



Article

Aggregated Genomic Data as Cohort-Specific Allelic Frequencies can Boost Variants and Genes Prioritization in Non-Solved Cases of Inherited Retinal Dystrophies

Ionut-Florin Iancu ^{1,2}, Irene Perea-Romero ^{1,2}, Gonzalo Núñez-Moreno ^{1,2,3} , Lorena de la Fuente ^{1,3}, Raquel Romero ^{1,2}, Almudena Ávila-Fernandez ^{1,2}, María José Trujillo-Tiebas ^{1,2}, Rosa Riveiro-Álvarez ^{1,2}, Berta Almoguera ^{1,2} , Inmaculada Martín-Mérida ^{1,2}, Marta Del Pozo-Valero ^{1,2}, Alejandra Damián-Verde ¹ , Marta Cortón ^{1,2}, Carmen Ayuso ^{1,2,*} and Pablo Minguez ^{1,2,3,*}

- ¹ Department of Genetics, Health Research Institute–Fundación Jiménez Díaz University Hospital, Universidad Autónoma de Madrid (IIS-FJD, UAM), 28049 Madrid, Spain; ionut.iancu@quironsalud.es (I.-F.I.); irene.perea@quironsalud.es (I.P.-R.); gonzalo.nunezm@quironsalud.es (G.N.-M.); ldelafuente.lorena@gmail.com (L.d.l.F.); raquel.romerof@quironsalud.es (R.R.); aavila@quironsalud.es (A.Á.-F.); mjtrujillo@fjd.es (M.J.T.-T.); rriveiro@fjd.es (R.R.-Á.); balmoguera@quironsalud.es (B.A.); inmaculada.martinm@quironsalud.es (I.M.-M.); martapova21@gmail.com (M.D.P.-V.); alejandra.damian@quironsalud.es (A.D.-V.); mcorton@quironsalud.es (M.C.)
- ² Center for Biomedical Network Research on Rare Diseases (CIBERER), Instituto de Salud Carlos III (ISCIII), 28040 Madrid, Spain
- ³ Bioinformatics Unit, Health Research Institute–Fundación Jiménez Díaz University Hospital, Universidad Autónoma de Madrid (IIS-FJD, UAM), Cantoblanco, 28049 Madrid, Spain
- * Correspondence: cayuso@fjd.es (C.A.); pablo.minguez@quironsalud.es (P.M.)



Citation: Iancu, I.-F.; Perea-Romero, I.; Núñez-Moreno, G.; de la Fuente, L.; Romero, R.; Ávila-Fernandez, A.; Trujillo-Tiebas, M.J.; Riveiro-Álvarez, R.; Almoguera, B.; Martín-Mérida, I.; et al. Aggregated Genomic Data as Cohort-Specific Allelic Frequencies can Boost Variants and Genes Prioritization in Non-Solved Cases of Inherited Retinal Dystrophies. *Int. J. Mol. Sci.* **2022**, *23*, 8431. <https://doi.org/10.3390/ijms23158431>

Academic Editor: Stephanie C. Joachim

Received: 23 June 2022

Accepted: 26 July 2022

Published: 29 July 2022

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2022 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

Abstract: The introduction of NGS in genetic diagnosis has increased the repertoire of variants and genes involved and the amount of genomic information produced. We built an allelic-frequency (AF) database for a heterogeneous cohort of genetic diseases to explore the aggregated genomic information and boost diagnosis in inherited retinal dystrophies (IRD). We retrospectively selected 5683 index-cases with clinical exome sequencing tests available, 1766 with IRD and the rest with diverse genetic diseases. We calculated a subcohort's IRD-specific AF and compared it with suitable pseudocontrols. For non-solved IRD cases, we prioritized variants with a significant increment of frequencies, with eight variants that may help to explain the phenotype, and 10/11 of uncertain significance that were reclassified as probably pathogenic according to ACMG. Moreover, we developed a method to highlight genes with more frequent pathogenic variants in IRD cases than in pseudocontrols weighted by the increment of benign variants in the same comparison. We identified 18 genes for further studies that provided new insights in five cases. This resource can also help one to calculate the carrier frequency in IRD genes. A cohort-specific AF database assists with variants and genes prioritization and operates as an engine that provides a new hypothesis in non-solved cases, augmenting the diagnosis rate.

Keywords: genetic rare diseases; inherited retinal dystrophies; variant prioritization; gene prioritization; variants of uncertain significance; carrier frequency

1. Introduction

Rare diseases are chronically debilitating or life-threatening and have a prevalence in Europe of less than 1 in every 2000 people [1]. Inherited retinal dystrophies (IRD) are a group of rare diseases with a degenerative and progressive course and are caused by the primary affection of photoreceptors and retinal pigmentary epithelial [2]. All together, they affect 1 in every 3000–4000 people in the western world [3]. They are clinically heterogeneous, covering several syndromes (e.g., Usher, Bardet-Biedl -BBS-,

or Joubert) [4–6], as well as non-syndromic forms such as retinitis pigmentosa [2] and macular dystrophies [7]. They have overlapped phenotypes and display any form of inherited patterns.

During the last two decades, next-generation sequencing (NGS) techniques have transformed research on genetic rare diseases, causing a substantial increase in the volume of available genomic data and knowledge generated [8,9]. Although several sequencing tests are available, a widespread approach in genetic diagnosis is to sequence the coding region of known clinically relevant genes (~4500), the so-called clinical exome (CE). A CE test detects several thousand variants, which need to be filtered and prioritized in order to highlight those responsible for the phenotype [10]. Regarding the task of filtering, in low-prevalence diseases, apart from quality filters, a small population frequency is one of the first requirements to purge non-causal variants [11]. Several global genomic initiatives [12,13] provide allele frequencies in large populations, although frequencies from local cohorts provide a better estimation of real variant prevalence [14,15] and have been proven to identify rare pathogenic variants [15,16].

In the absence of genomic information of a priori healthy people, Mendelian diseases can provide a good estimate of allele frequencies in the general population as pseudocontrols (PC) for other non-related diseases [17,18]. In the same terms, PC can also be applied to the calculation of carrier frequencies (CF) [16,19] of causal variants of non-related diseases in genes with a recessive inheritance pattern, as well as to the analysis of trios [20,21]. In more complex scenarios, where modifying and risk/protective variants may tune the effect of causal variants, the mutational landscape of a disease may help identify: (i) the genetic pleiotropy, together with causal variants in recessive forms [22]; (ii) digenic inheritance [23]; or (iii) disease-associated triallelic sites as in BBS [24].

With all these premises, we hypothesize that a database of variant allele frequencies calculated over a heterogeneous cohort of genetic diseases enriched in IRD cases can help to improve the detection of previously unnoticed, underrated, or unknown causal variants and gene-disease associations. Additionally, it can help one to uncover disease cases with overlapping phenotypes, as well as to describe the carrier frequencies of recessive variants. Thus, we built a database with the genomic data of a large cohort with various genetic diseases and developed methods to compare IRD-specific and PC frequencies. This tool is used as a global reanalysis platform to study frequent variants and over-mutated genes in IRD non-solved cases.

2. Results

2.1. A Multi-Disease Cohort Database of Variant Frequencies to Study the Aggregated Signal in Ird Genomic Landscape

We compiled a heterogeneous cohort of 5683 patients with genetic diseases referred to the Genetics Department of the UH-FJD and with a clinical exome sequencing test available (see Section 4). The cases were distributed into three groups of diseases as inherited retinal dystrophies (IRDs), with 1766 cases; other eye-related diseases (OERDs), with 386 cases; and non-eye-related diseases (NRD), with 3531 cases (Figure 1A). Additionally, IRD cases were classified according to their diagnostic status as solved ($N = 955$ cases, 54%) and non-solved ($N = 811$ cases, 46%). Within the cohort of non-solved IRD patients, we had 447 cases with no candidate variants (25% of the total), that is, excluding cases with a pathogenic/likely pathogenic variant reported in recessive cases (named partially solved cases) and cases with 1–2 VUS reported in dominant/recessive cases, respectively (named VUS cases) (Figure 1B).

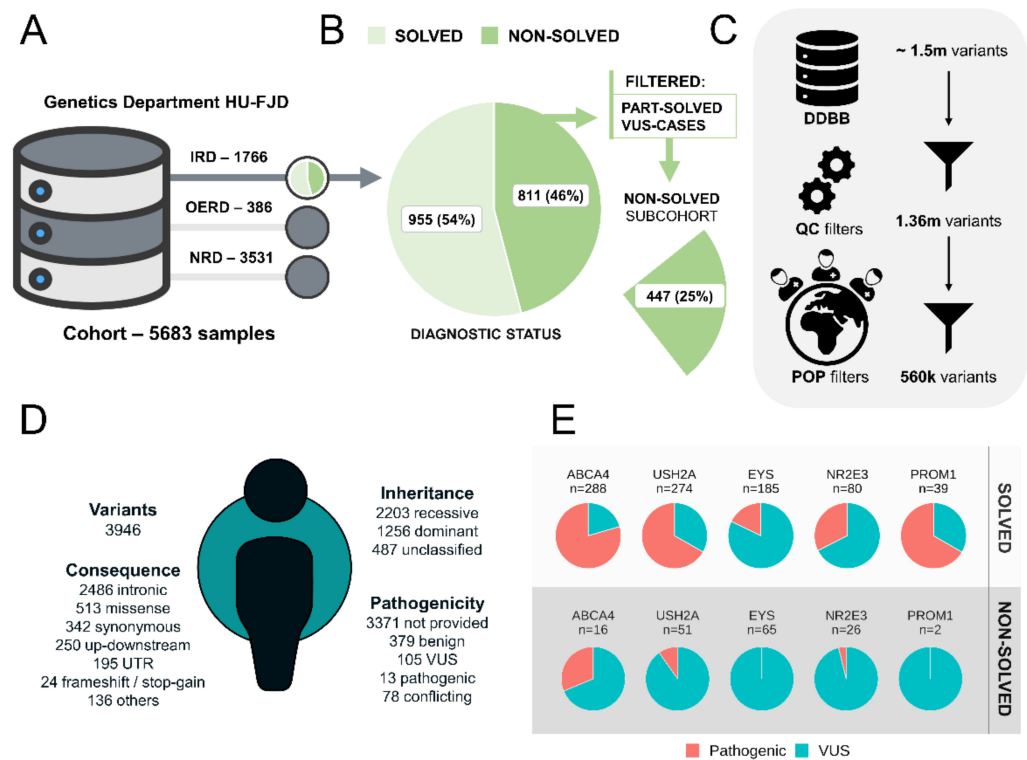


Figure 1. A heterogeneous cohort of rare diseases with the diagnostic status and variant composition in IRD cases. (A) The cohort of patients with suspected rare genetic diseases at the Genetics Department of HU-FJD was divided into three subcohorts. An IRD subcohort of 1766 samples, other-eye-related diseases (OERD) of 386 cases, and a pseudocontrol subcohort of non-eye-related diseases (NRD) with 3531 samples. (B) IRD diagnostic status of the samples included in IRD subcohort was solved and non-solved. The subcohort of non-solved IRD cases with no candidate variants represent the 25%. (C) Flow chart of different filters applied to variants according to quality control (QC) and population (POP) filters. (D) Summary of the variants included in the database in an average IRD case; values represent the average of all IRD samples. (E) Proportion of pathogenic and VUS variants detected in IRD cases in the genes with more pathogenic variants in IRD solved cases; top 5 genes are shown.

CEs were reanalyzed and variants detected and extracted for each case. We performed a quality filter and removed variants with a potential sub-population bias between IRD and other cases (Figure 1C and Section 4). Thus, for the 5683 samples together, around 560 K unique variants in 5046 genes were left for further analyses.

Using this set of variants, an average IRD individual has approximately 4 K non-polymorphic variants (Allele frequency, AF < 0.1) within the specified sequenced region (Figure 1D), being mostly intronic ($N = 2486$, 63%), missense ($N = 513$, 13%), and synonymous ($N = 342$, 9%). According to the inheritance pattern observed for the genes based on their associated diseases, 2203 variants (56%) have a recessive pattern; 1256 (32%) have a dominant pattern; and 487 (12%) either do not have a clear pattern, are undetermined, or are X-linked. Regarding their pathogenicity, 379 (10%) are benign or likely benign, 105 (3%) VUS, and 13 (0,3%) pathogenic or likely pathogenic (according to ClinVar). However, there are still a large percentage of variants with missing or conflicting annotation ($N = 3449$, 87%).

From a cohort perspective, the number of pathogenic (including likely pathogenic) variants and VUS are unequally distributed over IRD-associated genes when comparing solved-IRD and non-solved-IRD cases (Figure 1D). The top 5 genes with more pathogenic variants (including likely pathogenic) in solved-cases have a higher ratio of pathogenic variants/VUS for IRD solved cases than for IRD non-solved cases.

2.2. IRD-Specific Highly Frequent Variants

In the diagnosis of rare diseases using NGS data, variants must be prioritized in order to facilitate the detection of the causal mutations within the large amount of variation found in a single experiment. Here, we aimed to test if the comparison between the frequency of variants in cases and controls was able to extract deleterious variants. Thus, we focused on the IRD subcohort and treated solved and non-solved cases separately. Partially solved and VUS cases were excluded from the non-solved subcohort in order to work with a subcohort with no candidate variants. As controls for both subcohorts, we used the set of patients with non-related diseases, from now on called pseudocontrols (PC) (see Section 4). Allele frequencies (AF), allele numbers (AN), and allele counts (AC) were calculated for solved-IRD cases, non-solved-IRD cases (free of candidate variants), and PC cases.

Next, for both solved and non-solved IRD cases independently, we compared each variant's AF (solved-AF or non-solved-AF) with its AF in the PC subcohort (PC-AF), see Section 4. We defined the "most frequent variants" in a IRD subcohort (IRD-MFVs) as those within the top 10% with the highest \log_2 of the fold change, $\log_2(\text{FC})$ values in the comparison performed. The cut offs for the $\log_2(\text{FC})$ were 3.12 and 3.92 in solved and non-solved IRD cases, respectively (Figure 2A,B). The distribution of $\log_2(\text{FC})$ values in solved and non-solved IRD cases is shown in the Supplementary Figure S1. MFVs are prioritized for a posterior reevaluation of cases. Non-prioritized variants are defined as those with a higher frequency in IRD ($\text{FC} > 0$) but below the significant threshold. Classifying IRD-MFVs according to their clinical relevance and removing those not informative (see Section 4), we found IRD-MFVs enriched in deleterious variants in both solved and non-solved cases compared to non-prioritized variants (Fisher's exact test, p -values = 4.77×10^{-56} and 1.69×10^{-32} , respectively; Figure 2C,D and Supplementary Tables S1 and S2).

Focusing on the type of genes where the IRD-MFVs are located, we divided IRD-MFVs as present in IRD-associated genes, OERD-related genes, and NRD-associated genes. In IRD-solved cases, regarding IRD-MFVs in IRD-genes only, the ratio deleterious/benign is 78% (177/256), which is significantly higher than the same proportion in non-prioritized variants (14%, 593/1926, Figure 2E, Supplementary Table S1). A different trend was observed in non-solved IRD cases, where we found no significant differences in the percentages of deleterious/benign variants in IRD-genes between IRD-MFVs and non-prioritized variants (Figure 2F, Supplementary Table S2). In IRD-MFVs in OERD-genes, we also found more deleterious variants in our prioritized set in both solved and non-solved cases (Fisher's exact test, p -value = 1.71×10^{-6} and p -value = 3.25×10^{-6} , respectively; Figure 2E,F, Supplementary Tables S1 and S2). Finally, we also observed an enrichment of deleterious variants in the IRD-MFVs located in NRD-genes in solved and non-solved IRD cases (p -value = 2.75×10^{-8} and p -value = 1.90×10^{-28} , respectively; Figure 2E,F, Supplementary Tables S1 and S2).

Furthermore, solved and non-solved cases were divided into disease sub-categories as syndromic, non-syndromic, and macular dystrophy forms (Supplementary Figure S2), and the analysis was repeated for each sub-group. Thus, as for all solved cases considered as a whole, syndromic, macular dystrophy, and non-syndromic forms behave very similarly, with more deleterious variants in the prioritized sets for IRD genes (Supplementary Figure S3A–C, Supplementary Tables S3 and S4). For non-solved cases, we observed an increase in deleterious variants in the prioritized sets for macular dystrophy and non-syndromic forms (Supplementary Figure S3D–F and Supplementary Tables S3 and S4).

Next, we performed a reevaluation of non-solved IRD cases carrying an IRD-MFV and found eight variants that provide additional insights in 8 cases (Table 1). In this reevaluation, the diagnosis status of the cases is classified as (1) "solved", when the case has a conclusive diagnosis thanks to the variant(s) identified; or (2) "non-solved", in any other case. Non-solved cases were annotated according to findings that may help when making a future diagnosis as (2.1) "partially solved", if a heterozygous pathogenic or likely pathogenic variant is found within a recessive gene that fits the observed phenotype; and (2.2) "with evidence", when the variants identified in the analysis are: (2.2.1)

pathogenic/likely pathogenic variants in genes not yet associated with the disease but with some evidence published in the literature or with overlapping phenotypes, or (2.2.2) the VUS in a gene associated with the phenotype when one VUS in dominant or 1–2 VUS in recessive genes were found.

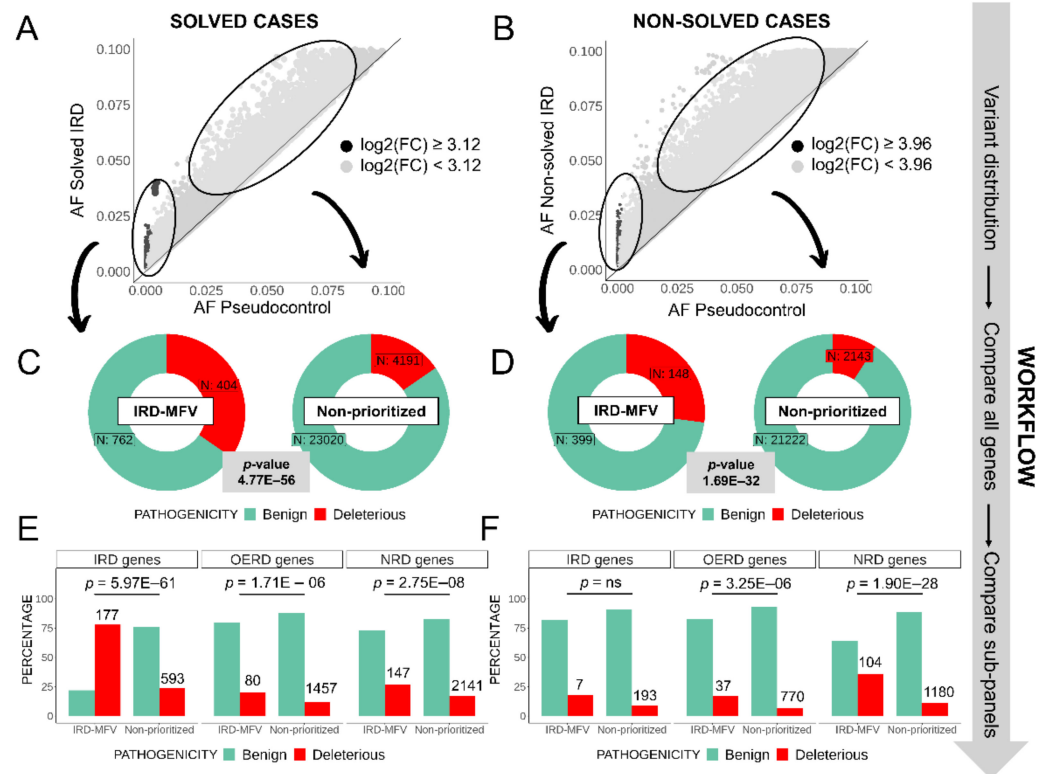


Figure 2. Comparison of the percentage of deleterious variants in prioritized and non-prioritized variants in solved and non-solved IRD patients. Variant AFs are compared in inherited retinal dystrophies (IRD) solved (A) and non-solved (B) subcohorts against pseudocontrols (PC) using fold changes (FC). IRD more frequent variants (IRD-MFVs), highlighted in dark, were defined using FC thresholds at 0.90 percentiles of all FC in each comparison (A,B). Proportion of deleterious and benign variants in both solved (C) and non-solved IRD cases (D) and the *p*-values representing the enrichment of deleterious variants in IRD-MFVs. Enrichment analyses are also performed dividing the IRD-MFVs according to the genes in which they are located; they are grouped into IRD genes, other eye-related diseases (OERD genes), and other non-eye-related diseases (NRD genes) (E,F). Total number of deleterious variants in each group is noted at the top of the red bars. Non-significant *p*-values are marked as “ns”.

Another direct application of the IRD relative variant frequencies is the reevaluation of the clinical significance of VUS. We extracted the FC of IRD-AF compared to PC-AF for a set of manually curated VUS whose reclassification could contribute to a conclusive diagnosis of an IRD case in our cohort and that is present in the final dataset (*N* = 63). Of them, six VUS are IRD-MFVs ($\log_2(\text{FC}) \geq 2.48$, see Section 4) and 11 VUS are more frequent in IRD cases than in PCs with a $\log_2(\text{FC}) \geq 1.5$. ACMG classification was performed for the 11 VUS adding the ACMG criteria PS4 (“The prevalence of the variant in affected individuals is significantly increased compared with the prevalence in controls”) as true. Of them, 10 (91%) were reclassified as likely pathogenic/pathogenic (Table 2). This analysis was able to solve two cases: one with the variant in the dominant gene *COL11A1*, and one with the variant in *OFD1* that has X-linked inheritance (Table 2).

Table 1. Identification of most frequent variants in IRD (IRD-MFV) allowed for the reevaluation of non-solved IRD cases. A total of 8 cases (one per row) gained knowledge in their reassessment. The gene panel column refers to the type of diseases that the gene has been associated with, inherited retinal dystrophies (IRD), other eye-related diseases (OERD), and other non-eye-related diseases (NRD). In the phenotype, abbreviations are RCD: rod-cone dystrophy, RP: retinitis pigmentosa, and MD: macular dystrophy. In variant column: HGVS_c, HGVS_p, and ACMG classification are included. In “Other Variants”, for those biallelic cases, the second variant found after reviewing the case is included. Genotype column has information about the genotype of the variant(s) in the sample. Diagnostic status column has the new diagnostic status of the case after the reanalysis, as (1) “solved”, when the case has a conclusive diagnosis thanks to the variant(s) identified; or (2) “non-solved”, in any other case. Non-solved cases were annotated according to findings that may help in a future diagnosis as (2.1) “partially solved”, if a heterozygous pathogenic or likely pathogenic variant is found within a recessive gene that fits the observed phenotype; and (2.2) “with evidence”, when the variants identified in the analysis are: (2.2.1) pathogenic/likely pathogenic variants in genes not yet associated with the disease but with some evidence published in the literature or with overlapping phenotypes, or (2.2.2) VUS in a gene associated with the phenotype when one VUS in dominant or 1–2 VUS in recessive genes were found.

Gene	Gene Panel	Phenotype	Prioritized Variant (HGVS _c , HGVS _p , and ACMG Classification)	Other variants	Genotype	Diagnostic Status
<i>CDH23</i>	IRD	RCD	NM_022124.6: c.6050–9G > A pathogenic	-	Monoallelic 0/1	Partially solved
<i>MYO7A</i>	IRD	Usher syndrome	NM_000260.4: c.1996C > T (p. Arg666Ter); pathogenic	NM_000260.4: c.3764del (p. Lys1255ArgfsTer8) (pathogenic)	Biallelic 0/1, 0/1	Solved
<i>IFT88</i>	ONRD	RP	NM_001318491.2: c.538G > T (p. Val180Phe); pathogenic	-	Monoallelic 0/1	With evidence
<i>KIAA2022</i>	NRD	MD	NM_001008537.3: c.4385del (p. Cys1462LeufsTer24); likely pathogenic	-	Monoallelic 0/1	With evidence
<i>TTPA</i>	IRD	RP	NM_000370.3: c.227_235del (p. Trp76Ter); likely pathogenic	-	Monoallelic 0/1	With evidence
<i>TTPA</i>	IRD	MD	NM_000370.3: c.227_235del (p. Trp76Ter); likely pathogenic	-	Monoallelic 0/1	With evidence
<i>CDHR1</i>	IRD	RP	NM_033100.4: c.2410_2485del (p. Thr804ProfsTer12); likely pathogenic	-	Monoallelic 0/1	Partially solved
<i>CDHR1</i>	IRD	RP	NM_033100.4: c.2410_2485del (p. Thr804ProfsTer12); likely pathogenic	NM_033100.4: c.476C > A (p. Ala159Glu) (VUS)	Biallelic 0/1, 0/1	Partially solved

2.3. Prioritization of Candidate Genes Based on Weighted Cohort-Specific Frequency of Pathogenic and Benign Variants in Non-Solved Ird Cases

In order to detect genes with an accumulated high pathogenicity in solved and non-solved IRD cases as good candidates to be involved in IRD phenotypes, for each IRD subcohort, we extracted deleterious and benign variants and calculated the FC for their AF compared to the AF in PCs. For every subcohort and gene, the distributions of log₂(FCs) in deleterious and benign variants were compared using a Wilcoxon rank sum test. Genes with a significant higher frequency in IRD cases of deleterious variants compared to benign variants (*p*-value < 0.05) were selected. This revealed a number of genes with an accumulated pathogenicity in solved and non-solved IRD cases that we divided, as before, into three groups: IRD-genes, OERD-genes, and NRD-genes. Actionable genes defined by ACMG were removed from the analysis (Supplementary Table S5). Thus, in

IRD solved cases, we found 56 genes enriched in deleterious mutations. Of them, 22 (39%) were IRD genes, and the top 5 genes with more deleterious variants were *ABCA4*, *USH2A*, *MYO7A*, *EYS*, and *ADGRV1*, (Supplementary Table S6, and Figure 3A). In addition, 24 (43%) were OERD genes; the top 5 are highlighted here with more deleterious variants *NEB*, *PAH*, *DNAH11*, *DNAH5*, and *ATM* (Supplementary Table S6 and Figure 3A), and 10 (18%) were NRD genes, including the top 5 genes *OBSCN*, *DYSF*, *SPTBN5*, *OTOF*, and *SPTB*. Regarding non-solved cases, we found 18 genes with an accumulated pathogenicity, with IRD-genes less represented (22%) and OERD-genes more present (55%) than in solved cases. Finally, 22% of the prioritized genes were NRD-genes (Figure 3B, and Table 3). Interestingly, there was a high overlap between IRD-genes and OERD-genes prioritized in solved and non-solved cases (75% and 70% of the smallest group (non-solved cases), respectively). However, the overlap in NRD-genes was smaller, with only one gene (25% of genes in non-solved cases) found in both IRD subcohorts (Supplementary Figure S4).

Table 2. Reclassification of VUS, with information added about the difference in their frequency in cases and controls. The ACMG classification of VUS with log₂ fold change (FC) higher than 1.5 is reevaluated. Gene code, variants in nucleotide code, and protein code are provided. The inheritance mode is annotated as autosomal recessive (AR), autosomal dominant (AD), and x-linked (XL). In the “Varsome” column, we add the automatic ACMG classification of the variants provided by the Varsome database at the time of writing. Previous ACMG criteria fulfilled by the cases are annotated in the “Criteria” column. Column “PS4” provides the new ACMG classification reached after including PS4 criteria. In column “Status”, reclassified variants are marked.

Gene	HGVSc	HGVSp	Inheritance	log ₂ (FC)	Varsome	Criteria	PS4	Status
<i>CACNA1F</i>	NM_001256789.3: c.4009-3C>G	-	XL	2.7	3	PM2, BP4	4	Reclassified
<i>CDH23</i>	NM_022124.5: c.4231G>A	p. Glu1411Lys	AR	2.0	3	PM2, PP3	4	Reclassified
<i>CDHR1</i>	NM_001171971.3: c.1589C>G	p. Thr530Ser	AR	2.6	3	PM1, PM2, and BP1	4	Reclassified
<i>COL11A1</i>	NM_080629.2: c.4838C>A	p. Thr1613Asn	AD	1.6	3	PM2, PP2	4	Reclassified
<i>GDF6</i>	NM_001001557.4: c.125G>T	p. Gly42Val	AR/AD	2.0	3	PM1, PM2, PP5, and BP6	4	Reclassified
<i>IMPG2</i>	NM_016247.4: c.1460A>T	p. His487Leu	AR/AD	2.7	3	PM2, BP4	4	Reclassified
<i>MERTK</i>	NM_006343.3: c.2435A>G	p. Tyr812Cys	AR	2.0	3	PM1, PM2, PP3, and BP6	4	Reclassified
<i>NYX</i>	NM_022567.2: c.505A>G	p. Asn169Asp	XL	2.9	3	PM2	3	Not reclassified
<i>OFD1</i>	NM_003611.2: c.87T>G	p. Asp29Glu	XL	2.7	3	PM2, PP3, and BP1	4	Reclassified
<i>RP1</i>	NM_006269.2: c.2497T>C	p. Phe833Leu	AR/AD	1.7	3	PM2, BP4	4	Reclassified
<i>WFS1</i>	NM_001145853.1: c.1597C>T	p. Pro533Ser	AR/AD	3.7	3; 4	PM2, PP3	4	Reclassified

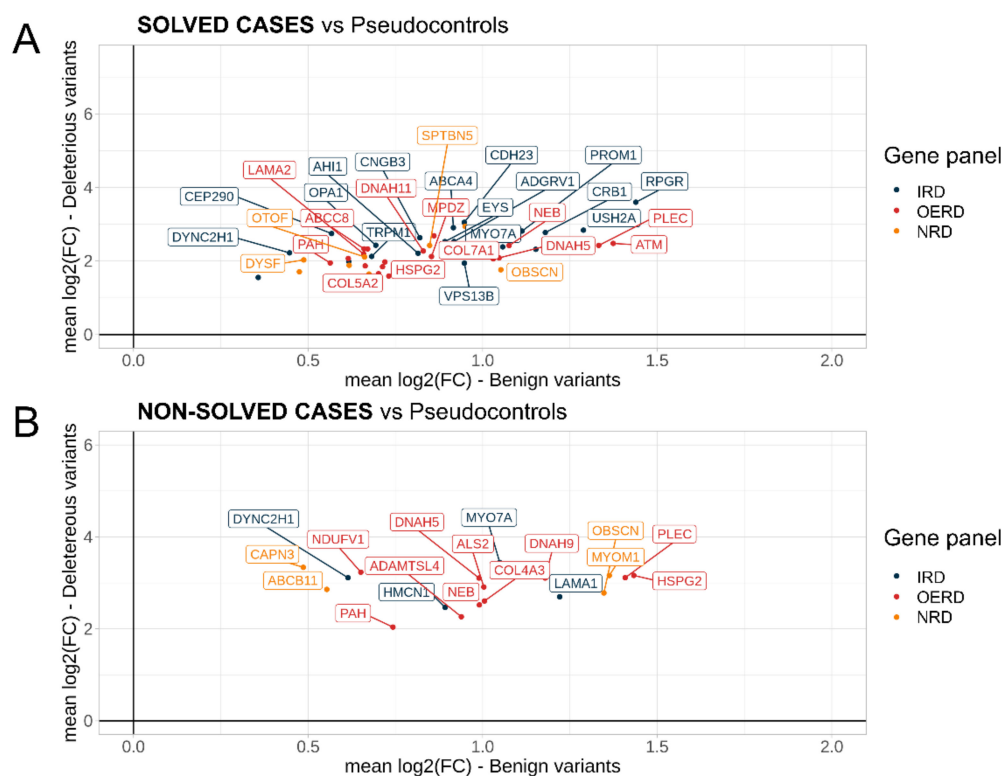


Figure 3. Genes with a higher accumulated pathogenicity in solved and non-solved IRD cases compared to pseudocontrols. Mean fold changes (FCs) in log2 scale for deleterious (Y-axis) and benign variants (X-axis) are shown for each gene. For IRD solved cases, significant genes with at least eight deleterious variants are shown in the plot (A). For IRD non-solved cases, all significant genes are shown (B).

Table 3. Genes with more frequent deleterious variants in IRD non-solved cases than in pseudocontrols. The genes are sorted by the gene panel where they are present (IRD, OERD, and NRD) and by number of deleterious variants. The gene panel column refers to the type of diseases that the gene has been associated with, inherited retinal dystrophies (IRD), other eye-related diseases (OERD), and other non-eye-related diseases (NRD). Number of deleterious and benign variants in IRD non-solved cases is shown. FDR stands for false discovery rate of the Wilcoxon rank sum test performed. Inheritance mode associated with each gene is annotated in the “Inheritance” column as recessive, dominant, or R/D (either recessive or dominant).

Gene	Deleterious	Benign	FDR	Gene Panel	Inheritance
<i>MYO7A</i>	10	42	1.08×10^{-3}	IRD	Recessive
<i>DYNC2H1</i>	9	43	9.81×10^{-4}	IRD	Recessive
<i>LAMA1</i>	7	22	3.56×10^{-2}	IRD	Recessive
<i>HMCN1</i>	6	43	2.75×10^{-2}	IRD	Dominant
<i>NEB</i>	12	59	4.55×10^{-3}	OERD	Recessive
<i>PAH</i>	11	12	1.32×10^{-2}	OERD	Recessive
<i>ALS2</i>	8	11	1.87×10^{-2}	OERD	Recessive
<i>DNAH9</i>	7	13	1.32×10^{-2}	OERD	Recessive
<i>HSPG2</i>	6	44	1.87×10^{-2}	OERD	R/D
<i>DNAH5</i>	6	43	1.87×10^{-2}	OERD	Recessive
<i>PLEC</i>	6	80	1.87×10^{-2}	OERD	Dominant

Table 3. Cont.

Gene	Deleterious	Benign	FDR	Gene Panel	Inheritance
<i>ADAMTSL4</i>	5	16	2.70×10^{-2}	OERD	Recessive
<i>NDUFV1</i>	5	8	1.32×10^{-2}	OERD	Recessive
<i>COL4A3</i>	5	15	1.87×10^{-2}	OERD	R/D
<i>OBSCN</i>	9	91	1.32×10^{-2}	NRD	Recessive
<i>CAPN3</i>	7	7	1.28×10^{-2}	NRD	Recessive
<i>MYOM1</i>	6	22	1.32×10^{-2}	NRD	Recessive
<i>ABCB11</i>	6	7	1.87×10^{-2}	NRD	Recessive

In a reevaluation of the non-solved cases with pathogenic, likely pathogenic, or VUS variants in the prioritized genes, we found five cases carrying mutations, possibly associated with the phenotype in two IRD associated genes, and one OERD gene (Table 4). Variants found in the gene *MYO7A* explained this case phenotype and helped to fully characterize it. Regarding OERD genes, the two pathogenic variants found in the gene *ADAMTSL4* helped to solve lens luxation phenotype of this case.

Table 4. New diagnostic status produced by the reassessment of non-solved IRD cases with deleterious variants in the genes prioritized. Details about five IRD non-solved cases with variants adding knowledge to the phenotype. The gene panel column refers to the type of diseases that the gene has been associated with, inherited retinal dystrophies (IRD), and other eye-related diseases (OERD). In the phenotype, abbreviations are MD: macular dystrophy. In variant column, HGVS_c, HGVS_p, and ACMG classification are included. Genotype column has information about the genotype of the variant(s) in the sample. Diagnostic status column has the new diagnostic status of the case after the reanalysis, as (1) “solved”, when the case has a conclusive diagnosis thanks to the variant(s) identified; or (2) “non-solved”, in any other case. Non-solved cases were annotated according to findings that may help in a future diagnosis as (2.1) “partially solved”, if a heterozygous pathogenic or likely pathogenic variant is found within a recessive gene that fits the observed phenotype; and (2.2) “with evidence”, when the variants identified in the analysis are: (2.2.1) pathogenic/likely pathogenic variants in genes not yet associated with the disease but with some evidence published in the literature or with overlapping phenotypes, or (2.2.2) VUS in a gene associated with the phenotype when one VUS in dominant or 1–2 VUS in recessive genes were found.

Gene	Gene Panel	Phenotype	Variant (HGVS _c , HGVS _p , ACMG Classification)	Genotype	Diagnostic Status
<i>YNC2H1</i>	IRD	MD	NM_001080463.2: c.3793C>T (p. Arg1265Cys), VUS; NM_001080463.2: c.1468C>T (p. Arg490Cys), VUS	Biallelic 0/1, 0/1	With evidence
<i>DYNC2H1</i>	IRD	MD	NM_001080463.2: c.988C>T p. Arg330Cys, pathogenic	Monoallelic 0/1	Partially solved
<i>DYNC2H1</i>	IRD	MD	NM_001080463.2: c.7966C>T p. Arg2656Cys, likely pathogenic	Monoallelic 0/1	Partially solved
<i>MYO7A</i>	IRD	Usher syndrome	NM_000260.4: c.1996C>T (p. Arg666Ter), pathogenic; NM_000260.4: c.3764del (p. Lys1255ArgfsTer8), pathogenic	Biallelic 0/1,0/1	Solved
<i>ADAMTSL4</i>	OERD	Lens luxation	NM_001288607.2: c.2594G>A (p. Ser865Asn), pathogenic; NM_001288607.2: c.767_786del (p. Gln256ProfsTer3), pathogenic	Biallelic 0/1,0/1	Solved

2.4. Carrier's Frequency of Recessive Diseases from a Multi-Disease Cohort

We calculated the carrier frequency (CF) for 69 genes involved in recessive non-syndromic IRDs using the frequency of pathogenic heterozygous variants located in IRD genes in the pseudocontrol subcohort. We found three genes with a CF $\geq 1\%$, which represents a total of $\sim 4\%$ of the total analyzed (Supplementary Table S12). We highlight *ABCA4* and *USH2A* with a CF of $\sim 7\%$ and $\sim 3\%$, respectively (Figure 4), with the first being responsible of Stargardt disease and the later causing Usher syndrome. In the case of *ABCA4*, if hypomorphic variants are excluded as in a previous work performed by Hanany and collaborators [19], CF is reduced to $\sim 4\%$. These genes are also the two most frequent in our IRD subcohort, found causally in 21% and 15% of the cases, respectively (data not shown). The most frequent variants for *ABCA4* in the IRD subcohort and PC subcohort were NM_000350.3: c.3386G>T and NM_000350.3: c.3113C>T, respectively, while *USH2A* had variant NM_007123.5: c.2276G>T as the most frequent in the two subcohorts. The CF calculated using deleterious heterozygous/hemizygous variants on dominant/X-linked genes is available in Supplementary Table S12.

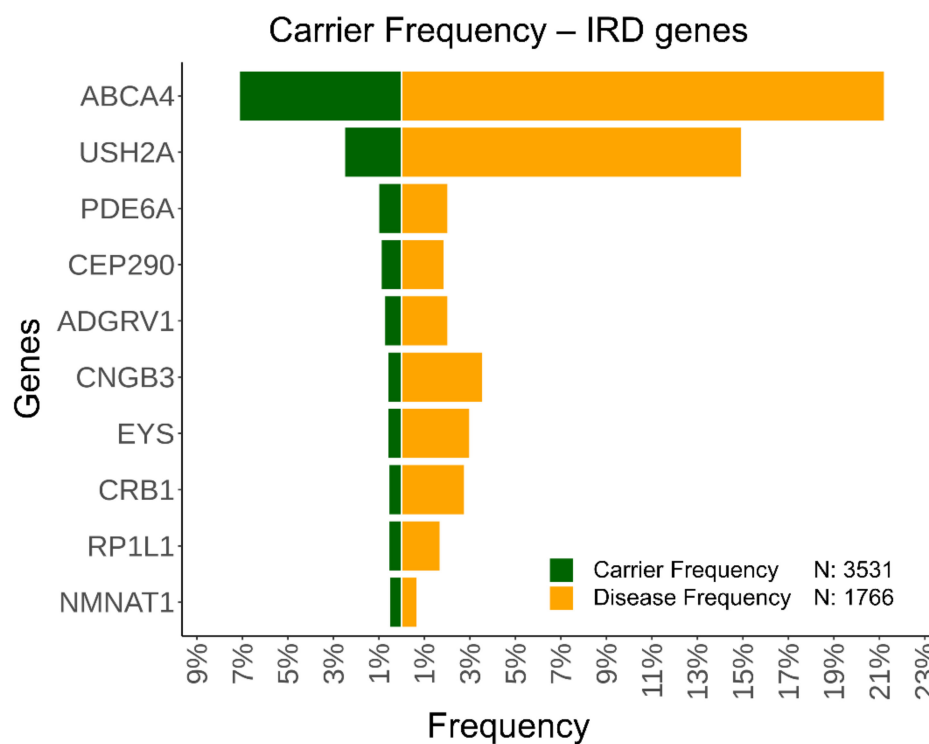


Figure 4. Carrier frequency of pathogenic variants in recessive IRD genes in pseudocontrols. The carrier frequency (CF) was calculated for all IRD genes with a recessive inheritance pattern, and at least one solved case in our cohort. Green represents the CF of the genes and orange the frequency in our IRD subcohort. The pseudocontrol subcohort was composed of 3531 cases, and the IRD subcohort had 1766 cases. The top 10 genes with higher CF are ordered decreasingly.

3. Discussion

The recent application of NGS techniques has considerably increased our competence to study and diagnose rare diseases. Regarding IRDs, although new associated genes are still being discovered (see RetNet, <https://sph.uth.edu/retnet>, accessed on 21 June 2022), the diagnostic ratios need to be boosted since up to $\sim 46\%$ remain unsolved (data from our own IRD cohort at the time of writing). At the same time, at the diagnostic setting, the genomic information from patients is being accumulated as databases or annotation systems, which can also be privative if only commercial solutions are used. Initiatives such as gnomAD [17] or the collaborative Spanish variability server [18] are acting as a crowdsourcing to recover this data from the community and offer it back as aggregated

allele frequencies. Allelic frequency calculations have been proven to be useful in studying the prevalence of rare diseases [25]; integrated in the diagnostic analysis routines, as an additional annotation source to detect technical biases [26]; or used to interpret and classify variants [27]. These resources have also been used to study the CF of deleterious variants in recessive genes in a particular population [16]. In this study, we propose to develop a framework of methods that are able to prioritize variants and genes in order to explore new associations with a specific disease. Thus, the reuse of the genomic data may provide new discovery capabilities to a heterogeneous cohort of genetic diseases, an extra value to the resources invested previously, as well as allowing the patients to contribute to their own or others' future diagnosis. The main hypothesis behind this work is that, aside from a few causal mutations, the genome of patients with Mendelian diseases behaves similarly to those of the healthy population, so patients with non-related diseases can act as pseudocontrols between each other. The fact of focusing on a heterogeneous but single-center cohort has two main advantages against using controls from larger genomic population databases [12,13,17]: first they can be technically more similar to the sequencing produced on the patients of interest; second, the geographical origin bias can be better controlled [15]; last, the phenotyping of patients can be fine-tuned by experts in order to create subcohorts of interest. A concept introduced here is the significantly different frequent variants, called MFVs, which are variants more present in cases than in pseudocontrols, with statistical support. As a proof of concept, we found that considering the whole IRD cohort, the IRD-MFVs are enriched in deleterious mutations, suggesting that the prioritization is effective there. Indeed, variants selected on solved IRD cases in IRD genes are mostly pathogenic (>78%) compared to benign variants, due to the detection of prevalent causal variants. In contrast, non-solved IRD cases display no differences in the proportion of deleterious/benign variants in prioritized and non-prioritized variants in IRD-related genes, indicating that there are not many described pathogenic variants left out during diagnosis. In spite of this, our approach was able to solve or partially solve five cases with prioritized variants in IRD genes. In this paper, we also highlight the significant accumulation of deleterious variants in those prioritized located in OERD and NRD genes, in both IRD solved and non-solved cases. The rationale behind it might be different, though. While frequent pathogenic mutations in OERD and NRD genes in solved cases may suggest a more complex genotype scenario for Mendelian diseases [28,29], in non-solved cases we have to add the possibility of needing a disease re-evaluation, and the causal mutation being present in a not yet associated IRD gene. The exploration of syndromic and non-syndromic forms provides the same signal but with lower p -values for syndromic cases. Up to eight non-solved cases gained new insights due to the reevaluation of our prioritized variants. Our variant prioritization approach was also applied for VUS reclassification, a major challenge to unlock the diagnosis of rare pathologies [10,30–32]. Indeed, in our IRD subcohort, the top five genes with more deleterious variants in IRD solved cases present a higher degree of uncertainty in variant annotation (the proportion of VUS and deleterious variants) in unsolved IRD cases. For an initial list of 63 VUS whose reclassification may solve a case from the IRD cohort, we found 11 VUS more frequently in IRD cases than in pseudocontrols, and 10 of them (~91%) changed their classification to likely pathogenic or pathogenic by the application of the ACMG PS4 criteria. This reclassification was able to solve two cases and provide new evidence to the other 10 patients.

In addition to the revised cases reported in this work, the disease-specific AF and its comparison to pseudocontrols have been implemented in the variant annotation task of our reanalysis pipeline [33] that performs periodic reanalysis of non-solved cases, as well as of WES and WGS analysis. Thus, the database is expected to contribute to the diagnosis of more patients over time. Furthermore, and in order to facilitate the implementation of the database and adjust it to the cohorts of other clinical settings, we have code and instructions to build an in-house database available at <https://github.com/TBLabFJD/BbofAFs>, accessed on 21 June 2022.

In parallel, we also aimed to highlight genes besides variants. Our method extracts genes having more frequent deleterious variants in IRD cases than in pseudocontrols, weighted by the relative frequency of benign variants in order to increase the disease association signal. Several discovery scenarios may fit into the results of this proposal: first, finding underrated IRD genes with a role in IRD cases, such as *DYNC2H1* and *MYO7A* genes that may carry not previously inspected pathogenic variants (thus, we were able provide new evidence in four cases); next, by providing extra findings in complex cases, either syndromic cases, gene modifiers, or dual diagnosis. For instance, we made a dual diagnosis in a previously non-solved syndromic IRD case with rod-cone dystrophy (HP:0000510) and lens luxation (HP:0012019), among other systemic findings. This case includes a partially solved IRD phenotype with a heterozygous variant in gene *CDH23* (Table 1), and a solved Lens Luxation phenotype with two pathogenic variants in compound heterozygosity in gene *ADAMTSL4*. In non-classical Mendelian scenarios, the exploration of the mutational landscape of IRD may help, for instance, to identify more complex cases as (i) genetic pleiotropy together with causal variants in recessive forms [22], (ii) digenic inheritance [23], or (iii) triallelic sites associated with BBS [24]. Last, there is also the possibility of detecting genes not yet associated with IRD, that is, candidate genes that may become IRD genes if further analyses are performed.

As a result of the reevaluation of IRD non-solved cases using the methods described in this work, we were able to solve four cases (one of them with a dual diagnosis still unsolved) and provide new insights in 20 cases, 15 of them providing a single variant for a recessive case that fits the phenotype (partially solved) and five with candidate variants in genes not yet fully associated with the phenotype (marked here as “with evidence”). Further studies are needed to solve cases annotated as partially solved and with evidence.

An additional interesting use of an internal database of allele frequencies is to have a cohort-specific CF estimation. This analysis can provide a better understanding on how deleterious variants are distributed in a general population and is relevant for their use in a public health strategy for genetic counselling. For instance, in the case of IRD, the gene with a higher carrier frequency is *ABCA4*, with carrier variants in ~7% of the population, which is in line with previous estimations [34]. Considering the curated set of pathogenic variants used by Hanany et al. [19] for *ABCA4*, we obtained a similar CF (~6%). Nevertheless, excluding hypomorphic variants, as recommended in this study, CF drops to ~4%. The high CF obtained for this gene in our cohort can be explained partially by these variants.

It is reasonable to state that although pseudocontrols are suitable for providing a good estimation of general allele frequencies, the lack of healthy controls can be seen as a limitation. However, availing of such a control sample set is not always feasible for clinical settings and should be provided under the umbrella of national plans. The major constraint in the discovery capability of our database is that we are restricted to the ~5000 genes targeted in the clinical exome approaches, and thus implementation using data from whole exomes would be optimal. Our intention is to maintain and expand the database in a number of cases but also in genomic regions. We should mention that the IRD non-solved cases presented in the cohort can also have causal variants in non-coding regions, which are not covered with the clinical exome approach.

Although in this work we focus on IRD as the larger group of diseases in our cohort, the same methodology can be applied to other genetic rare diseases in our cohort as well as in other settings. In addition, the disease-specific subcohort that is the subject of analysis can also be tuned in its composition in order to provide different capabilities to the discovery process. For instance, although our subcohort of IRD unsolved cases did not include suspected recessive cases with only one detected candidate variant in heterozygosity (the so-called monoallelic cases), they could be included in this discovery subcohort, so these candidate variants are also evaluated.

In conclusion, our cohort-specific database of allele frequencies has proven to be able to diagnose non-solved IRD cases, reclassify VUS, propose candidate genes, and calculate CF on genes of interest. We believe that the results shown here can highlight the importance

of the reuse of genomic data produced in clinical settings, where the phenotyping is usually exhaustive and the patients waiting for a diagnosis or a genetic counselling can be directly benefited.

4. Materials and Methods

4.1. Ethics Approval and Consent to Participate

The project was reviewed and approved by the Research Ethics Committee of UH-FJD (Ref. 2016/59) and fulfills the principles of the Declaration of Helsinki and subsequent reviews. All patients signed an informed consent before participating. All samples included in this work were pseudonymized, and genomic data were only treated in aggregation.

4.2. Cohort Description

We retrospectively selected all index cases ($N = 5683$) with a clinical exome test performed as a first-tier approach at the Genetics and Genomics Department of the University Hospital Fundación Jiménez Díaz (UH-FJD, Madrid, Spain) from September 2015 to May 2021. The cohort included patients suffering from genetic diseases classified in 14 categories, with the largest disease group being IRD ($N = 1766$) (Supplementary Table S7). The rest of the diseases were grouped in “other eye-related diseases” (OERD) and “non-related diseases” (NRD). Based on the diagnostic status set by the molecular geneticists after CE inspection, all IRD cases were classified as solved and non-solved. From the non-solved cases, we extract those annotated as “partially solved” or “VUS-cases”.

4.3. Sequencing Tests

Samples were analyzed using targeted DNA sequencing with two different commercial sequencing panels: TruSightOne Sequencing Panel kit (TSO, Illumina, San Diego, CA, USA) and Clinical Exome Solution Sequencing Panel kit (CES, Sophia Genetics, Boston, MA, USA). The CES panel targets a total of 4828 genes and regulatory regions and the TSO panel targets 4813 genes, with an overlap of 3567 genes between both panels (Supplementary Figure S5).

4.4. Bioinformatics Reanalysis

In order to have a homogeneous variant calling and annotation of all sequencing tests, all sequenced data were reanalyzed using a custom bioinformatics pipeline for both single nucleotide variants (SNVs) and small insertions and deletions (indels) [33]. This pipeline included exonic, intronic, and UTR analysis. For variant calling, we included 1000 base pairs’ padding for each target region, for both TSO and CES clinical exome tests. The pipeline is based on the GATK 4.1 variant caller [35] and uses the BWA-MEM aligner [36] to the GRCh37/hg19 reference genome. The following databases were used for annotation: (i) allele frequency: gnomAD [17], 1000 genomes [37], and Kaviar [38]; (ii) pathogenicity prediction: SIFT [39], PolyPhen [40], CADD [41], LRT [42], M-CAP [43], MetaLR [44], MetaSVM [44], MutationAssessor [45], MutationTaster [46], PROVEAN [47], and FATHMM [48]; (iii) splicing prediction: ada_score [49] and rf_score [49]; (iv) ClinVar [50]; (v) conservation: phastCons20way [51,52] and phyloP20way [51,52]; (vi) gene tolerance to loss of function (LoF) variants: LoFtool [53], and ExACpLI [12]; (vii) constrained coding regions by means of gnomAD_CCR [54]; and (viii) potential loss of heterozygosity regions, annotated with PLINK [55]. The pipeline is available at <https://github.com/TBLLabFJD/VariantCallingFJD> (accessed on 1 March 2021).

4.5. Detection and Removal of Sample Duplicates and Cryptic Relatedness

All known sample duplicates and relatives were removed prior to frequency calculation. In order to detect other possible sample duplicates and relatives, a PLINK whole-genome association analysis toolkit [55] was used to calculate inbreeding coefficients (identity-by-descent, IBD). First, single nucleotide polymorphisms (SNPs) pruning was performed removing SNPs covered in less than 95% of the samples (PLINK parameter: geno 0.05), with less than 5% allelic frequency (PLINK parameter: maf 0.05), and in linkage

disequilibrium (PLINK parameters: indep-pairwise 50 5 0.5). With the resulting SNPs, the IBD was calculated for all sample combinations (PLINK parameter: genome). All samples with a PI_HAT score higher than 0.35 were removed.

4.6. Variant Frequency Calculation for IRD Patients and Pseudocontrols

After identifying and removing sample duplicates and relatedness ($N = 5683$), variants (in vcf-format) from the index-cases were processed together and merged into a multi-vcf file. Sequencing coverage was also calculated for each sample to distinguish between non-covered and non-mutated sites.

We developed in-house routines for allelic frequency calculation based on the Hail python library for genomics data exploration and analysis (<https://hail.is>, accessed on 1 March 2021). The allele frequency (AF), allele-number (AN), allele-count (AC), and homozygotes-count were obtained for the general cohort and for several subcohorts composed of IRD cases: (1) all IRD cases, (2) non-solved IRD cases, (3) solved-IRD cases, (4) syndromic-IRD cases, (5) non-syndromic IRD cases, and (6) macular dystrophy cases. To define the subcohort of IRD PC for all IRD subcohorts, we took samples from the subcohort NRD ($N = 3531$). AF, AN, and AC were calculated from this PC subcohort (PC-AF, PC-AN, and PC-AC).

4.7. Definition of Genes Associated with IRD, OERD, and NRD

Three disease-specific gene panels were used in the inspection of variants and genes: (i) the IRD gene panel (244 genes, including 136 genes for syndromic-IRD and 108 genes for non-syndromic IRD) as the virtual gene panel used in the diagnosis of IRD cases in the Genetics and Genomics Department of the UH-FJD, and extracted using RetNet, HGMD, and literature searches (Supplementary Table S8); (ii) OERD genes, including non-IRD genes with ocular phenotype (all genes linked with the HPO term “Eye Disease”—HP:0000478, $N = 1542$ genes, Supplementary Table S9); and (iii) NRD genes (the rest of the genes included in TSO/CES panels, not related to eye diseases, $N = 3260$, Supplementary Table S10). Genes that ACMG recommends reporting in case of secondary findings [56] (Supplementary Table S5) were excluded from the gene panels and analyses, except for the gene *RPE65* that belongs to the IRD panel.

4.8. Variants Discarded for Analysis

Variants detected in the 5683 samples from our general cohort were further filtered out using two criteria: (i) quality filtering—we removed 5% of variants with lowest AN; and (ii) population filtering—in order to discard a population origin bias in our IRD subcohort compared to the rest of the cohort, we keep variants present in IRD solved or non-solved cases, and assuming no differences in population origin between IRD solved and non-solved cases, having a fold change between non-solved-AF and solved-AF > 90% percentile, from them we rescue those having a non-solved-AF and solved-AF < 0.1.

4.9. Determination of Differentially Frequent Variants in IRD Subcohorts Compared to Pseudocontrols

We define differentially frequent variants as those that have a higher frequency in a subcohort compared with a control subcohort. In order to extract variants that are differentially frequent in the IRD subcohorts (solved-IRD, non-solved IRD, syndromic-IRD, non-syndromic IRD, and macular dystrophies) compared to the IRD PC subcohort, we calculated the FC of the AF in the IRD subcohort compared to the PC-AF for each of the variants. Based on the distribution of the log₂ of FCs (log₂(FC)) of all variants, we selected those above the 90% as the significant differential frequent variants in a subcohort, tagged as IRD-MFV (IRD most frequent variants) for any IRD subcohort. Variants annotated by ClinVar as “pathogenic” or “likely pathogenic” or with a CADD_PHRED ≥ 30 (top 0.1% most deleterious variants according to CADD) were classified as deleterious, and variants annotated by ClinVar as “benign” or “likely benign” were classified as benign. We compared the proportion of deleterious variants in the IRD-MFV group with the

rest of the variants (non-prioritized variants) for solved-IRD and non-solved-IRD cases. Furthermore, we performed this comparison grouping IRD cases as syndromic forms, non-syndromic forms, and macular dystrophies. A Fisher's exact test was applied to compare the proportion of deleterious variants in these groups, and a p -value < 0.05 was taken as significant.

4.10. VUS Reclassification

We selected VUS whose reclassification can determine the diagnosis of an IRD case in our cohort. These VUS are reported in the diagnostic process at the Genetics Department of the UH-FJD if no pathogenic or likely-pathogenic variant is found to be associated with the phenotype. Variants are classified using ACMG guides. In the IRD subcohort, 100 VUS were reported to fulfil these criteria [10] (Supplementary Table S11). Of these, 63 fulfilled criteria to be within the database generated in this work and were still classified as VUS according to ACMG (information taken from VarSome at the time of the analysis). For these VUS, we annotated IRD-AF (general IRD cohort) and PC-AF frequencies, calculated the FC for these two frequencies, and selected two sets: 1) VUS with a $\log_2(\text{FC}) \geq 1.5$ ($N = 11$), and VUS with a $\log_2(\text{FC}) \geq 2.48$ (value of the 90th percentile of the distributions of the $\log_2(\text{FC})$, $N = 6$). For all the selected VUS ($N = 11$), we marked the specific ACMG criterion PS4 for which "the prevalence of the variant in affected individuals is significantly increased compared with the prevalence in controls" and applied previous evidence to obtain a new ACMG classification.

4.11. Gene Prioritization for IRD Association

To prioritize genes in non-solved IRD cases, we selected for each gene included in the database (i) deleterious variants annotated by ClinVar as "pathogenic" or "likely-pathogenic" or with a $\text{CADD_PHRED} \geq 30$; and (ii) benign variants annotated by ClinVar as "benign" or "likely benign". Genes with at least five deleterious and five benign variants were selected for further analysis. For each selected variant, a $\log_2(\text{FC})$ was calculated between non-solved IRD-AF and PC-AF. Finally, we applied the Wilcoxon rank sum test to the distribution of $\log_2(\text{FC})$ for deleterious and benign variants in each gene. P -values were adjusted using FDR, and genes with an adjusted p -value < 0.05 were considered significant. This list of significant genes was classified into three different gene panels according to the relation degree with IRDs: (i) IRD gene panel, (ii) OERD gene panel, and (iii) NRD gene panel. This analysis was also performed for solved-IRD cases (Supplementary Figure S6).

4.12. Carrier Frequency Calculation

Carrier frequency (CF) was calculated for genes in the non-syndromic IRD gene panel with at least three solved cases in our cohort. Genes were classified as having autosomal recessive or dominant inheritance patterns using the software DOMINO [57] and OMIM database [58]. In genes annotated as recessive, CF was calculated including the variants classified: (i) "pathogenic" or "likely pathogenic" in ClinVar; (ii) "pathogenic" or "likely pathogenic" in LOVD database; (iii) with a CADD_PHRED score ≥ 30 ; (iv) or frameshift/stop-gain variants. The total AC of the variants selected was divided by the AN, and the result was multiplied by 2 (two alleles) and by 100 to represent the result as a percentage (0–100%), according to Equation (1).

Equation (1). Carrier frequency (CF) calculation.

$$CF = 2 \left(\frac{\sum AC}{\max AN} \right) \quad (1)$$

For the gene *ABCA4*, the CF was also calculated excluding hypomorphic variants, as described in Hanany et al. [19].

Supplementary Materials: The following supporting information can be downloaded at: <https://www.mdpi.com/article/10.3390/ijms23158431/s1>.

Author Contributions: Study concept and design: P.M., C.A. and I.-F.I. Bioinformatics methods: I.-F.I., G.N.-M., L.d.l.F., R.R. and P.M. Software implementation: I.-F.I., G.N.-M., L.d.l.F. and R.R. Data analysis and interpretation: I.-F.I., I.P.-R., A.Á.-F., M.J.T.-T., R.R.-Á., B.A., I.M.-M., M.D.P.-V., M.C., C.A. and P.M. Drafting of the manuscript: I.-F.I. and P.M. Manuscript reviewing and editing: A.D.-V., I.-F.I., I.P.-R., G.N.-M., L.d.l.F., R.R., A.Á.-F., M.J.T.-T., R.R.-Á., B.A., I.M.-M., M.D.P.-V., M.C., C.A. and P.M. Resources: C.A. and P.M. All authors have read and agreed to the published version of the manuscript.

Funding: This work was supported by the Comunidad de Madrid (CAM, RAREGenomics Project, B2017/BMD-3721), Instituto de Salud Carlos III (ISCIII) of the Spanish Ministry of Health (FIS; PI18/00579, PI19/00321, PI20/00851), Ramón Areces Foundation (4019/012), Centro de Investigación Biomédica en Red Enfermedades Raras (CIBERER, 06/07/0036), IIS-FJD BioBank (PT13/0010/0012), the Organización Nacional de Ciegos Españoles (ONCE), the European Regional Development Fund (FEDER), and the University Chair UAM-IIS-FJD of Genomic Medicine. I.-F.I. is supported by a grant from the Comunidad de Madrid (CAM, PEJ-2017-AI/BMD7256) and ISCIII (IMPACT-Data; IMP/00019); I.P.-R. is supported by a PhD studentship from the predoctoral program from ISCIII (FI17/00192); G.N.-M. is supported by a grant from the Comunidad de Madrid (PEJ-2020-AI/BMD-18610); L.d.l.F. is supported by the platform technician contract of ISCIII (CA18/00017); R.R. is supported by a postdoctoral fellowship of the Comunidad de Madrid (2019-T2/BMD-13714); B.A. is supported by a Juan Rodes program from ISCIII (JR17/00020); A.D.-V. is supported by a PhD studentship from the predoctoral program from ISCIII (FI18/00123); and P.M. is supported by a Miguel Servet program contract from ISCIII (CP16/00116, CPII21/00015).

Institutional Review Board Statement: The study was conducted in accordance with the Declaration of Helsinki and approved by the Research Ethics Committee of UH-FJD (Ref. 2016/59).

Informed Consent Statement: Informed consent was obtained from all subjects involved in the study.

Data Availability Statement: The data underlying this article will be shared on reasonable request to the corresponding author.

Acknowledgments: We thank the patients for consenting to the use of their data for the study. We also thank all technical staff in the Genetics Department of the UH-FJD for conducting the sequencing and further analysis.

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Richter, T.; Nestler-Parr, S.; Babela, R.; Khan, Z.M.; Tesoro, T.; Molsen, E.; Hughes, D.A. Rare Disease Terminology and Definitions—A Systematic Global Review: Report of the ISPOR Rare Disease Special Interest Group. *Value Health* **2015**, *18*, 906–914. [[CrossRef](#)] [[PubMed](#)]
2. Ayuso, C.; Millan, J.M. Retinitis pigmentosa and allied conditions today: A paradigm of translational research. *Genome Med.* **2010**, *2*, 34. [[CrossRef](#)] [[PubMed](#)]
3. Wright, A.F.; Chakarova, C.F.; Abd El-Aziz, M.M.; Bhattacharya, S.S. Photoreceptor degeneration: Genetic and mechanistic dissection of a complex trait. *Nat. Rev. Genet.* **2010**, *11*, 273–284. [[CrossRef](#)] [[PubMed](#)]
4. Sanchez-Navarro, I.; Da Silva, L.R.J.R.J.; Blanco-Kelly, F.; Zurita, O.; Sanchez-Bolivar, N.; Villaverde, C.; Lopez-Molina, M.I.I.; Garcia-Sandoval, B.; Tahsin-Swafiri, S.; Minguez, P.; et al. Combining targeted panel-based resequencing and copy-number variation analysis for the diagnosis of inherited syndromic retinopathies and associated ciliopathies. *Sci. Rep.* **2018**, *8*, 5285. [[CrossRef](#)] [[PubMed](#)]
5. Mockel, A.; Perdomo, Y.; Stutzmann, F.; Letsch, J.; Marion, V.; Dollfus, H. Retinal dystrophy in Bardet–Biedl syndrome and related syndromic ciliopathies. *Prog. Retin. Eye Res.* **2011**, *30*, 258–274. [[CrossRef](#)] [[PubMed](#)]
6. Gana, S.; Serpieri, V.; Valente, E.M. Genotype-phenotype correlates in Joubert syndrome: A review. *Am. J. Med. Genet. C Semin. Med. Genet.* **2022**, *90*, 72–88. [[CrossRef](#)] [[PubMed](#)]
7. Del Pozo-Valero, M.; Riveiro-Alvarez, R.; Martin-Merida, I.; Blanco-Kelly, F.; Swafiri, S.; Lorda-Sanchez, I.; Trujillo-Tiebas, M.J.; Carreño, E.; Jimenez-Rolando, B.; Garcia-Sandoval, B.; et al. Impact of Next Generation Sequencing in Unraveling the Genetics of 1036 Spanish Families With Inherited Macular Dystrophies. *Investig. Ophthalmol. Vis. Sci.* **2022**, *63*, 11. [[CrossRef](#)] [[PubMed](#)]
8. Martin-Merida, I.; Aguilera-Garcia, D.; Fernandez-San Jose, P.; Blanco-Kelly, F.; Zurita, O.; Almoguera, B.; Garcia-Sandoval, B.; Avila-Fernandez, A.; Arteche, A.; Minguez, P.; et al. Toward the mutational landscape of autosomal dominant retinitis pigmentosa: A comprehensive analysis of 258 Spanish families. *Investig. Ophthalmol. Vis. Sci.* **2018**, *59*, 2345–2354. [[CrossRef](#)] [[PubMed](#)]

9. Martin-Merida, I.; Avila-Fernandez, A.; Del Pozo-Valero, M.; Blanco-Kelly, F.; Zurita, O.; Perez-Carro, R.; Aguilera-Garcia, D.; Riveiro-Alvarez, R.; Arteché, A.; Trujillo-Tiebas, M.J.; et al. Genomic Landscape of Sporadic Retinitis Pigmentosa: Findings from 877 Spanish Cases. *Ophthalmology* **2019**, *126*, 1181–1188. [[CrossRef](#)]
10. Iancu, I.F.F.; Avila-Fernandez, A.; Arteché, A.; Trujillo-Tiebas, M.J.; Riveiro-Alvarez, R.; Almoguera, B.; Martin-Merida, I.; Del Pozo-Valero, M.; Perea-Romero, I.; Corton, M.; et al. Prioritizing variants of uncertain significance for reclassification using a rule-based algorithm in inherited retinal dystrophies. *NPJ Genom. Med.* **2021**, *6*, 18. [[CrossRef](#)] [[PubMed](#)]
11. Richards, S.; Aziz, N.; Bale, S.; Bick, D.; Das, S.; Gastier-Foster, J.; Grody, W.W.; Hegde, M.; Lyon, E.; Spector, E.; et al. Standards and guidelines for the interpretation of sequence variants: A joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet. Med.* **2015**, *17*, 405–424. [[CrossRef](#)] [[PubMed](#)]
12. Lek, M.; Karczewski, K.J.; Minikel, E.V.; Samocha, K.E.; Banks, E.; Fennell, T.; O'Donnell-Luria, A.H.; Ware, J.S.; Hill, A.J.; Cummings, B.B.; et al. Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **2016**, *536*, 285–291. [[CrossRef](#)] [[PubMed](#)]
13. Auton, A.; Abecasis, G.R.; Altshuler, D.M.; Durbin, R.M.; Bentley, D.R.; Chakravarti, A.; Clark, A.G.; Donnelly, P.; Eichler, E.E.; Flück, P.; et al. A global reference for human genetic variation. *Nature* **2015**, *526*, 68–74. [[PubMed](#)]
14. Mathieson, I.; Reich, D. Differences in the rare variant spectrum among human populations. *PLoS Genet.* **2017**, *13*, e1006581. [[CrossRef](#)] [[PubMed](#)]
15. Dopazo, J.; Amadoz, A.; Bleda, M.; Garcia-Alonso, L.; Aleman, A.; Garcia-Garcia, F.; Rodríguez, J.A.; Daub, J.; Muntane, G.; Rueda, A.; et al. 267 Spanish Exomes Reveal Population-Specific Differences in Disease-Related Genetic Variation. *Mol. Biol. Evol.* **2016**, *33*, 1205–1218. [[CrossRef](#)] [[PubMed](#)]
16. Fridman, H.; Yntema, H.G.G.; Mägi, R.; Andreson, R.; Metspalu, A.; Mezzavilla, M.; Tyler-Smith, C.; Xue, Y.; Carmi, S.; Levy-Lahad, E.; et al. The landscape of autosomal-recessive pathogenic variants in European populations reveals phenotype-specific effects. *Am. J. Hum. Genet.* **2021**, *108*, 608–619. [[CrossRef](#)]
17. Karczewski, K.J.; Francioli, L.C.; Tiao, G.; Cummings, B.B.; Alfoldi, J.; Wang, Q.; Collins, R.L.; Laricchia, K.M.; Ganna, A.; Birnbaum, D.P.; et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **2020**, *581*, 434–443. [[CrossRef](#)]
18. Peña-Chilet, M.; Roldán, G.; Perez-Florido, J.; Ortuño, F.M.; Carmona, R.; Aquino, V.; Lopez-Lopez, D.; Loucera, C.; Fernandez-Rueda, J.L.; Gallego, A.; et al. CSVS, a crowdsourcing database of the Spanish population genetic variability. *Nucleic acids research* **2021**, *49*, D1130–D1137. [[CrossRef](#)] [[PubMed](#)]
19. Hanany, M.; Rivolta, C.; Sharon, D. Worldwide carrier frequency and genetic prevalence of autosomal recessive inherited retinal diseases. *Proc. Natl. Acad. Sci. USA* **2020**, *117*, 2710–2716. [[CrossRef](#)]
20. Wang, C.; Sun, L.; Zheng, H.; Hu, Y.-Q. Detecting multiple variants associated with disease based on sequencing data of case-parent trios. *J. Hum. Genet.* **2016**, *61*, 851–860. [[CrossRef](#)]
21. Cordell, H.J.; Barratt, B.J.; Clayton, D.G. Case/pseudocontrol analysis in genetic association studies: A unified framework for detection of genotype and haplotype associations, gene-gene and gene-environment interactions, and parent-of-origin effects. *Genet. Epidemiol.* **2004**, *26*, 167–185. [[CrossRef](#)] [[PubMed](#)]
22. Molinari, E.; Srivastava, S.; Sayer, J.A.; Ramsbottom, S.A. From disease modelling to personalised therapy in patients with CEP290 mutations. *F1000Research* **2017**, *6*, 669. [[CrossRef](#)] [[PubMed](#)]
23. Gao, F.-J.; Zhang, S.-H.; Chen, J.-Y.; Xu, G.-Z.; Wu, J.-H. Digenic heterozygous mutations in EYS/LRP5 in a Chinese family with retinitis pigmentosa. *Int. J. Ophthalmol.* **2017**, *10*, 325–328. [[PubMed](#)]
24. Katsanis, N.; Ansley, S.J.; Badano, J.L.; Eichers, E.R.; Lewis, R.A.; Hoskins, B.E.; Scambler, P.J.; Davidson, W.S.; Beales, P.L.; Lupski, J.R. Triallelic Inheritance in Bardet-Biedl Syndrome, a Mendelian Recessive Disorder. *Science* **2001**, *293*, 2256–2259. [[CrossRef](#)] [[PubMed](#)]
25. Perea-Romero, I.; Gordo, G.; Iancu, I.F.I.F.F.; Del Pozo-Valero, M.; Almoguera, B.; Blanco-Kelly, F.; Carreño, E.; Jimenez-Rolando, B.; Lopez-Rodriguez, R.; Lorda-Sanchez, I.; et al. Genetic landscape of 6089 inherited retinal dystrophies affected cases in Spain and their therapeutic and extended epidemiological implications. *Sci. Rep.* **2021**, *11*, 1526. [[CrossRef](#)] [[PubMed](#)]
26. Maffucci, P.; Bigio, B.; Rapaport, F.; Cobat, A.; Borghesi, A.; Lopez, M.; Patin, E.; Bolze, A.; Shang, L.; Bendavid, M.; et al. Blacklisting variants common in private cohorts but not in public databases optimizes human exome analysis. *Proc. Natl. Acad. Sci. USA* **2019**, *116*, 950–959. [[CrossRef](#)]
27. Musacchia, F.; Ciolfi, A.; Mutarelli, M.; Bruselles, A.; Castello, R.; Pinelli, M.; Basu, S.; Banfi, S.; Casari, G.; Tartaglia, M.; et al. VarGenius executes cohort-level DNA-seq variant calling and annotation and allows to manage the resulting data through a PostgreSQL database. *BMC Bioinform.* **2018**, *19*, 477. [[CrossRef](#)]
28. Riordan, J.D.; Nadeau, J.H. From Peas to Disease: Modifier Genes, Network Resilience, and the Genetics of Health. *Am. J. Hum. Genet.* **2017**, *101*, 177–191. [[CrossRef](#)] [[PubMed](#)]
29. Katsanis, N. The continuum of causality in human genetic disorders. *Genome Biol.* **2016**, *17*, 233. [[CrossRef](#)]
30. Mahecha, D.; Nuñez, H.; Lattig, M.C.; Duitama, J. Machine learning models for accurate prioritization of variants of uncertain significance. *Hum. Mutat.* **2022**, *43*, 449–460. [[CrossRef](#)]
31. Nicora, G.; Zucca, S.; Limongelli, I.; Bellazzi, R.; Magni, P. A machine learning approach based on ACMG/AMP guidelines for genomic variant classification and prioritization. *Sci. Rep.* **2022**, *12*, 2517. [[CrossRef](#)] [[PubMed](#)]

32. Castillo-Guardiola, V.; Rosado-Jiménez, L.; Sarabia-Meseguer, M.D.; Marín-Vera, M.; Macías-Cerrolaza, J.A.; García-Hernández, R.; Zafra-Poves, M.; Sánchez-Henarejos, P.; Moreno-Locubiche, M.Á.; Cuevas-Tortosa, E.; et al. Next step in molecular genetics of hereditary breast/ovarian cancer: Multigene panel testing in clinical actionably genes and prioritization algorithms in the study of variants of uncertain significance. *Eur. J. Med. Genet.* **2022**, *65*, 104468. [[CrossRef](#)] [[PubMed](#)]
33. Romero, R.; de la Fuente, L.; Del Pozo-Valero, M.; Riveiro-Álvarez, R.; Trujillo-Tiebas, M.J.; Martín-Mérida, I.; Ávila-Fernández, A.; Iancu, I.-F.; Perea-Romero, I.; Núñez-Moreno, G.; et al. An evaluation of pipelines for DNA variant detection can guide a reanalysis protocol to increase the diagnostic ratio of genetic diseases. *NPJ Genom. Med.* **2021**, *7*, 7. [[CrossRef](#)]
34. Riveiro-Álvarez, R.; Aguirre-Lamban, J.; Lopez-Martinez, M.A.; Trujillo-Tiebas, M.J.; Cantalapiedra, D.; Vallespin, E.; Avila-Fernandez, A.; Ramos, C.; Ayuso, C. Frequency of ABCA4 mutations in 278 Spanish controls: An insight into the prevalence of autosomal recessive Stargardt disease. *Br. J. Ophthalmol.* **2009**, *93*, 1359–1364. [[CrossRef](#)] [[PubMed](#)]
35. McKenna, A.; Hanna, M.; Banks, E.; Sivachenko, A.; Cibulskis, K.; Kernytzky, A.; Garimella, K.; Altshuler, D.; Gabriel, S.; Daly, M.; et al. The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **2010**, *20*, 1297–1303. [[CrossRef](#)] [[PubMed](#)]
36. Li, H.; Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **2009**, *25*, 1754–1760. [[CrossRef](#)] [[PubMed](#)]
37. Durbin, R.M.; Altshuler, D.L.; Abecasis, G.R.; Bentley, D.R.; Chakravarti, A.; Clark, A.G.; Collins, F.S.; De La Vega, F.M.; Donnelly, P.; Egholm, M.; et al. A map of human genome variation from population-scale sequencing. *Nature* **2010**, *467*, 1061–1073.
38. Glusman, G.; Caballero, J.; Mauldin, D.E.; Hood, L.; Roach, J.C. Kaviar: An accessible system for testing SNV novelty. *Bioinformatics* **2011**, *27*, 3216–3217. [[CrossRef](#)]
39. Ng, P.C.C.; Henikoff, S. SIFT: Predicting amino acid changes that affect protein function. *Nucleic Acids Res.* **2003**, *31*, 3812–3814. [[CrossRef](#)]
40. Adzhubei, I.; Jordan, D.M.; Sunyaev, S.R. Predicting functional effect of human missense mutations using PolyPhen-2. *Curr. Protoc. Hum. Genet.* **2013**, *76*, 7.20.1–7.20.41. [[CrossRef](#)]
41. Rentzsch, P.; Witten, D.; Cooper, G.M.; Shendure, J.; Kircher, M. CADD: Predicting the deleteriousness of variants throughout the human genome. *Nucleic Acids Res.* **2019**, *47*, D886–D894. [[CrossRef](#)]
42. Chun, S.; Fay, J.C. Identification of deleterious mutations within three human genomes. *Genome Res.* **2009**, *19*, 1553–1561. [[CrossRef](#)] [[PubMed](#)]
43. Jagadeesh, K.A.; Wenger, A.M.; Berger, M.J.; Guturu, H.; Stenson, P.D.; Cooper, D.N.; Bernstein, J.A.; Bejerano, G. M-CAP eliminates a majority of variants of uncertain significance in clinical exomes at high sensitivity. *Nat. Genet.* **2016**, *48*, 1581–1586. [[CrossRef](#)] [[PubMed](#)]
44. Dong, C.; Wei, P.; Jian, X.; Gibbs, R.; Boerwinkle, E.; Wang, K.; Liu, X. Comparison and integration of deleteriousness prediction methods for nonsynonymous SNVs in whole exome sequencing studies. *Hum. Mol. Genet.* **2015**, *24*, 2125–2137. [[CrossRef](#)] [[PubMed](#)]
45. Reva, B.; Antipin, Y.; Sander, C. Determinants of protein function revealed by combinatorial entropy optimization. *Genome Biol.* **2007**, *8*, R232. [[CrossRef](#)]
46. Schwarz, J.M.; Rödelsperger, C.; Schuelke, M.; Seelow, D. MutationTaster evaluates disease-causing potential of sequence alterations. *Nat. Methods* **2010**, *7*, 575–576. [[CrossRef](#)]
47. Choi, Y.; Sims, G.E.; Murphy, S.; Miller, J.R.; Chan, A.P. Predicting the functional effect of amino acid substitutions and indels. *PLoS ONE* **2012**, *7*, e46688. [[CrossRef](#)] [[PubMed](#)]
48. Shihab, H.A.; Gough, J.; Cooper, D.N.; Stenson, P.D.; Barker, G.L.A.; Edwards, K.J.; Day, I.N.M.; Gaunt, T.R. Predicting the functional, molecular, and phenotypic consequences of amino acid substitutions using hidden Markov models. *Hum. Mutat.* **2013**, *34*, 57–65. [[CrossRef](#)]
49. Jian, X.; Boerwinkle, E.; Liu, X. In silico prediction of splice-altering single nucleotide variants in the human genome. *Nucleic Acids Res.* **2014**, *42*, 13534–13544. [[CrossRef](#)]
50. Landrum, M.J.; Lee, J.M.; Benson, M.; Brown, G.; Chao, C.; Chitipiralla, S.; Gu, B.; Hart, J.; Hoffman, D.; Hoover, J.; et al. ClinVar: Public archive of interpretations of clinically relevant variants. *Nucleic Acids Res.* **2016**, *44*, D862–D868. [[CrossRef](#)]
51. Pollard, K.S.; Hubisz, M.J.; Rosenbloom, K.R.; Siepel, A. Detection of nonneutral substitution rates on mammalian phylogenies. *Genome Res.* **2010**, *20*, 110–121. [[CrossRef](#)] [[PubMed](#)]
52. Siepel, A.; Bejerano, G.; Pedersen, J.S.; Hinrichs, A.S.; Hou, M.; Rosenbloom, K.; Clawson, H.; Spieth, J.; Hillier, L.W.; Richards, S.; et al. Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.* **2005**, *15*, 1034–1050. [[CrossRef](#)] [[PubMed](#)]
53. Fadista, J.; Oskolkov, N.; Hansson, O.; Groop, L. LoFtool: A gene intolerance score based on loss-of-function variants in 60 706 individuals. *Bioinformatics* **2017**, *33*, 471–474. [[CrossRef](#)] [[PubMed](#)]
54. Havrilla, J.M.; Pedersen, B.S.; Layer, R.M.; Quinlan, A.R. A map of constrained coding regions in the human genome. *Nat. Genet.* **2018**, *51*, 88–95. [[CrossRef](#)] [[PubMed](#)]
55. Purcell, S.; Neale, B.; Todd-Brown, K.; Thomas, L.; Ferreira, M.A.R.; Bender, D.; Maller, J.; Sklar, P.; de Bakker, P.I.W.; Daly, M.J.; et al. PLINK: A tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **2007**, *81*, 559–575. [[CrossRef](#)]

56. Miller, D.T.; Lee, K.; Chung, W.K.; Gordon, A.S.; Herman, G.E.; Klein, T.E.; Stewart, D.R.; Amendola, L.M.; Adelman, K.; Bale, S.J.; et al. ACMG SF v3.0 list for reporting of secondary findings in clinical exome and genome sequencing: A policy statement of the American College of Medical Genetics and Genomics (ACMG). *Genet. Med.* **2021**, *23*, 1381–1390. [[CrossRef](#)] [[PubMed](#)]
57. Quinodoz, M.; Royer-Bertrand, B.; Cisarova, K.; Di Gioia, S.A.; Superti-Furga, A.; Rivolta, C. DOMINO: Using Machine Learning to Predict Genes Associated with Dominant Disorders. *Am. J. Hum. Genet.* **2017**, *101*, 623–629. [[CrossRef](#)] [[PubMed](#)]
58. Hamosh, A.; Scott, A.F.; Amberger, J.; Bocchini, C.; Valle, D.; McKusick, V.A. Online Mendelian Inheritance in Man (OMIM), a knowledgebase of human genes and genetic disorders. *Nucleic Acids Res.* **2002**, *30*, 52–55. [[CrossRef](#)]