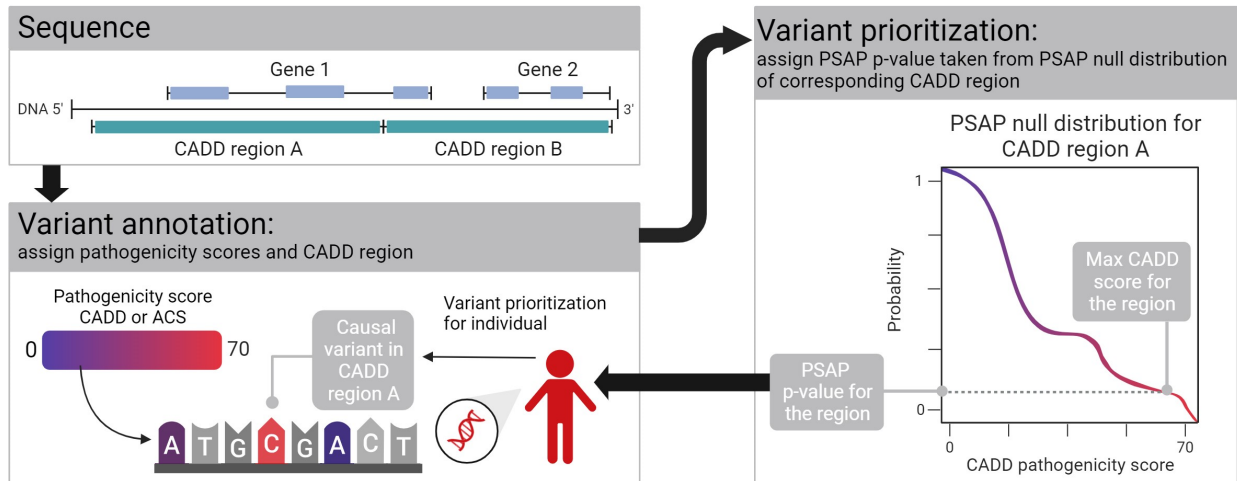
90 devised an evaluation protocol using artificially-generated disease exomes and genomes, obtained by  
91 inserting coding and non-coding ClinVar (25) variants in general population whole genomes from the 1000  
92 Genomes Project (26) and exomes from the FrEnch EXome (FREX) project (27). We show the consistent  
93 improvement in prioritization by using PSAP-genomic-regions over pathogenicity scores alone for non-  
94 coding and then coding variants. For coding variants, we also demonstrate the good performance of PSAP-  
95 genomic-regions compared to PSAP-genes. On real-life data, we illustrate the power of PSAP-genomic-  
96 regions on WES data from six resolved cases of Cerebral Small Vessel Disease (CSVD) and WGS data from  
97 three families affected by male infertility. These two diseases are particularly relevant to test our method,  
98 monogenic forms of CSVD (28) and male infertility (29) being extremely heterogeneous.

## 99 Results

### 100 Construction of PSAP null distribution in coding and non-coding regions

101 The idea behind the original PSAP method, referred to as PSAP-genes, relies on the calculation of  
102 gene-specific null distributions of CADD pathogenicity scores. More precisely, for an individual exome or  
103 genome and in a given gene, PSAP-genes considers the genotype with the highest CADD score and  
104 evaluates the probability to observe such a high CADD score in this gene in the general population. PSAP-  
105 genes deals separately with heterozygote and homozygote variants in the autosomal dominant (AD) and  
106 the autosomal recessive (AR) models respectively. Here, we will focus on homozygote variants for the  
107 recessive model. As a result, PSAP-genes gives a p-value to the genotype with the highest CADD score in  
108 the gene for each gene, model, and individual. PSAP can also score compound heterozygote variants, i.e.  
109 two heterozygote variants in the same gene, thus also giving a PSAP p-value to the genotype with the  
110 second highest CADD score in the gene. This p-value allows the ranking of the genes for an individual  
111 exome or genome. The PSAP principle can be generalized to any genomic unit.

112 Here, with PSAP-genomic-regions, we extended the PSAP method to analyze whole-genome data  
113 using predefined CADD regions as testing units instead of genes (Fig 1). The same principle as before is  
114 employed, with the difference being that the genotype with the highest CADD score in the region can be  
115 coding or non-coding. We thus constructed PSAP-genomic-regions null distributions using the CADD  
116 pathogenicity score (PHRED scaled across the whole genome). Our novel strategy will be referred to as  
117 PSAP-genomic-regions-CADD. We also explored the use of another pathogenicity score, the ACS (22)  
118 (PHRED scaled CADD scores by “coding”, “regulatory” and “intergenic” regions) to mitigate the higher  
119 CADD scores of coding variants (PSAP-genomic-regions-ACS strategy). The PSAP-genomic-regions  
120 strategies were compared to the initial PSAP-genes strategy, also referred to as PSAP-genes-CADD.



121

122 **Fig 1. Description of the PSAP-genomic-regions strategy.**

123 We calculated PSAP null distributions for SNVs in genes and CADD regions, in the hg19 and hg38  
 124 assemblies of the human genome. In hg19, PSAP null distributions were obtained for 19,283 genes and  
 125 119,695 CADD regions. In hg38 PSAP null distributions were obtained for 18,395 genes and 123,991 CADD  
 126 regions. PSAP null distributions and their parameters (unit of testing, allele frequencies and pathogenicity  
 127 score) can be found in S1 Table.

128

129 **Evaluating the performance of PSAP-genomic-regions on artificially-**  
 130 **generated disease exomes and genomes using ClinVar variants**

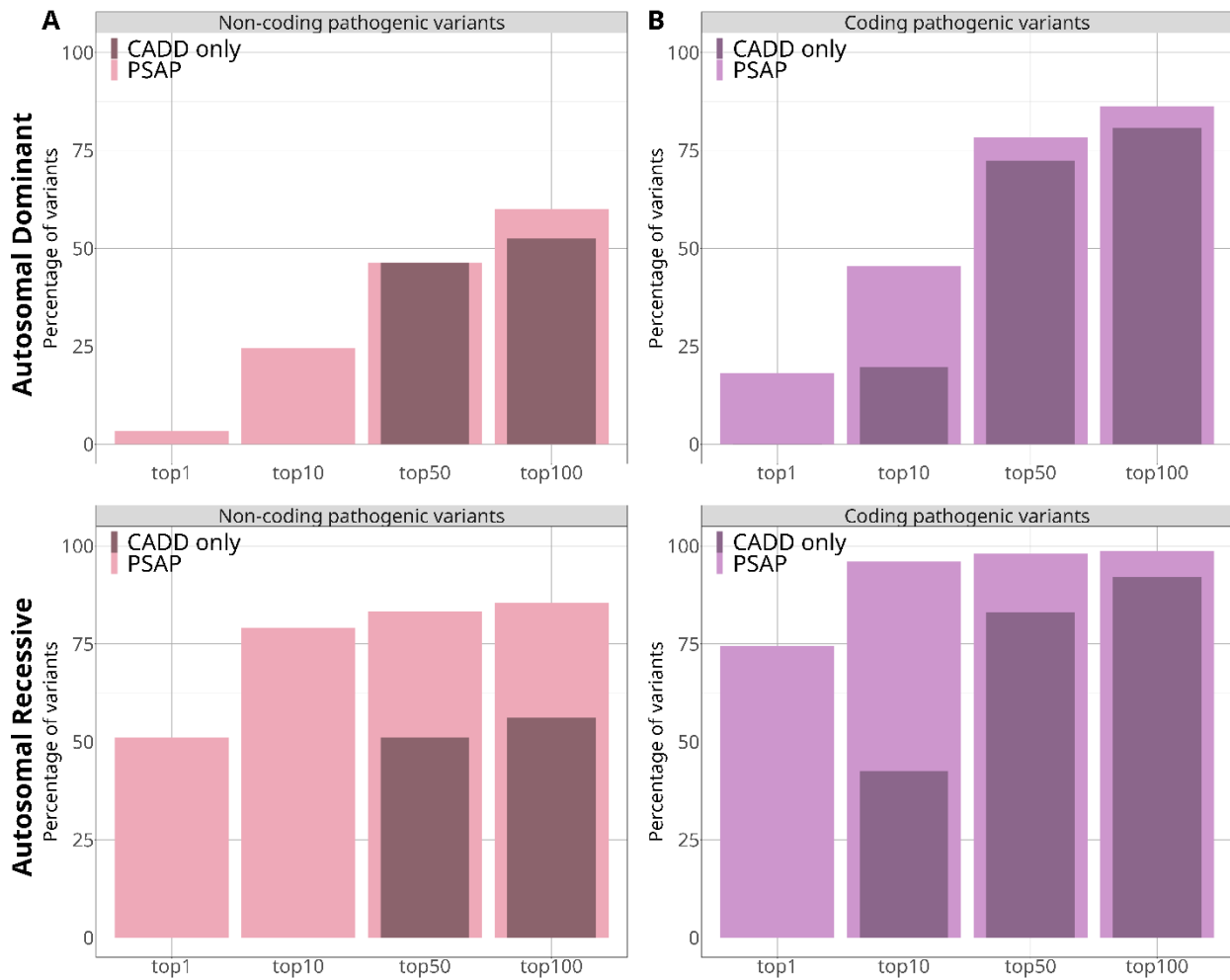
131 **Prioritization of non-coding pathogenic variants**

132 First, to evaluate how PSAP-genomic-regions performed to prioritize non-coding pathogenic variants,  
 133 we used artificially-generated disease genomes created by inserting non-coding ClinVar variants in the  
 134 Non-Finnish Europeans (NFE) from the 1000 Genomes Project phase 4 (NFE genomes) (see Material &  
 135 Methods and S1 File for the list of variants). Because the 1000 Genomes project is population-based, we



136 expect that some individuals might carry one or a few pathogenic variants in their genome. These  
137 pathogenic variants are characterized by a high CADD score and a low PSAP p-value. Thus, in order to  
138 summarize the rank of a ClinVar variant in an evaluation setting, we considered the best rank reached by  
139 the variant in at least 90% of the individuals.

140 Most of the NFE genomes carried a variant with a higher pathogenicity score or a lower PSAP p-value  
141 than most of the ClinVar variants (S1 Fig). We thus compared the percentage of the non-coding pathogenic  
142 variants ranked among the top N (N = 1, 10, 50 and 100) in at least 90% of the NFE genomes. The ranking  
143 at the individual level was done among all heterozygous variants for the ClinVar variants under the AD  
144 model, and across homozygous variants for the ClinVar variants under the AR model. Our main strategy  
145 PSAP-genomic-regions-CADD performed systematically better than using the CADD score alone (Fig 2A).  
146 The improvement was especially large for the top 10 ranking: 24.6% and 79.2% of ClinVar variants reached  
147 the top 10 with PSAP-genomic-regions-CADD for the AD and AR models, respectively, while no ClinVar  
148 variant reached the top 10 with CADD scores alone. For the prioritization of coding variants, the PSAP-  
149 genomic-regions-CADD strategy always outperformed PSAP-genomic-regions-ACS (S2 Fig).



150

151 **Fig 2. Comparison of the PSAP-genomic-regions-CADD strategy versus the CADD score alone in**

152 **artificially-simulated disease genomes.** Percentage of non-coding and coding pathogenic ClinVar variants

153 reaching the top N of variants in at least 90% of NFE genomes, with PSAP-genomic-regions (darker shade

154 of pink or purple) or the CADD score alone (lighter shade of pink or purple) (A) N = 175 non-coding AD

155 variants and N = 96 non-coding AR variants (B) N = 4,965 coding AD variants and N = 2,680 coding AR

156 variants.

157

158 Using the ACS scores improved the performance to detect non-coding-variants for the AD model (S2

159 Fig): 56.6% and 24.6% of variants reached the top 10 with PSAP-genomic-regions-ACS and PSAP-genomic-

160 regions-CADD, respectively. The gain in performance with PSAP-genomic-regions-ACS compared to PSAP-  
161 genomic-regions-CADD is not significant for the AR model for the top 10, top 50 and top 100. Nonetheless,  
162 we can note the pattern is different for the top 1 for the AR model: 51% with PSAP-genomic-regions-CADD  
163 to 5.5% with PSAP-genomic-regions-ACS. Indeed, switching from CADD score to ACS score has lowered  
164 the PSAP p-value of non-coding variants shared by more than 10% of NFE genomes. This led to a defect  
165 of the top rank reached by the ClinVar variants, as we considered the lowest rank reached in at least 90%  
166 of individuals. For instance, a variant in the CADD region R109138 shared by 70 of the NFE genomes went  
167 from a CADD score of 18.1 and a PSAP-genomic-regions-CADD p-value of 0.1 to an ACS of 22.2 and a PSAP-  
168 genomic-regions-ACS p-value of  $5.18 \times 10^{-10}$ . Thus, the ClinVar variants inserted in these individuals having  
169 a higher p-value than  $5.18 \times 10^{-10}$  do not rank first. Considering the overall more consistent performance  
170 of the PSAP-genomic-regions-CADD strategy, we chose to focus on this strategy, although we provide  
171 comparison with the PSAP-genomic-regions-ACS strategy which can have advantages for non-coding  
172 variants.

173 We further explored PSAP results for splicing ClinVar variants versus other type of non-coding ClinVar  
174 variants. Indeed, we observed that splicing variants are the major type of non-coding ClinVar variants.  
175 These splicing variants often had a very good ranking, especially with PSAP-genomic-regions-ACS (n=115  
176 splicing variants among 175 non-coding AD variants and n=72 splicing variants among 96 non-coding AR  
177 variants; S3 Table; Panel A in S3 Fig). Splicing ClinVar variants have a much higher ACS than CADD scores  
178 (Panel B in S3 Fig) which results in better ranking than for other types of non-coding ClinVar variants using  
179 PSAP-genomic-regions-ACS p-values (Panel C in S3 Fig). As a consequence, the percentage of splicing  
180 ClinVar variants ranked in the top 10 was largely improved when using PSAP-genomic-regions-ACS, for the  
181 AD model especially which was less powerful with PSAP-genomic-regions-CADD to begin with (Panel D in  
182 S3 Fig).

183 The full results of ranking by PSAP-genomic-regions-CADD and PSAP-genomic-regions-ACS for the  
184 non-coding non-splicing pathogenic ClinVar variants can be found in S2 File. With PSAP-genomic-regions-  
185 ACS, around half of the non-coding non-splicing variants are ranked in the top 50 and top 10 of variants  
186 for more than 90% of NFE genomes for the AD and AR models, respectively (19 out of 45 variants for the  
187 AD model and 8 out of 21 variants for the AR model). The other half of variants present a less significant  
188 PSAP-genomic-regions-ACS p-value and a poorer ranking. PSAP-genomic-regions-CADD achieves a similar  
189 ranking of AR non-coding non-splicing variants (7 out of 21 variants) but a decreased prioritization for AD  
190 non-coding non-splicing variants (6 out of 45 variants). To confirm this pattern of ranking for non-coding  
191 non-splicing pathogenic variants on another set of variants, we evaluated with our artificially generated  
192 disease genomes protocol 320 non-coding SNVs used to train Genomiser (30). These variants were not  
193 associated with a mode of inheritance. Hence, we inserted them in the NFE genomes and scored them  
194 with both AD and AR PSAP-genomic-regions-ACS null distributions. Among the 320 non-coding variants,  
195 169 reached the top 100 in at least 90% of NFE genomes, with either the AD or AR model (S3 File). This  
196 can be explained by the distributions of CADD scores compared to ACS scores for the ClinVar variants: the  
197 non-coding variants that do not reach the top 100 have a significantly lower CADD and ACS scores  
198 compared to all the other types of variants (S4 Fig). Overall, PSAP-genomic-regions prioritizes around half  
199 of non-coding ClinVar and Genomiser training variants in the top 100 of NFE genomes. The ones who have  
200 a higher ranking present much lower CADD and ACS scores and would never be well-ranked by any PSAP  
201 strategy.

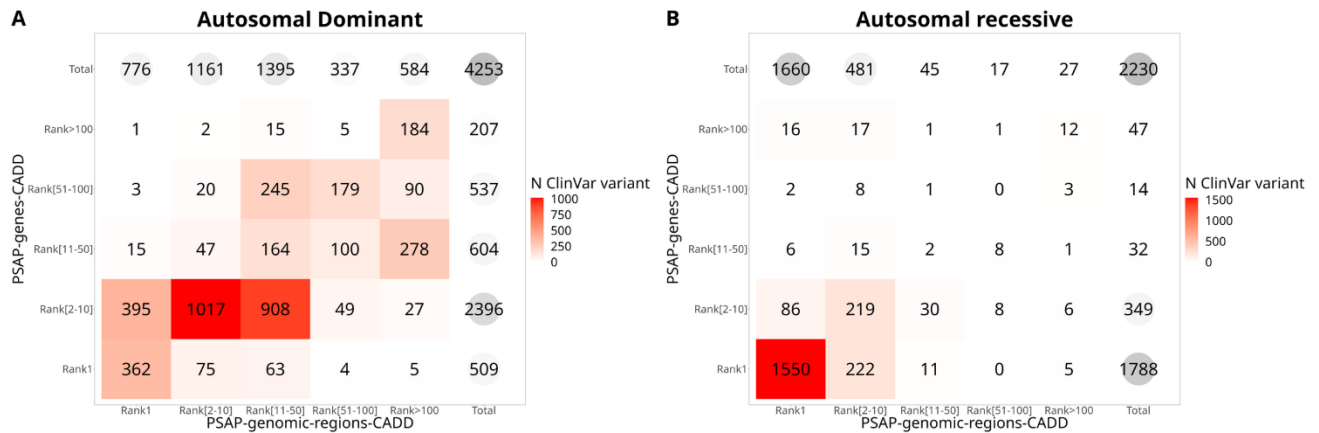
202 PSAP-genomic-region is also relevant for the analysis of exome data. Indeed, exome sequencing  
203 captures variants outside of the bounds of coding regions (31), such as intronic variants. We explored the  
204 prioritization of non-coding ClinVar variants located within the WES-targeted regions of the FREX  
205 individuals using our artificially-generated disease exomes protocol (N=48 variants for the AD model and  
206 N=64 variants for the AR model, Panel A in S5 Fig). For both PSAP-genomic-regions-CADD and PSAP-

207 genomic-regions-ACS, there was a large increase in prioritization performance compared to using only the  
208 pathogenicity scores. Because there are fewer variants in an exome background than in a genome  
209 background, the rankings of these non-coding ClinVar variants were better in FREX than in NFE genomes.  
210 The best ranking was achieved using PSAP-genomic-regions-ACS, with 82% and 90.3% of variants reaching  
211 the top 10 for the AD and AR models, respectively, whilst PSAP-genomic-regions-CADD achieved a similar  
212 ranking for AR variants. Most of these non-coding pathogenic variants were splicing variants (40 out of 73  
213 variants for the AD model and 56 out of 64 variants for the AR model), and half of them were considered  
214 as having a functional “HIGH IMPACT” (26 variants for the AD model and 22 variants for the AR model).  
215 Hence, prioritizing variants with PSAP-genomic-regions allows identifying more variants even in exome  
216 data, that are in addition functionally-relevant.

217

## 218 [Prioritization of coding pathogenic variants](#)

219 Similar evaluations were performed for ClinVar coding variants inserted in either WGS from  
220 1000G NFE individuals or WES from FREX. As observed for non-coding pathogenic variants, PSAP-genomic-  
221 regions outperformed the pathogenicity scores alone (Fig 2B, Panel B in S5 Fig). However, in the context  
222 of coding pathogenic ClinVar variants, we observed that the strategy of PSAP-genomic-regions-CADD  
223 provided better prioritization compared with the PSAP-genomic-regions-ACS strategy. We observed that  
224 18.2% and 74.6% of the coding variants reached the top 1 in at least 90% of genomes backgrounds with  
225 the PSAP-genomic-regions-CADD for the AD and AR model respectively, against no variants with the CADD  
226 score alone, and against 5.3% and 2.5% reaching the top 1 with PSAP-genomic-regions-ACS. In the exome  
227 background and with PSAP-genomic-regions-CADD, 38.7% and 89.8% of AD variants reached the top 1  
228 and top 50, respectively; 80.3% and 97.9% of AR variants reached the top 1 and the top 50, respectively.



229

230 **Fig.3. Comparison of PSAP-genomic-regions-CADD and PSAP-genes-CADD strategies in artificially-**  
 231 **simulated disease genomes.** Number of coding pathogenic ClinVar variants reaching rank [x-y] of variants  
 232 in at least 90% of 1000 Genomes Project NFE individuals for each strategy.

233

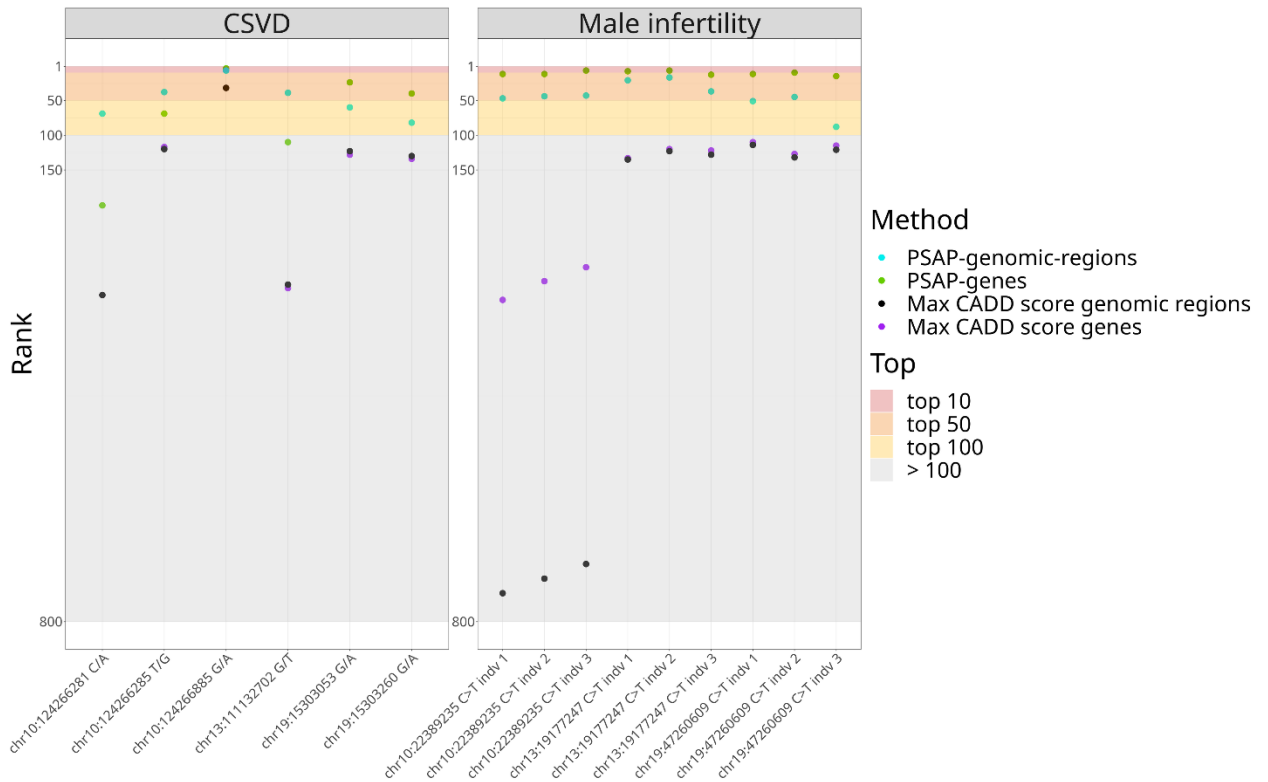
234 We also compared the number of coding ClinVar variants reaching the tops in NFE genomes between  
 235 PSAP-genomic-regions-CADD strategy and the initial PSAP-genes-CADD strategy (Fig 3). More differences  
 236 were observed across the two PSAP strategies for the AD than for the AR model (Fig 3A). There were 362  
 237 variants ranked first and 1,017 variants ranked [2-10] in common between the two strategies. However,  
 238 908 variants that were ranked [2-10] with PSAP-genes-CADD were [11-50] with PSAP-genomic-regions-  
 239 CADD, and 395 variants that were ranked [2-10] with PSAP-genes-CADD were ranked first with PSAP-  
 240 genomic-regions-CADD. Regarding variants that are ranked more than a 100 with PSAP-genomic-regions-  
 241 CADD, 278 of them are ranked [11-50] and 90 are ranked [51-100] by PSAP-genes-CADD. Regarding the  
 242 AR model (Fig 3B), PSAP-genomic-regions-CADD performed similarly to PSAP-genes-CADD, and the  
 243 majority of variants were ranked first with both strategies (1,550 variants). Even more promising results  
 244 can be found when looking at the same comparison of ranks within the FREX exomes (S6 Fig). For instance,  
 245 in the AD model, 592 variants that were ranked [2-10] with PSAP-genes-CADD are ranked first with PSAP-

246 genomic-regions-CADD, against 115 variants ranked [2-10] with PSAP-genomic-regions-CADD that  
247 become first with PSAP-genes-CADD.

248

## 249 **Application of PSAP-genomic-regions to real data with different modes** 250 **of inheritance**

251 To illustrate our method in real-life settings, we analyzed two datasets (S4 Table), one with an AD  
252 mode of inheritance and the other with an AR mode of inheritance. The first dataset consisted of WES  
253 data for six individuals affected by monogenic forms of CSVD (32). Using PSAP-genomic-regions-CADD, all  
254 of the causal variants were ranked at least in the top 100 in each patient (Fig 4). The contribution of CADD  
255 regions as a unit of testing was especially visible for the variant in *COL4A2* and one variant in *HTRA1* which  
256 were not well-ranked using genes as testing unit (rank 110 and 193 respectively with genes, and rank 3  
257 and 69 with CADD regions). Using their maximal CADD score by gene or CADD region alone, these variants  
258 would not have been prioritized in the top 100 for five out of six individuals.



259

260 **Fig. 4. Prioritization of 6 known CSVD mutations and 3 male infertility candidate variants with PSAP-**  
 261 **genomic-regions-CADD, PSAP-genes-CADD and the maximal CADD score on genes or CADD regions.**

262

263 The second dataset consisted of WGS data for 9 individuals from three families with clinically  
 264 diagnosed male infertility (33). All causal variants fell within the top 20 of variants with prioritization by  
 265 PSAP-genes-CADD, and within the top 50 for at least one case per family with PSAP-genomic-regions-  
 266 CADD (within top 100 for all cases, Fig 4). PSAP-genomic-regions-CADD did not improve the ranking of  
 267 these coding variants, which was expected considering the large number of variants in a WGS analysis  
 268 (see S4 Table for the total number of variants in each analysis). The prioritization from PSAP-genomic-  
 269 regions-CADD was still interesting to narrow the set of candidates for causal variants. In clinics when the  
 270 CADD score alone is used, these variants would not have been prioritized (CADD score  
 271 < 25, and rank > 100 with the maximal CADD score strategy). PSAP-genomic-regions-CADD thus allow a



272 relevant prioritization of coding pathogenic variants in WGS sequencing and an unbiased exploratory  
273 analysis at the scale of the whole genome.

274 Using PSAP-genomic-regions-ACS or the ACS score alone, almost all of the CSVD and male infertility  
275 coding pathogenic variants had a rank greatly exceeding the top 100 (S4 Table). The only exception is one  
276 variant in *HTRA1* (10:124266885 G/A) that was ranked 3 by PSAP-genomic-regions-ACS and 10 by the  
277 maximal ACS score alone. This *HTRA1* variant was a splicing variant, which confirms the good performance  
278 of the PSAP-genomic-regions-ACS strategy on this type of variant.

## 279 Discussion

280 Variant prioritization, especially in the case of very heterogeneous rare diseases, is a clinically-  
281 relevant methodological challenge for both clinicians and researchers. Mounting evidence suggests that  
282 current methods of analysis and their restriction to the coding genome are a hindrance to the discovery  
283 of new genetic variants implicated in rare diseases (16). We have developed PSAP-genomic-regions, an  
284 extension of the PSAP method to the whole genome using functionally-relevant genomic regions. PSAP-  
285 genomic-regions broadens the scope of variants evaluated by PSAP and addresses the issue of variant  
286 prioritization at an individual whole-genome scale.

287 PSAP-genomic-regions has been thoroughly tested and validated by using simulations emulating real-  
288 life scenarios of causal variant prioritization. PSAP-genomic-regions achieves a prioritization of coding  
289 pathogenic SNVs in the top 100 variants of an exome or genome which is a relevant number of variants  
290 to analyze for clinicians. Without use of prior knowledge on the disease, PSAP-genomic-regions achieves  
291 relevant variant prioritization within millions of variants to analyze, which is illustrated by the ranking of  
292 6 variants involved in CSVD and 3 variants involved in familial cases of male infertility in the top 100 of  
293 WES and WGS data respectively. PSAP-genomic-regions thus helps with the diagnosis of such  
294 heterogeneous diseases in conjunction with other relevant information like the mode of transmission,  
295 prevalence or type of variant involved.

296 PSAP-genomic-regions also allows the scoring of variants otherwise discarded from the analysis, like  
297 splicing variants with a high predicted functional impact, and other non-coding variants of proven clinical  
298 significance. The only scenario for which PSAP-genomic-regions is not advantageous compared to the  
299 PSAP-genes strategy is for prioritizing coding variants in WGS data. In that case, using coding CADD  
300 regions, i.e. the coding parts of CADD regions for the analysis still yields better results compared to PSAP-  
301 genes (S7 Fig). Our simulations using known pathogenic variants have shown which PSAP strategy

302 performs the best depending on the type of data and variant expected to be involved in the disease  
303 mechanism. If there is no expected type of variant: we advise on the use of the PSAP-genomic-regions-  
304 CADD strategy, which gives the overall best results. For coding variants prioritization specifically, PSAP-  
305 genomic-regions-CADD gives the best results in WES, and PSAP-coding-genomic-regions-CADD performs  
306 best in WGS data. Finally, if no coding variant of interest for the disease is found with PSAP-genomic-  
307 regions-CADD or PSAP-coding-genomic-regions-CADD, PSAP-genomic-regions-ACS can be applied to look  
308 for non-coding variants of interest especially for an AD expected model of transmission.

309 To the best of our knowledge, there is no other score of predicted pathogenicity for all possible SNVs  
310 comparable to CADD. The main pathogenicity prediction scores developed to date were described and  
311 compared in a recent review (34). Multiple benchmarks on the subject show conflicting conclusions  
312 depending on the variant testing set (35,36). A significant limitation of some of the most popular tools,  
313 such as SIFT (37), PolyPhen-2 (38), VEST (39), and REVEL (40), is their restriction to analyzing only missense  
314 variants. In contrast, CADD stands out as it is a meta-predictor, integrating scores from SIFT, PolyPhen-2,  
315 phyloP (41), and GERP (42), and enabling the scoring of any SNV with pre-computed scores, and any InDel  
316 in the genome with on-request scores. Additionally, CADD is trained on a much larger number of variants  
317 compared to other machine-learning methods, while using a relatively modest number of features. Similar  
318 types of methods aim at prioritizing more constrained regions in the non-coding genome (18,20) or  
319 distinguishing deleterious non-coding variants from neutral ones (18,43). However, most of these  
320 prediction methods either do not provide a variant-specific score, or are not defined in both coding and  
321 non-coding parts of the genome. Other well-known methods for identification of pathogenic variants in  
322 exome and genome data rely on the use of HPO terms to make a prediction, like Exomiser (44) or  
323 Genomiser (30), making in comparison PSAP-genomic-regions an unmatched prioritization tool. As any  
324 other bioinformatics variant prioritization method, it has to be used in conjunction with other lines of  
325 evidence like the expression of the associated gene in a tissue of interest or segregation of variants if

326 familial data is available to ultimately lead to any genetic diagnosis of a patient. PSAP-genomic-regions  
327 does not make assumption on the type of variants and does explore the whole genome. The ranking by  
328 p-values coming from the application of PSAP-genomic-regions to an individual's variants is a useful way  
329 to narrow-down the list of variants to further investigate for both researchers and clinicians in different  
330 scenarios. Other criteria for variant filtering will depend heavily on the type of disease studied. The clinical  
331 interpretation of pathogenicity for non-coding variants is more challenging than for coding variants and  
332 can be improved by applying modified ACMG guidelines which can help pick out potentially candidate  
333 regulatory elements to explain the phenotype of the patient (45).

334 The method most comparable to the strategy followed by PSAP-genomic-regions is the recently-  
335 developed machine-learning algorithm FINSURF (46). FINSURF aims to predict the functional impact of  
336 non-coding variants in regulatory regions and has been applied to known pathogenic variants inserted in  
337 WGS data like we did. Nonetheless it has been difficult to compare properly the two methods considering  
338 FINSURF only scores non-coding variants in predefined regulatory regions, and the set of variants used to  
339 train the method is not available.

340 The main limitation of PSAP-genomic-regions comes from the score used to calibrate null  
341 distributions, namely the CADD score. We have observed that known pathogenic non-coding ClinVar  
342 variants that were not well-ranked by PSAP-genomic-regions had significantly lower CADD and ACS scores  
343 compared to splicing and better-ranked non-coding variants. Because such CADD score is likely to be seen  
344 in the general population, PSAP-genomic-regions will not be able to prioritize such a variant with at a low  
345 rank. We also observed that some CADD regions were badly-calibrated and resulted in the assignment of  
346 very low PSAP-genomic-regions p-values to putatively neutral variants in the 1000 Genomes Project. As  
347 allele frequencies from larger databases and more accurate pathogenicity scores become available, this  
348 will lead to an improvement of the PSAP method as well. The most recent release of the CADD score v1.7

349 (47) notably integrates regulatory annotations and may further improve the prioritization of non-coding  
350 pathogenic variants when integrated in PSAP-genomic-regions.

351 Many avenues of further development and improvement are open for PSAP-genomic-regions,  
352 including the inclusion and scoring of InDel variations and structural variants. Exploring the combination  
353 of the PSAP-genomic-regions p-values with other metrics or information coming from omics analysis could  
354 also improve prediction. Finally, the flexibility of the PSAP method makes it potentially adaptable to other  
355 more complex models like digenic and oligogenic models of inheritance, considering the increasing  
356 availability of information coming from gene networks and biological pathways.

357

358

## 359 **Materials and Methods**

### 360 **Construction of PSAP null distributions**

361 The first parameter is the units in which to construct the PSAP null distribution. Here we considered  
362 two unit strategies: the genes and the CADD regions (S1 Table). For the genes, the coding regions of genes  
363 were defined based on the biomaRt R package: the gene coding sequences were retrieved from Ensembl  
364 (48) by requesting the “genomic\_coding\_start” and “genomic\_coding\_end”, on both the hg19 and hg38  
365 builds. To account for splicing regions, the coding regions were extended by two bases on both sides of  
366 the gene coding regions. In total, 19,780 genes were retrieved in hg19 and 23,163 in the hg38 build. For  
367 the CADD regions, their coordinates were downloaded from <https://lysine.univ-brest.fr/RAVA-FIRST/> for  
368 the hg19 build and were lifted over to hg38 using the Ensembl Assembly Converter. CADD regions  
369 coordinates in hg38 are available on Easy-PSAP GitHub (<https://github.com/msogloblinsky/Easy-PSAP>).  
370 There were 135,224 CADD regions in hg19 and 131,970 in hg38. For the coding CADD regions, i.e. the

371 coding parts of CADD regions, we considered the intersection of the CADD regions and the gene coding  
372 regions for each build, which yielded 37,978 coding CADD regions in hg19 and 52,340 in hg38.

373 The second parameter is the allele frequencies database. Here we considered the global allele  
374 frequencies from the gnomAD database to calibrate the PSAP null distributions: gnomAD genome r2.0.1  
375 for hg19 and gnomAD V3 (49) for hg38. For our purpose, we considered only single nucleotide variants  
376 (SNVs) annotated as PASS by the Variant Quality Score Recalibration (VQSR) of GATK (50) and located in  
377 well-covered regions. Well-covered regions in gnomAD genome were defined as regions for which 90% of  
378 individuals have coverage at depth 10. Variants not seen in gnomAD genome, not annotated as PASS or  
379 not located in well-covered regions (gnomAD genome version according to the build) have a frequency of  
380 0 and thus did not contribute to the construction of the null distributions.

381 To ensure reliability of PSAP null distribution, it is crucial that the units are well covered in the  
382 database from which the allele frequencies are taken. Thus, we only considered units for which at least  
383 half of the unit was well-covered (as defined previously) in gnomAD genome (version according to the  
384 build). Coding regions of genes and well-covered regions in gnomAD genome were intersected to get the  
385 percentage of each gene's coding regions that were well-covered in the database. The same steps were  
386 carried out with CADD regions as genomic units for PSAP, for hg19 and hg38 builds. PSAP null distributions  
387 were thus constructed for 19,283 and 18,395 genes in hg19 and hg38 respectively, 119,695 and 123,991  
388 CADD regions, and 34,397 and 35,226 coding CADD regions in hg19 and hg38 respectively.

389 The third parameter is the pathogenicity score. Here, for the evaluation of PSAP on coding variants,  
390 we used the version 1.6 of CADD (51) for each build, accessible on the CADD website  
391 (<https://cadd.gs.washington.edu/>). For the evaluation on non-coding variants, which tend to have lower  
392 CADD scores than coding variants (52), we followed the strategy described in Bocher et al.(22) to adjust

393 the RAW CADD score v1.6 of all possible SNVs on a PHRED scale stratifying by type of genomic regions:  
394 “coding”, “regulatory” and “intergenic”, resulting in “adjusted CADD scores”, referred to as “ACS”.

395 Easy-PSAP (<https://github.com/msogloblinsky/Easy-PSAP>) was used to generate null distributions  
396 according to the previously described input files and parameters. This resulted in 4 sets of null  
397 distributions for the AD and AR models for both hg19 and hg38 assemblies (S1 Table).

398

## 399 **Evaluating the performance of PSAP-genomic-regions using artificially-** 400 **generated disease exomes and genomes**

401 To evaluate the ability of PSAP-genomic-regions to prioritize known pathogenic variants in an  
402 individual, we leveraged artificially-generated disease exomes and genomes using available general  
403 population cohorts. These different PSAP strategies (see Table 1) were compared in terms of their  
404 performances to prioritize the known pathogenic variants.

405 The pathogenic ClinVar (25) SNVs with coordinates in hg19 and hg38 were downloaded from the NCBI  
406 website (<https://www.ncbi.nlm.nih.gov/clinvar/>, accessed on the 3rd of June 2022). Some of these ClinVar  
407 variants had an annotated mode of inheritance ("moi autosomal recessive" and "moi autosomal  
408 dominant"). From ClinVar, there were 14,056 variants annotated as AD and 12,758 variants annotated as  
409 AR. Variants were filtered out to keep only autosomal pathogenic SNVs having as review status either  
410 “reviewed by expert panel” or “criteria provided, multiple submitters, no conflicts”, which are the two  
411 best review status in ClinVar. There were 1,518 AD and 1,118 AR variants meeting these criteria.

412 For variants which did not have an annotated mode of inheritance, we used a curated version of the  
413 database OMIM, hOMIM (53) to retrieve a mode of inheritance, and kept variants that were always  
414 associated with an AD or AR mode of inheritance in hOMIM. The same filtering was applied, which left

415 3,641 additional variants for the AD and 1,706 for the AR model. In total, we had a set of 5,159 variants  
416 for the AD model and 2,824 variants for the AR model. Among these ClinVar variants, 4,965 and 2,680  
417 variants were coding SNVs respectively for the AD and AR models. Similarly, 175 and 96 variants were  
418 non-coding variants for the AD model and AR models, among which 48 variants for the AD model and 64  
419 for AR model fell within the boundaries covered by FREX exomes. The list of pathogenic ClinVar variants  
420 and their mode of inheritance can be found in S1 File.

421 We inserted each variant from our curated list of pathogenic ClinVar variants successively in each of  
422 the 533 high coverage NFE genomes and each of the 574 exomes from the FREX project. An individual-  
423 focused QC was applied on both datasets using the RAVAQ R package (54): we performed a genotype and  
424 variant QC with default parameters corresponding to standard GATK hard filtering criteria, mean allele  
425 balance computed across heterozygous genotypes and call rates, except for MAX\_AB\_GENO\_DEV = 0.25,  
426 MAX\_ABHET\_DEV, MIN\_CALLRATE and MIN\_FISHER\_CALLRATE "disabled".

427 We conducted the artificially-generated disease genome and exome evaluation with PSAP null  
428 distributions in hg19 and hg38 respectively, to match with the build of the data. We then applied the 3  
429 PSAP strategies mentioned previously (PSAP-genes-CADD, PSAP-genomic-regions-CADD and PSAP-  
430 genomic-regions-ACS). For each strategy, we kept the maximal pathogenicity score (CADD or ACS) for each  
431 unit (gene or CADD regions) and then ranked the units according to their PSAP p-value or to their  
432 pathogenicity score alone within each genome or exome. We compared the PSAP-genes-CADD and PSAP-  
433 genomic-regions-CADD strategies to using the maximal CADD score alone by gene or CADD regions,  
434 respectively; and the PSAP-genomic-regions-ACS strategy to using the maximal ACS score by CADD region.  
435 For each ClinVar variant, we retrieved its rank within each genome or exome. Coding ClinVar variants were  
436 evaluated with the 3 PSAP strategies whereas non-coding ClinVar variants were evaluated with the novel  
437 PSAP-genomic-regions-CADD and PSAP-genomic-regions-ACS strategies (see S2 Table for more details).



438

## 439 Patient data analysis

440 The PSAP strategies were applied to real WES data from six unrelated patients affected by a CSVD for  
441 which the causal variant is known, which allowed a comparison of performance between the different  
442 strategies. The full description of the dataset can be found in [Aloui et al. 2021] (32), with the exception  
443 of the QC process. For this analysis, the same QC as for the FREX and 1000 Genomes Project datasets was  
444 performed. We applied PSAP-genes-CADD and PSAP-genomic-regions-CADD in hg19 to the six resolved  
445 CSVD patients' exome data. The other PSAP parameters were the ones by default as described previously.  
446 Two of the individuals had a causal pathogenic variant in the gene *NOTCH3* (19:15303053 G/A and  
447 19:15303260 G/A), one individual in the gene *COL4A2* (13:111132702 G/T) and three individuals in the  
448 gene *HTRA1* (10:124266285 T/G, 10:124266281 C/A and 10:124266885 G/A). The rank of the known CSVD  
449 variants among other heterozygote variants in the patient's exome according to its PSAP p-value for the  
450 2 strategies was then retrieved.

451 The PSAP strategies were also applied to WGS data of three families with clinically diagnosed forms  
452 of male infertility (33) and for which a pathogenic recessive variant was prioritized using a computational  
453 pipeline featuring the initial PSAP-genes implementation. Three affected individuals were analyzed for  
454 each family. The description of the whole dataset and candidate variant filtering process can be found in  
455 [Khan and Akbari et al. 2023] (33), except for the QC that was performed in the same way as for the CSVD  
456 data. Two other families were resolved from the same dataset, but considering that the causal variants  
457 were deletions we did not include them in the current analysis. The prioritized pathogenic variants were  
458 in the genes: *SPAG6* (chr10:22389235 C/T) for family 3, *TUBA3C* (chr13:19177247 C/T) for family 7 and  
459 *CCDC9* (chr19:47260609 C/T) for family 4. We applied PSAP-genes-CADD and PSAP-genomic-regions-

460 CADD in hg38 to the 9 cases and retrieved the rank of the known male infertility variants among other  
461 homozygote variants in the patient's genomes according to its PSAP p-value for the 2 strategies.

## References

- 463 1. Sequeira AR, Mentzakis E, Archangelidi O, Paolucci F. The economic and health impact of rare  
464 diseases: A meta-analysis. *Health Policy and Technology*. 2021 Mar 1;10(1):32–44.
- 465 2. Amberger J, Bocchini CA, Scott AF, Hamosh A. McKusick's Online Mendelian Inheritance in Man  
466 (OMIM®). *Nucleic Acids Research*. 2009 Jan 1;37(suppl\_1):D793–6.
- 467 3. Amberger JS, Bocchini CA, Schiettecatte F, Scott AF, Hamosh A. OMIM.org: Online Mendelian  
468 Inheritance in Man (OMIM®), an online catalog of human genes and genetic disorders. *Nucleic Acids  
469 Research*. 2015 Jan 28;43(D1):D789–98.
- 470 4. Ehrhart F, Willighagen EL, Kutmon M, van Hoften M, Curfs LMG, Evelo CT. A resource to explore the  
471 discovery of rare diseases and their causative genes. *Sci Data*. 2021 May 4;8(1):124.
- 472 5. Wright CF, FitzPatrick DR, Firth HV. Paediatric genomics: diagnosing rare disease in children. *Nat Rev  
473 Genet*. 2018 May;19(5):253–68.
- 474 6. Boycott KM, Rath A, Chong JX, Hartley T, Alkuraya FS, Baynam G, et al. International Cooperation to  
475 Enable the Diagnosis of All Rare Genetic Diseases. *The American Journal of Human Genetics*. 2017  
476 May 4;100(5):695–705.
- 477 7. Chong JX, Buckingham KJ, Jhangiani SN, Boehm C, Sobreira N, Smith JD, et al. The Genetic Basis of  
478 Mendelian Phenotypes: Discoveries, Challenges, and Opportunities. *The American Journal of Human  
479 Genetics*. 2015 Aug 6;97(2):199–215.
- 480 8. Wilfert AB, Chao KR, Kaushal M, Jain S, Zöllner S, Adams DR, et al. Genomewide significance testing  
481 of variation from single case exomes. *Nat Genet*. 2016 Dec;48(12):1455–61.
- 482 9. Kircher M, Witten DM, Jain P, O'Roak BJ, Cooper GM, Shendure J. A general framework for  
483 estimating the relative pathogenicity of human genetic variants. *Nat Genet*. 2014 Mar;46(3):310–5.
- 484 10. Wyrwoll MJ, Temel ŞG, Nagirnaja L, Oud MS, Lopes AM, van der Heijden GW, et al. Bi-allelic  
485 Mutations in M1AP Are a Frequent Cause of Meiotic Arrest and Severely Impaired Spermatogenesis  
486 Leading to Male Infertility. *The American Journal of Human Genetics*. 2020 Aug 6;107(2):342–51.
- 487 11. Kasak L, Punab M, Nagirnaja L, Grigorova M, Minajeva A, Lopes AM, et al. Bi-allelic Recessive Loss-  
488 of-Function Variants in FANCM Cause Non-obstructive Azoospermia. *The American Journal of  
489 Human Genetics*. 2018 Aug 2;103(2):200–12.
- 490 12. Salas-Huetos A, Tüttelmann F, Wyrwoll MJ, Kliesch S, Lopes AM, Conçaves J, et al. Disruption of  
491 human meiotic telomere complex genes TERB1, TERB2 and MAJIN in men with non-obstructive  
492 azoospermia. *Hum Genet*. 2021 Jan;140(1):217–27.
- 493 13. Kasak L, Rull K, Yang T, Roden DM, Laan M. Recurrent Pregnancy Loss and Concealed Long-QT  
494 Syndrome. *J Am Heart Assoc*. 2021 Aug 16;10(17):e021236.

- 495 14. Bustamante-Marin XM, Horani A, Stoyanova M, Charng WL, Bottier M, Sears PR, et al. Mutation of  
496 CFP57, a protein required for the asymmetric targeting of a subset of inner dynein arms in  
497 *Chlamydomonas*, causes primary ciliary dyskinesia. *PLoS Genet*. 2020 Aug 7;16(8):e1008691.
- 498 15. Hindorff LA, Sethupathy P, Junkins HA, Ramos EM, Mehta JP, Collins FS, et al. Potential etiologic and  
499 functional implications of genome-wide association loci for human diseases and traits. *Proceedings*  
500 *of the National Academy of Sciences*. 2009 Jun 9;106(23):9362–7.
- 501 16. Posey JE. Genome sequencing and implications for rare disorders. *Orphanet Journal of Rare*  
502 *Diseases*. 2019 Jun 24;14(1):153.
- 503 17. Buniello A, MacArthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangone C, et al. The NHGRI-EBI  
504 GWAS Catalog of published genome-wide association studies, targeted arrays and summary  
505 statistics 2019. *Nucleic Acids Res*. 2019 Jan 8;47(Database issue):D1005–12.
- 506 18. Gussow AB, Copeland BR, Dhindsa RS, Wang Q, Petrovski S, Majoros WH, et al. Orion: Detecting  
507 regions of the human non-coding genome that are intolerant to variation using population genetics.  
508 *PLOS ONE*. 2017 Aug 10;12(8):e0181604.
- 509 19. Huang YF, Gulko B, Siepel A. Fast, scalable prediction of deleterious noncoding variants from  
510 functional and population genomic data. *Nat Genet*. 2017 Apr;49(4):618–24.
- 511 20. Vitsios D, Dhindsa RS, Middleton L, Gussow AB, Petrovski S. Prioritizing non-coding regions based on  
512 human genomic constraint and sequence context with deep learning. *Nat Commun*. 2021 Mar  
513 8;12(1):1504.
- 514 21. Robinson PN, Köhler S, Bauer S, Seelow D, Horn D, Mundlos S. The Human Phenotype Ontology: A  
515 Tool for Annotating and Analyzing Human Hereditary Disease. *Am J Hum Genet*. 2008 Nov  
516 17;83(5):610–5.
- 517 22. Bocher O, Ludwig TE, Oglobinsky MS, Marenne G, Deleuze JF, Suryakant S, et al. Testing for  
518 association with rare variants in the coding and non-coding genome: RAVA-FIRST, a new approach  
519 based on CADD deleteriousness score. *PLOS Genetics*. 2022 Sep 16;18(9):e1009923.
- 520 23. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, et al. The mutational constraint  
521 spectrum quantified from variation in 141,456 humans. *Nature*. 2020 May;581(7809):434–43.
- 522 24. Köster J, Rahmann S. Snakemake—a scalable bioinformatics workflow engine. *Bioinformatics*. 2012  
523 Oct 1;28(19):2520–2.
- 524 25. Landrum MJ, Lee JM, Benson M, Brown GR, Chao C, Chitipiralla S, et al. ClinVar: improving access to  
525 variant interpretations and supporting evidence. *Nucleic Acids Res*. 2018 Jan 4;46(Database  
526 issue):D1062–7.
- 527 26. Auton A, Abecasis GR, Altshuler DM, Durbin RM, Abecasis GR, Bentley DR, et al. A global reference  
528 for human genetic variation. *Nature*. 2015 Oct;526(7571):68–74.

- 529 27. Génin E, Redon R, Deleuze J, Campion D, Lambert J, Dartigues J, et al. The French Exome (FREX)  
530 Project: A Population-based panel of exomes to help filter out common local variants. *Genetic*  
531 *Epidemiology*. 2017;41(7):691–691.
- 532 28. Rannikmäe K, Henshall DE, Thrippleton S, Ginja Kong Q, Chong M, Grami N, et al. Beyond the Brain.  
533 *Stroke*. 2020 Oct;51(10):3007–17.
- 534 29. Houston BJ, Riera-Escamilla A, Wyrwoll MJ, Salas-Huetos A, Xavier MJ, Nagirnaja L, et al. A  
535 systematic review of the validated monogenic causes of human male infertility: 2020 update and a  
536 discussion of emerging gene–disease relationships. *Human Reproduction Update*. 2022 Feb  
537 1;28(1):15–29.
- 538 30. Smedley D, Schubach M, Jacobsen JOB, Köhler S, Zemojtel T, Spielmann M, et al. A Whole-Genome  
539 Analysis Framework for Effective Identification of Pathogenic Regulatory Variants in Mendelian  
540 Disease. *Am J Hum Genet*. 2016 Sep 1;99(3):595–606.
- 541 31. Guo Y, Long J, He J, Li CI, Cai Q, Shu XO, et al. Exome sequencing generates high quality data in non-  
542 target regions. *BMC Genomics*. 2012 May 20;13:194.
- 543 32. Aloui C, Hervé D, Marenne G, Savenier F, Le Guennec K, Bergametti F, et al. End-Truncated LAMB1  
544 Causes a Hippocampal Memory Defect and a Leukoencephalopathy. *Annals of Neurology*.  
545 2021;90(6):962–75.
- 546 33. Khan MR, Akbari A, Nicholas TJ, Castillo-Madeen H, Ajmal M, Haq TU, et al. Genome sequencing of  
547 Pakistani families with male infertility identifies deleterious genotypes in SPAG6, CCDC9, TKTL1,  
548 TUBA3C, and M1AP. *Andrology*. 2023 Dec 10;
- 549 34. Garcia FA de O, Andrade ES de, Palmero EI. Insights on variant analysis in silico tools for  
550 pathogenicity prediction. *Frontiers in Genetics* [Internet]. 2022 [cited 2024 Feb 14];13. Available  
551 from: <https://www.frontiersin.org/journals/genetics/articles/10.3389/fgene.2022.1010327>
- 552 35. Li J, Zhao T, Zhang Y, Zhang K, Shi L, Chen Y, et al. Performance evaluation of pathogenicity-  
553 computation methods for missense variants. *Nucleic Acids Res*. 2018 Sep 6;46(15):7793–804.
- 554 36. Anderson D, Lassmann T. An expanded phenotype centric benchmark of variant prioritisation tools.  
555 *Hum Mutat*. 2022 May;43(5):539–46.
- 556 37. Ng PC, Henikoff S. SIFT: predicting amino acid changes that affect protein function. *Nucleic Acids*  
557 *Res*. 2003 Jul 1;31(13):3812–4.
- 558 38. Adzhubei I, Jordan DM, Sunyaev SR. Predicting Functional Effect of Human Missense Mutations  
559 Using PolyPhen-2. *Current Protocols in Human Genetics*. 2013;76(1):7.20.1-7.20.41.
- 560 39. Carter H, Douville C, Stenson PD, Cooper DN, Karchin R. Identifying Mendelian disease genes with  
561 the variant effect scoring tool. *BMC Genomics*. 2013;14 Suppl 3(Suppl 3):S3.
- 562 40. Ioannidis NM, Rothstein JH, Pejaver V, Middha S, McDonnell SK, Baheti S, et al. REVEL: An Ensemble  
563 Method for Predicting the Pathogenicity of Rare Missense Variants. *The American Journal of Human*  
564 *Genetics*. 2016 Oct 6;99(4):877–85.

- 565 41. Pollard KS, Hubisz MJ, Rosenbloom KR, Siepel A. Detection of nonneutral substitution rates on  
566 mammalian phylogenies. *Genome Res.* 2010 Jan;20(1):110–21.
- 567 42. Davydov EV, Goode DL, Sirota M, Cooper GM, Sidow A, Batzoglou S. Identifying a high fraction of  
568 the human genome to be under selective constraint using GERP++. *PLoS Comput Biol.* 2010 Dec  
569 2;6(12):e1001025.
- 570 43. Caron B, Luo Y, Rausell A. NCBoost classifies pathogenic non-coding variants in Mendelian diseases  
571 through supervised learning on purifying selection signals in humans. *Genome Biology.* 2019 Feb  
572 11;20(1):32.
- 573 44. Smedley D, Jacobsen JOB, Jager M, Köhler S, Holtgrewe M, Schubach M, et al. Next-generation  
574 diagnostics and disease-gene discovery with the Exomiser. *Nat Protoc.* 2015 Dec;10(12):2004–15.
- 575 45. Ellingford JM, Ahn JW, Bagnall RD, Baralle D, Barton S, Campbell C, et al. Recommendations for  
576 clinical interpretation of variants found in non-coding regions of the genome. *Genome Medicine.*  
577 2022 Jul 19;14(1):73.
- 578 46. Moyon L, Berthelot C, Louis A, Nguyen NTT, Crollius HR. Classification of non-coding variants with  
579 high pathogenic impact. *PLOS Genetics.* 2022 Apr 29;18(4):e1010191.
- 580 47. Schubach M, Maass T, Nazaretyan L, Röner S, Kircher M. CADD v1.7: using protein language models,  
581 regulatory CNNs and other nucleotide-level scores to improve genome-wide variant predictions.  
582 *Nucleic Acids Research.* 2024 Jan 5;52(D1):D1143–54.
- 583 48. Cunningham F, Allen JE, Allen J, Alvarez-Jarreta J, Amode MR, Armean IM, et al. Ensembl 2022.  
584 *Nucleic Acids Research.* 2022 Jan 7;50(D1):D988–95.
- 585 49. Chen S, Francioli LC, Goodrich JK, Collins RL, Kanai M, Wang Q, et al. A genome-wide mutational  
586 constraint map quantified from variation in 76,156 human genomes [Internet]. *bioRxiv*; 2022 [cited  
587 2023 Aug 30]. p. 2022.03.20.485034. Available from:  
588 <https://www.biorxiv.org/content/10.1101/2022.03.20.485034v2>
- 589 50. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytzky A, et al. The Genome Analysis  
590 Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.*  
591 2010 Sep;20(9):1297–303.
- 592 51. Rentzsch P, Schubach M, Shendure J, Kircher M. CADD-Splice—improving genome-wide variant  
593 effect prediction using deep learning-derived splice scores. *Genome Medicine.* 2021 Feb  
594 22;13(1):31.
- 595 52. Rentzsch P, Witten D, Cooper GM, Shendure J, Kircher M. CADD: predicting the deleteriousness of  
596 variants throughout the human genome. *Nucleic Acids Research.* 2019 Jan 8;47(D1):D886–94.
- 597 53. Blekhman R, Man O, Herrmann L, Boyko AR, Indap A, Kosiol C, et al. Natural selection on genes that  
598 underlie human disease susceptibility. *Curr Biol.* 2008 Jun 24;18(12):883–9.

599 54. Marenne G, Ludwig TE, Bocher O, Herzig AF, Aloui C, Tournier-Lasserre E, et al. RAVAQ: An  
600 integrative pipeline from quality control to region-based rare variant association analysis. *Genetic*  
601 *Epidemiology*. 2022;46(5–6):256–65.

602

603