

Gene-specific metrics to facilitate identification of disease genes for molecular diagnosis in patient genomes: a systematic review

Journal:	<i>Briefings in Functional Genomics</i>
Manuscript ID	BFGP-18-0028.R1
Manuscript Type:	Review Paper
Date Submitted by the Author:	30-Aug-2018
Complete List of Authors:	Alyousfi, Dareen; University of Southampton, Human Genetics Collins, Andrew; University of Southampton, Human Genetics Baralle, Diana; University of Southampton, Human Genetics and Genomic Medicine Group
Keywords:	gene-specific metrics, disease genome, gene-level scores, gene essentiality, gene-specific filtering

SCHOLARONE™
Manuscripts

1
2
3 **Gene-specific metrics to facilitate identification of disease genes for**
4 **molecular diagnosis in patient genomes: a systematic review**

5 Dareen Alyousfi¹, Diana Baralle², *Andrew Collins¹
6
7
8
9

10 1. Genetic Epidemiology and Bioinformatics Research Group
11 University of Southampton, Southampton SO17 1BJ, UK
12

13
14 ²[Human Development and Health, Faculty of Medicine, University of Southampton, SO16 6YD, UK and Wessex Clinical
15 Genetics Service, Princess Anne Hospital, Southampton, SO16 5YA.]
16

17 University of Southampton,etc.
18

19
20 Word count: 4673 words
21
22
23
24

25
26 *Corresponding author email: arc@soton.ac.uk
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1 Abstract

The evolution of next-generation sequencing (NGS) technologies has facilitated the detection of causal genetic variants in diseases previously undiagnosed at a molecular level. However, in genome sequencing studies, the identification of disease genes among a candidate gene list is often difficult because of the large number of apparently damaging (but usually neutral) variants. A number of *variant* prioritization tools have been developed to help detect disease-causal sites. However, the results may be misleading as many variants scored as damaging by these tools are often tolerated, and there are inconsistencies in prediction results among the different variant-level prediction tools. Recently, studies have indicated that understanding *gene* properties might improve detection of genes liable to have associated disease variation and that this information improves molecular diagnostics. The purpose of this systematic review is to evaluate how understanding gene-specific properties might improve filtering strategies in clinical sequence data to prioritise potential disease variants. Improved understanding of the “disease genome”, which includes coding, non-coding and regulatory variation, might help resolve difficult cases. This review provides a comprehensive assessment of existing gene-level approaches, the relationships between measures of gene-pathogenicity and how use of these prediction tools can be developed for molecular diagnostics.

Key words: gene-specific metrics; disease genome; gene-level scores; gene essentiality; gene-specific filtering.

Author profiles

Dareen Alyousfi, Bachelor of Medicine and Bachelor of Surgery (MBBS), MSc Genomic Medicine, University of Southampton, UK and is undertaking a PhD in human genetics in development and application of methods for resolving molecular diagnoses from patient sequence data.

Diana Baralle is a Professor Genomic Medicine and Honorary Consultant in Clinical Genetics, Faculty of Medicine, University of Southampton and Programme Lead for the MSc in Genomic Medicine.

Andrew Collins is head of the Genetic Epidemiology and Bioinformatics Research Group at the University of Southampton and is involved in next generation sequencing-based studies in population genetics and research into the genetic basis of a number of diseases.

2 Introduction

The sequencing of whole genomes using next generation sequencing (NGS) yields vast datasets which present significant analytical challenges for identification of disease-causal variants. It is known that a subset of human genes contain, or are associated with, rare and/or common variation which have a role in disease processes (the “disease genome”). However, recognition of causal variants amongst many thousands of mostly neutral variants is a huge challenge and a pressing problem. For example, Chong et al (2015) state that the genes underlying ~50% of all Mendelian phenotypes remain unknown and many more Mendelian conditions are still to be described.

1
2
3 Alongside methods for predicting the potential pathogenicity of individual DNA variants a
4 number of gene-specific metrics (scores) have been developed in recent years which may
5 help facilitate recognition of disease causing variation. Understanding the properties of the
6 disease genome and integrating existing gene-specific predictors may help in classifying
7 genes based on their specific features to refine molecular diagnosis. Pathogenicity scores for
8 individual DNA variants are often inconsistent in that different methods can provide
9 conflicting evidence on potential pathogenicity. The degree of redundancy in the genome
10 makes the task of picking out causal variation particularly challenging. We propose an
11 integrated approach which evaluates evidence at both gene and variant levels. We recognize
12 that variant prediction tools alone are currently not conclusive and that evidence at the gene-
13 specific level has the potential to enhance the recognition of variant pathogenicity [1].
14
15
16

17
18 This systematic review considers the literature related to gene-specific scores and their
19 applicability to improve filtering of genome sequence data. We set out to achieve a
20 satisfactory answer to the following research question: “Can the use of gene-specific metrics
21 facilitate the identification of disease genes in patient genomes?” Details of the methodology
22 used in this systematic review are given in the Supplementary methods, Supplementary
23 Figures 1 and 2 and Supplementary Table 1.
24
25

26 **Findings: Key Models**

27
28 From a set of 20 papers yielded by the systematic review methods were classified into three
29 groups determined by the main focus of each method and the corresponding scores: (i)
30 Essentiality and conservation (ii) Haploinsufficiency (iii) Selection.
31
32

33 **4.1 Characteristics of essential and conserved genes.**

34
35
36 Essential and conserved genes encode proteins which have core biological functions essential
37 for an organism’s viability. Genes vary in their degree of essentiality and a number of
38 quantitative scores provide an approximation to essentiality. These include predictions of the
39 extent to which a gene is tolerant or intolerant of loss of function (LoF) mutations and
40 estimation of the expected rate of *de novo* mutations (Pengelly et al, 2017) [11].
41 Supplementary Table 2 outlines the key approaches in this category. The Residual Variation
42 Intolerance Score (RVIS) (Petrovski et al.) ranks genes by probability of carrying more, or
43 less, functional genetic variation than expected highlighting genes intolerant to common
44 functional variation [12]. Genes with positive scores have more common functional variation,
45 while negative scoring genes are less tolerant having reduced associated common functional
46 variation. Genes containing variation involved in monogenic diseases have lower RVIS
47 scores than other genes.
48
49
50

51
52 By examining the evolutionary conservation of protein sequences, Rackham et al. built the
53 **Evolutionary inTolerance score (EvoTol)** to identify genes which are intolerant to
54 mutation [13] [14]. Because only small areas of a gene may be intolerant, for example
55 protein-coding domains, these sub-regions may be considered particularly essential [14].
56
57
58
59
60

1
2
3 EvoTol allows identification of intolerant protein sub-domains alongside the identification of
4 intolerant genes more generally.
5

6
7 The development of NGS makes possible the identification of newly arising (*de novo*)
8 mutations (DNMs) and their potential roles in rare disease. Such mutations are not considered
9 to play a significant role in the pathogenesis of complex diseases [15]. To accurately
10 estimate the expected rate of *de novo* mutations in a given gene, careful assessment of gene
11 mutability is required. Gene length and local sequence context are essential factors
12 underlying mutation rate differences (11). Samocha et al. calculated per-gene probabilities of
13 mutation which are correlated with observed counts of rare missense variants in the Exome
14 Sequencing Project (ESP) data set. The Samocha et al. study extends a model which
15 investigated *de novo* mutations in epileptic encephalopathy patients (Epi4K consortium) by
16 considering depth of coverage (i.e., how many sequence reads were present on average per
17 base) and the regional divergence in genes between humans and Macaques. Significant
18 numbers of genes with missense variant deficits were observed, compared to expectation
19 from predicted mutation rates, suggesting strong evolutionary constraint removing variants
20 by negative selection [15] [16]. The Samocha et al. model utilizes exome sequence data to
21 evaluate the DNM rate by gene set and on a single gene basis [15], this score is referred to
22 as *de novo* excess (DNE). The metric is predictive of selective constraint in the human
23 genome and they identified 1,003 constrained genes known to cause severe human
24 disease[15]. It was found that constrained genes contain higher *de novo* LoF mutation rate
25 than expected by chance[15].
26
27
28
29
30

31 The LoFtool measures the ratio of LoF mutations to synonymous mutations for every gene.
32 The performance of LoFtool, compared to RVIS, DNE Z-score, and EvoTol, suggests
33 enhanced performance for predicting *de novo* haploinsufficient disease-causing genes. The
34 LoFtool represents values as intolerance percentiles: genes that are intolerant to LoF variation
35 have low LoFtool percentiles [13]. The four measures of genic intolerance outlined so far
36 were included by Bartha et al. who described them as essentiality scores [17].
37
38
39

40 In early 2016, using data from 1000 Genomes Project, Aggarwala et al. proposed the
41 Substitution Intolerance Score (SIS) as a gene-level measurement of essentiality. The
42 interpretation of this score is such that genes with high SIS scores are functionally
43 constrained, while genes which score low are tolerant of functional changes in the protein
44 which might arise through mutations in the DNA sequence [18].
45
46

47 Another scoring system by Gussow et al. evaluates intolerance in genic sub-regions
48 proposing that more conserved regions within a gene are expected to contain more variants
49 which are pathogenic [19]. Genes are divided into sub-regions and tiered by intolerance to
50 functional variation. This 'subRVIS' score ranks regions using RVIS but with the addition of
51 information on conservation. Regions intolerant to functional variation are scored low by the
52 subRVIS scoring system. The method utilizes the GERP++ score to evaluate evolutionary
53 constraint for bases in each sub-region [19].
54
55
56
57
58
59
60

1
2
3 The Loss Intolerance probability (pLI) score quantifies the likelihood that a gene is
4 intolerant to a mutation which produces LoF in the protein product [20]. The score is derived
5 using the Exome Aggregation Consortium (ExAC) database which is an extensive catalogue
6 of human genetic diversity. This catalogue identifies one variant every eight bases on average
7 in the exome providing a powerful filter for analysis of candidate deleterious variants in
8 severe Mendelian diseases [20]. Lek et al. proposed that genes with high pLI score (pLI \geq
9 0.9) are most intolerant of LoF variation. Genes in this category are the most evolutionarily
10 constrained. The least constrained genes (LoF tolerant) have low pLI scores (pLI \leq 0.1) and
11 typically contribute to the least constrained biological pathways, such as sensory perception,
12 where high haplotype diversity is potentially advantageous [20].
13
14

15
16 It is challenging to assess the relationship between the DNM rate and genes involved in
17 disease. In 2017, Jiang et al. utilized available DNM data to correct for the background
18 mutation rate seen as one of the main limitations in the Samocha et al.[15] work. The
19 problem arises because by sequencing more individuals, more DNMs are inevitably observed
20 in the same gene by chance. Therefore, in a given disease, if a *de novo* mutation is related to
21 pathogenesis, disease-genes might be expected to contain more DNMs than predicted from
22 background rates. This work includes the development of a database which describes the
23 background DNM rate (DNMR), acquired from population variation data [21].
24
25
26
27

28 **4.2 Characteristics of Haploinsufficient genes**

29
30 Haploinsufficiency (HI) occurs whenever there is a missing or damaged copy of a **gene**
31 leaving a single copy insufficient to maintain normal function [22]. Haploinsufficiency is
32 mostly caused by LoF mutations and results in dominant diseases. Recognition and prediction
33 of genes which are haploinsufficient can facilitate the filtering of disease genome data
34 wherever the phenotype is likely to have arisen through reduced levels of gene product.
35
36

37 In 2010, Haung et al. proposed a deletion-based HI score by identifying differences between
38 HI and haplosufficient (HS) genes, aiming to better distinguish pathogenic from benign
39 deletions which helps in variant prioritization [22]. The analysis develops a logarithm-of-
40 odds (LOD) score to estimate the probability of a deletion causing a HI phenotype. A high
41 LOD score suggests deletions are likely to be deleterious through HI and therefore potential
42 candidates for causing dominant traits. The score assumes there are no statistical interactions
43 between the genes [22]. Previously, and to try to assess the pathogenicity of a deletion,
44 clinicians considered the length of a deletion or the number of genes deleted. The Haung et
45 al. score provides a rational basis to classify pathogenic deletions by comparing deletions
46 seen in patients with deletions in controls and calculating the fraction of controls with a
47 deletion at least as deleterious as that seen in the patient [22].
48
49
50
51

52
53 To distinguish false-positive disease variants from the genuinely causal variants is crucial for
54 accurate molecular diagnoses. MacArthur et al. developed the REcessive (REC) score for
55 distinguishing genes involved in recessive diseases from genes which are LoF- variation
56
57
58
59
60

1
2
3 tolerant [23]. A “healthy” genome might contain 100 true LoF variants, the majority in a
4 heterozygous state. Evidence suggests that the average human carries five recessive lethal
5 alleles in single copy in their genome. Consequently, the majority of LoF variants are
6 considered common variants. However, these variants might still have a phenotypic effect
7 [23]. MacArthur et al. demonstrated differences in functional and evolutionary features
8 between recessive disease and LoF-tolerant genes, allowing for the development of a
9 predictive model to predict recessive disease variants [23].
10
11

12
13 Khurana et al. developed the “gene position in NETworks” (NET) indispensability score to
14 investigate relationships between degree of network centrality of a gene and selection within
15 biological networks [24]. They consider a range of biological networks (i.e., phosphorylation,
16 signaling, protein-protein interaction, regulatory and genetic networks). Genes which are
17 highly connected to many biological networks are the most functionally significant, therefore,
18 mutations in those genes might have serious consequences[24]. However, genes connected
19 to metabolic networks were found to have more duplicated copies through more paralogs
20 with more LoF mutations[24]. This score was included as a predictor of haploinsufficient
21 genes in the Hsu et al. study [2]
22
23

24
25 Ge et al. consider gene-specific pathogenicity using the ratio of non-synonymous to
26 synonymous substitution rates (dN/dS) for X-chromosome genes [25]. Genes with unusually
27 low ratios suggest intolerance to non-synonymous variation, suggesting these are susceptible
28 to disease related variation. They found correlation between genomic regions depleted for
29 missense variation with disease-causal variants [25].
30
31

32
33 Steinberg et al. proposed that study biases existing in many biological networks might affect
34 the ability of previous HI prediction scores to recognize the genuinely haploinsufficient
35 genes. For that reason they constructed a new, unbiased, HI score, the Genome-wide
36 HaploInsufficiency Score (GHIS) which replaces biological networks with co-expression
37 networks [26] [27]. They compared their model with the three pre-existing methods (i.e., HI
38 [22], NET [24] and RVIS [12]) and demonstrated that GHIS provides a score for many genes
39 not scored by other methods [26] with enhanced performance at classifying less well
40 studied genes [26].
41
42

43
44 Scores have been developed to recognize Mendelian genes with different modes of
45 inheritance. Hsu et al. considered Mendelian disease gene characteristics according to their
46 mode of inheritance. Haploinsufficiency is an essential characteristic of Mendelian disease
47 genes with an autosomal dominant (AD) mode of inheritance and sensitivity to *de novo*
48 mutations was recognized for this group of genes [2]. In contrast disease genes with
49 autosomal recessive (AR) modes of inheritance tend to have more non-synonymous variants
50 and regulatory transcript isoforms [2]. However, the X-linked (XL) pattern of inheritance is
51 associated with fewer non-synonymous and synonymous variants [2]. Based on these
52 findings they create a new approach to prioritize Mendelian disease genes based on their
53 mode of inheritance (AD, AR, and XL) termed Inheritance-mode Specific Pathogenicity
54 Prioritization (ISPP) [2]. This score integrates pre-existing gene-specific prediction methods
55 namely: HI (Huang et al., 2010) [23], REC (MacArthur et al., 2012) [24], RVIS (Petrovski et
56
57
58
59
60

1
2
3 al., 2013) [13], NET (Khurana et al., 2013) [25], DNE (Samocha et al., 2014) [16] and GDI
4 (Itan et al., 2015) [35] along with numerous genetic properties including global expression
5 from RNA-Seq data, DNA replication time and the noncoding (intronic region) mutation rate
6 [2].
7

8 Because the human genome contains an abundance of non-deleterious heterozygous variants,
9 the identification of dominant mutations for monogenic disorders is challenging. Quinodoz et
10 al. created DOMINO a method using machine learning to identify whether a given gene is
11 liable to carry dominant changes [28].
12
13

14 Inevitably, well-studied genes are over-represented in most biological networks used to
15 create scores that predict HI compared to less-studied genes, hence most biological networks
16 are affected by study bias. Therefore the creation of unbiased HI score becomes essential[27].
17 Recently, Shibab et al. produced an integrated machine learning approach called (HIPred)
18 merging functional annotations with genomic and evolutionary features to predict HI genes
19 without study bias using data from NIH Roadmap Epigenomics [29] and the ENCODE [30]
20 project. The performance of this approach is considered to exceed the pre-existing HI
21 predictors [27]. Supplementary Table 3 outlines the key approaches in this category.
22
23
24
25

26 **4.3 Characteristics of genes under selection.**

27
28
29 Genetic variants may be subject to positive selection whereby, if they are advantageous,
30 they may increase in frequency. Negative selection, in contrast, acts to remove deleterious
31 alleles. Scores which quantify the intensity of negative selection acting on genes provide
32 insights into which genes are more likely to have variation which may have damaging
33 consequences. The pattern is complex because some essential genes are not known to have
34 any associated disease variation and are perhaps subject to negative selection at
35 particularly high intensity [31].
36
37

38 Bustamante et al. calculate the extent and directionality of Selection operating on a given
39 gene, this score referred to here as “Sel”. They first compared fixed sequence differences,
40 both synonymous and non-synonymous, between humans in the sample and Chimpanzees
41 over 11.81 Mb region of aligned coding DNA. The ratio of non-synonymous to
42 synonymous differences (divergence) was 23.76%. In contrast the ratio of non-
43 synonymous to synonymous polymorphisms in the human subjects was 38.42%. This
44 shows a significant excess of amino acid variation, relative to divergence, consistent with
45 previous work stating that much amino acid variation in the human genome is slightly to
46 moderately damaging [32].
47
48
49

50 Eilertson et al. create a model to identify genes under natural selection with a non-parametric
51 approach (with no assumption of a specific population genetic model) which is robust to
52 demography [33]. This approach, called Selection Inference using Poisson Random Effects
53 (SnIPRE), utilizes polymorphism and divergence data from synonymous and non-
54 synonymous sites within genes [33].
55
56
57
58
59
60

1
2
3 The Gene-level Integrated Metric of negative Selection (GIMS) was created by combining
4 two meta-analyses into a single meta-analysis. The first meta-analysis combines comparative
5 genomic metrics (GERP++) and functional genomic metrics (Poly-phen2), and the second
6 meta-analysis combines mutation rates (as SNPs/kb) and allele frequencies (as % rare) from
7 the 1000 Genomes Project. Meta-analysis was achieved by combining those metrics into
8 GIMS scores for 20,079 genes [34]. Because the majority of genes are under purifying
9 selection, the aim was to quantify the degree of negative selection applied to genes.
10 Conservation and functional scores were initially combined as ‘functional genomic metrics’
11 integrated with mutation rates and fraction of rare variants as ‘population genetic metrics’.
12 The GIMS score combines these two metrics and provides a unified score per-gene. GIMS
13 gives a probability distribution across the entire genome in quantiles. Genes under negative
14 selection are scored low by GIMS [34].
15
16
17
18

19 The Gene Damage Index (GDI) is a gene-specific score which predicts the liability of a
20 human protein-coding gene to contain disease-causing mutations considering the influences
21 of selection and genetic drift. In GDI, Combined Annotation Dependent Depletion (CADD)
22 scores are used as the variant-level damage prediction method because this method is
23 efficient at distinguishing between benign and deleterious variants and is strongly dependent
24 on evolutionary conservation [35]. Moreover, CADD scores can assess most types of variants
25 while other methods, like Poly-Phen-2 and SIFT, can only predict missense variants. To
26 construct the GDI score the cumulative predicted damage in exonic regions of the gene is
27 calculated using the CADD score for each allele compared to the expected score for variants
28 with similar allele frequencies. The homogenized Phred I-score is calculated for each metric
29 to indicate the ranking of the targeted gene relative to all other genes. A low Phred score:
30 indicates a human gene with a low GDI and high Phred score indicates a gene susceptible to
31 contain damaging variation. Genes with high GDI tend to be under less intense purifying
32 selective pressure. A low GDI score is associated with highly conserved genes (including
33 genes enriched for ribosome, chemokine signaling proteasome and spliceosome functions)
34 reflecting essentiality. Such genes tend to be under stronger purifying selection than the
35 median selective pressure acting on human genes [35]. Supplementary Table 4 outlines the
36 key approaches in this category.
37
38
39
40
41
42
43

44 **3 Discussion**

45

46 Considering approaches which score genes according to essentiality and conservation the
47 DNE score offers some advantages. The main limitation of DNE is its validity only for
48 interpretation of *de novo* mutations [2] but considers more variables related to mutation rate
49 which goes beyond sequence context compared to other methods like RVIS and Sel. These
50 additional variables include consideration of sequence depth of coverage and regional
51 divergence in genes between humans and Macaques independently, which improve the
52 predictive value of this model [15]. The DNE score has been compared to the RVIS and
53 negative selection score Sel. The comparison showed that DNE and RVIS were equally
54 effective emphasizing the benefits predicted from combining the two scores [15].
55
56
57
58
59
60

1
2
3 The strength of Samocha et al. model is enhanced by incorporation of the depth of coverage
4 (i.e., how many sequence reads were present on average per base) and the regional
5 divergence in genes between humans and Macaques independently. These strengths play a
6 significant role in the improvement of their predictive model. The number of rare
7 synonymous variants in the Exome Sequencing Project (ESP) is shown to be highly
8 correlated with the probability of a synonymous mutation determined by their model.
9
10

11
12 EvoTol was compared to the RVIS and the DNE scores and shown to have increased
13 performance at classifying intolerant genes compared to RVIS. EvoTol was shown to be
14 highly sensitive and more powerful to characterize genes with high pathogenicity [14].
15 Although there was no significant correlation between RVIS and EvoTol, the application of
16 the two scores simultaneously will likely be advantageous [14].
17
18

19 Considering approaches for scoring genes for potential roles in haploinsufficiency
20 phenotypes the HIPred approach has been evaluated against five predictors (HI Score, NET,
21 RVIS, EvoTol and GHIS, Supplementary Tables 2 and 3). HIPred was found to outperform
22 all in predicting HI genes [27]. Using different perspectives across the 26 disease-associated
23 gene lists, Hsu et al. estimates the power of several methods that predict gene pathogenicity
24 showing a substantial positive correlation between HI and REC (correlation $r=0.77$) while
25 the six scores have a moderate relationship with each other ($r=0.46$) [2]. Among these gene
26 scores (DNE, GDI, HI, NET, RVIS, and REC) the best predictor of disease-predisposing
27 genes was the REC score [2]. The performance of ISPP score was significantly superior for
28 prioritizing AR and X-linked disease-associated genes [2]. The REC score is effective at
29 predicting disease-associated genes generally but less successful in discriminating recessive
30 and dominant disease genes [2].
31
32
33
34

35 DNE measures the rate of per-gene *de novo* mutation while RVIS ranks human genes based
36 on the strength and consistency of the purifying selection acting against functional variation.
37 Analysis has shown that GDI and RVIS capture unique sets of reciprocal information from
38 population genetic data [35]. In essence, RVIS reflects selective pressure while DNE is based
39 on *de novo* mutation rate estimates; both methods do not quantitatively estimate the
40 mutational load for a gene in a healthy human population. For this reason, these methods are
41 not optimal for filtering genes with high mutation rates and many residual false positives
42 might be expected. GDI has proved to be the most efficient approach for filtering out false
43 positive variants in genes known to contain damaging variation [35].
44
45
46

47 The Ge et al. X-linked scoring system is not limited by previous gene annotation and the
48 dN/dS ratio can be calculated for any protein-coding gene. This score applies to all X-
49 chromosome protein-coding genes and therefore can assess genes for multiple disease
50 phenotypes [25]. Because the intra-human dN/dS ratio is not specific to the X-chromosome
51 the analysis of more genomic data using dN/dS ratio is recommended for future studies to
52 identify genes which may have disease variation [25].
53
54
55
56
57
58
59
60

1
2
3 This work aims to bring together the growing evidence that gene properties, alongside variant
4 scoring systems, can play an important role in filtering disease sequence data. As healthy
5 individuals can have genetic variants that lead to disruption of protein-coding genes (with no
6 clinical phenotype) [26][27][23][36], challenges remain to distinguish which loss of function
7 variants are associated with disease phenotypes from those that do not cause any functional
8 disturbance [26]. Data from the 1000 Genomes Project show that on average a healthy person
9 might carry 250-300 LoF SNVs (1000 Genomes Project Consortium et al., 2010; The 1000
10 Genomes Project Consortium, 2012) [2].

11
12
13
14 The understanding of human genomes is advanced through the accumulation of sequence
15 data in publically available databases. The ExAC resource provides a potent filter to aid
16 recognition of pathogenic variants in severe Mendelian diseases. Using ExAC for filtering to
17 remove false positive, but plausibly pathogenic, variants decreases the number of candidate
18 protein-altering variants by 7-fold compared to the smaller Exome Sequencing Project
19 database (ESP) which has fewer exome sequences [20].

20
21
22
23 Coupled with the previous evidence, another study suggests that the missense Z score which
24 represents genes rather than variants adds more information than variant-specific Poly-phen2
25 and CADD classifications signifying that gene-level scores of constraints provide more
26 details to variant-level scores in evaluating pathogenicity [20]. Furthermore, Haung et al.
27 contend that variant level scores (e.g., SIFT and POLYPHEN) are limited by lacking the
28 capability to determine, from cross-species alignments, if negative selection at a given site is
29 acting in a recessive, additive or dominant mode [22].

30
31
32 The work proposed by Gussow et al., was based on dividing the genes into sub-regions to
33 identify exactly where the pathogenic mutations are likely to present [19]. This study
34 brought up an important question: is the whole gene the correct unit to judge patterns of
35 intolerance? Future analyses may consider refinements to gene-specific scores which
36 consider within-gene regional patterns of intolerance in more detail.

37
38
39 Another controversial issue is the difficulty in interpretation of benign LoF variants for which
40 the nomenclature is still not unified. It is important to realize that there are overlaps in the
41 interpretation of LoF variants in healthy people. In the literature, all the following categories
42 are represent LoF variants in healthy individuals: true variants that do not seriously disrupt
43 gene function, benign LoF variation in redundant genes, non-deleterious or less-deleterious
44 variants that have an impact on risk of phenotype or disease [23].

45
46
47 Because each genic scoring approach considers only a specific property of genetic
48 architecture, each individual score has limitations. For example, the (i) REC score does not
49 consider dominant disease-predisposing genes (ii) Non-CNV (Non-Copy Number Variation)
50 genetic variants were not included in HI prediction score. (iii) NET score lacks the systematic
51 comparison of different known disease-associated genes (iv) RVIS score does not consider
52 variations in allele frequencies across different populations (v)The DNE score has limited
53 applicability for testing *de novo* mutations. (vi) The GDI score only considers mutation
54
55
56
57
58
59
60

1
2
3 profiles [2]. Furthermore, a major limitation of the GHIS score is that the genetic
4 background in individuals is not considered, which is an important issue since genetic
5 variants do not act in isolation and disturbance of individual genes within a single biological
6 pathway might affect the risk of a disease [26]. Accordingly, this analysis which provides a
7 comprehensive review of each prediction scheme, may help establish new routes for
8 prioritizing disease-causal variants.
9

10
11 Presented here are a range of well-studied gene-specific predictors with various
12 independent genetic properties. Addressing the limitations of each score or perhaps
13 exploiting the developed scores of pathogenicity and combining these scores in an integrated
14 metrics might better predict disease-genes since there is currently no single method that is
15 reliably predictive of gene pathogenicity.
16
17

18
19 Many advances were developed to assess whether a gene is tolerant or intolerant to common
20 functional variation. Initially, scores were developed per gene then studies were published
21 showing that dividing the gene into sub-regions might help in allocating the mutation
22 accurately. At that time all scores that measure genic intolerance required disease knowledge,
23 this limitation was addressed by developing a tool with no prior disease knowledge required,
24 an essential step to better predict genic intolerance.
25
26

27
28 It is hoped that this review highlights existing work to identify and explain different gene-
29 specific pathogenicity predictors, while pointing to the gaps in disease-gene prioritization and
30 annotation issues to facilitate new scores and better prioritization of disease-causal genes.
31

32 **Key points**

- 33
34 1. A wide range of well-established models exist that prioritize genes based on their
35 associated disease variation potential.
- 36 2. Integration of these strategies to represent individual genes could have a significant
37 impact on our understanding of genic properties and the recognition of disease-related
38 functional variation.
- 39 3. Evaluation and comparison of these individual scores and the development of
40 integrated models to enhance NGS filtering strategies in disease genomes is a fertile
41 area for future studies.
42
43
44
45

46 **Funding**

47
48 DA is funded by the Saudi Arabia Cultural Bureau, London, UK.
49 DB is funded by a NIHR Research Professorship.
50
51

52 **References**

53
54
55 [1] X. Huang, J. Lin, and D. Demner-Fushman, "Evaluation of PICO as a knowledge
56 representation for clinical questions.," *AMIA Annu. Symp. Proc.*, pp. 359–63, 2006.
57
58
59

- 1
2
3 [2] I. T. Newsweekly, “Mendelian Disease ; Findings from University of Hong Kong
4 Yields New Data on Mendelian Disease [(ISPP) for human protein coding genes],”
5 vol. 32, no. 20, pp. 2016–2018, 2016.
- 6 [3] B. Kitchenham and S. Charters, “Guidelines for performing Systematic Literature
7 Reviews in Software Engineering,” *Engineering*, vol. 2, p. 1051, 2007.
- 8 [4] S. U. Khan, M. Niazi, and R. Ahmad, “Barriers in the selection of offshore software
9 development outsourcing vendors: An exploratory study using a systematic literature
10 review,” *Inf. Softw. Technol.*, vol. 53, no. 7, pp. 693–706, 2011.
- 11 [5] C. Jalal, Samireh and Wohlin, “Systematic Literature Studies: Database Searches vs.
12 Backward Snowballing Samireh.” .
- 13 [6] D. Badampudi, “Experiences from using snowballing and database searches in
14 systematic literature studies.” .
- 15 [7] J. Gehanno, L. Rollin, and S. Darmoni, “Is the coverage of google scholar enough to
16 be used alone for systematic reviews,” no. December 2009, pp. 0–4, 2013.
- 17 [8] S. Becker, A. Bryman, and H. (Thomas H. Ferguson, *Understanding research for
18 social policy and practice : themes, methods and approaches*. Policy, 2012.
- 19 [9] G. Craswell and M. Poore, *Writing for academic success*. SAGE, 2012.
- 20 [10] C. Thermes, “Ten years of next-generation sequencing technology,” *Trends Genet.*,
21 vol. 30, no. 9, pp. 418–426, 2014.
- 22 [11] R. J. Pengelly, A. Vergara-Lope, D. Alyousfi, M. R. Jabalameli, and A. Collins,
23 “Understanding the disease genome: gene essentiality and the interplay of selection,
24 recombination and mutation,” *Brief. Bioinform.*, no. June, pp. 1–7, 2017.
- 25 [12] S. Petrovski, Q. Wang, E. L. Heinzen, A. S. Allen, and D. B. Goldstein, “Genic
26 Intolerance to Functional Variation and the Interpretation of Personal Genomes,” *PLoS
27 Genet.*, vol. 9, no. 8, 2013.
- 28 [13] J. Fadista, N. Oskolkov, O. Hansson, and L. Groop, “LoFtool: A gene intolerance
29 score based on loss-of-function variants in 60 706 individuals,” *Bioinformatics*, vol.
30 33, no. 4, pp. 471–474, 2017.
- 31 [14] O. J. L. Rackham, H. A. Shihab, M. R. Johnson, and E. Petretto, “EvoTol : a protein-
32 sequence based evolutionary intolerance framework for disease-gene prioritization,”
33 vol. 43, no. 5, 2014.
- 34 [15] K. E. Samocha *et al.*, “A framework for the interpretation of de novo mutation in
35 human disease,” vol. 46, no. 9, pp. 944–950, 2015.
- 36 [16] A. S. Allen *et al.*, “De novo mutations in epileptic encephalopathies,” *Nature*, vol. 501,
37 no. 7466, pp. 217–221, 2013.
- 38 [17] A. Bartha, István & di Iulio, Julia & Venter, J & Telenti, “Human gene essentiality.
39 Nature Reviews Genetics.” 2017.
- 40 [18] V. Aggarwala and B. F. Voight, “An expanded sequence context model broadly
41 explains variability in polymorphism levels across the human genome,” *Nat. Genet.*,
42 vol. 48, no. 4, pp. 349–355, 2016.
- 43 [19] A. B. Gussow, S. Petrovski, Q. Wang, A. S. Allen, and D. B. Goldstein, “The
44 intolerance to functional genetic variation of protein domains predicts the localization
45 of pathogenic mutations within genes,” *Genome Biol.*, vol. 17, no. 1, pp. 1–11, 2016.
- 46 [20] Lek, “Analysis of protein-coding genetic variation in 60,706 humans,” *Nature*, vol.
47 536, no. 7616, pp. 285–291, 2017.
- 48 [21] Y. Jiang *et al.*, “MirDNMR: A gene-centered database of background de novo
49 mutation rates in human,” *Nucleic Acids Res.*, vol. 45, no. D1, pp. D796–D803, 2017.
- 50 [22] N. Huang, I. Lee, E. M. Marcotte, and M. E. Hurles, “Characterising and predicting
51 haploinsufficiency in the human genome,” *PLoS Genet.*, vol. 6, no. 10, pp. 1–11, 2010.
- 52 [23] D. MacArthur, S. Balasubramanian, and A. Frankish, “A Systematic Survey of Loss-
53
54
55
56
57
58
59
60

- of-Function Variants in Human Protein-Coding Genes,” *Science* (80-.), vol. 335, no. 6070, pp. 1–14, 2012.
- [24] E. Khurana, Y. Fu, J. Chen, and M. Gerstein, “Interpretation of Genomic Variants Using a Unified Biological Network Approach,” *PLoS Comput. Biol.*, vol. 9, no. 3, 2013.
- [25] X. Ge, P. Y. Kwok, and J. T. C. Shieh, “Prioritizing genes for X-linked diseases using population exome data,” *Hum. Mol. Genet.*, vol. 24, no. 3, pp. 599–608, 2015.
- [26] J. Steinberg, F. Honti, S. Meader, and C. Webber, “Haploinsufficiency predictions without study bias,” *Nucleic Acids Res.*, vol. 43, no. 15, pp. 1–9, 2015.
- [27] H. A. Shihab, M. F. Rogers, C. Campbell, and T. R. Gaunt, “HIPred: An integrative approach to predicting haploinsufficient genes,” *Bioinformatics*, vol. 33, no. 12, pp. 1751–1757, 2017.
- [28] M. Quinodoz *et al.*, “REPORT DOMINO : Using Machine Learning to Predict Genes Associated with Dominant Disorders,” *Am. J. Hum. Genet.*, vol. 101, no. 4, pp. 623–629, 2017.
- [29] Roadmap Epigenomics Consortium *et al.*, “Integrative analysis of 111 reference human epigenomes,” *Nature*, vol. 518, no. 7539, pp. 317–329, 2015.
- [30] Encode Consortium, N. Carolina, and C. Hill, “For Junk DNA,” *Nature*, vol. 489, no. 7414, pp. 57–74, 2013.
- [31] N. Spataro, J. A. Rodríguez, A. Navarro, and E. Bosch, “Properties of human disease genes and the role of genes linked to Mendelian disorders in complex disease aetiology,” *Hum. Mol. Genet.*, vol. 26, no. 3, pp. 489–500, 2017.
- [32] C. D. Bustamante *et al.*, “Natural selection on protein-coding genes in the human genome,” *Nature*, vol. 437, no. 7062, pp. 1153–1157, 2005.
- [33] K. E. Eilertson, J. G. Booth, and C. D. Bustamante, “SnIPRE: Selection Inference Using a Poisson Random Effects Model,” *PLoS Comput. Biol.*, vol. 8, no. 12, 2012.
- [34] M. G. Sampson, C. E. Gillies, W. Ju, M. Kretzler, and H. M. Kang, “Gene-level integrated metric of negative selection (GIMS) prioritizes candidate genes for nephrotic syndrome,” *PLoS One*, vol. 8, no. 11, pp. 1–9, 2013.
- [35] Y. Itan *et al.*, “The human gene damage index as a gene-level approach to prioritizing exome variants,” *Proc. Natl. Acad. Sci.*, vol. 112, no. 44, pp. 13615–13620, 2015.
- [36] P. C. Ng *et al.*, “Genetic variation in an individual human exome,” *PLoS Genet.*, vol. 4, no. 8, 2008.

1
2
3 **Gene-specific metrics to facilitate identification of disease genes for**
4 **molecular diagnosis in patient genomes: a systematic review**

5 Dareen Alyousfi¹, Diana Baralle², *Andrew Collins¹
6
7
8
9

10 1. Genetic Epidemiology and Bioinformatics Research Group
11 Human Development and Health,
12 Faculty of Medicine,
13 University of Southampton,
14 SO16 6YD, UK
15

16 2. Human Development and Health,
17 Faculty of Medicine,
18 University of Southampton, SO16 6YD, UK and Wessex Clinical Genetics Service,
19 Princess Anne Hospital,
20 Southampton, SO16 5YA.
21
22

23 Word count: 4586 words
24
25

26 *Corresponding author email: arc@soton.ac.uk
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1 Abstract

The evolution of next-generation sequencing (NGS) technologies has facilitated the detection of causal genetic variants in diseases previously undiagnosed at a molecular level. However, in genome sequencing studies, the identification of disease genes among a candidate gene list is often difficult because of the large number of apparently damaging (but usually neutral) variants. A number of *variant* prioritization tools have been developed to help detect disease-causal sites. However, the results may be misleading as many variants scored as damaging by these tools are often tolerated, and there are inconsistencies in prediction results among the different variant-level prediction tools. Recently, studies have indicated that understanding *gene* properties might improve detection of genes liable to have associated disease variation and that this information improves molecular diagnostics. The purpose of this systematic review is to evaluate how understanding gene-specific properties might improve filtering strategies in clinical sequence data to prioritize potential disease variants. Improved understanding of the “disease genome”, which includes coding, non-coding and regulatory variation, might help resolve difficult cases. This review provides a comprehensive assessment of existing gene-level approaches, the relationships between measures of gene-pathogenicity and how use of these prediction tools can be developed for molecular diagnostics.

Key words: gene-specific metrics; disease genome; gene-level scores; gene essentiality; gene-specific filtering.

Author profiles

Dareen Alyousfi , is a Bachelor of Medicine and Bachelor of Surgery (MBBS) and has an MSc in Genomic Medicine from the University of Southampton, UK and is undertaking a PhD in human genetics studying development and application of methods for resolving molecular diagnoses from patient sequence data.

Diana Baralle is Professor Genomic Medicine and Honorary Consultant in Clinical Genetics, Faculty of Medicine, University of Southampton and Programme Lead for the MSc in Genomic Medicine.

Andrew Collins is head of the Genetic Epidemiology and Bioinformatics Research Group at the University of Southampton and is involved in next generation sequencing-based studies in population genetics and research into the genetic basis of a number of diseases.

2 Introduction

The sequencing of whole genomes using next generation sequencing (NGS) yields vast data-sets which present significant analytical challenges for identification of disease-causal variants. It is known that a subset of human genes contain, or are associated with, rare and/or common variation which have a role in disease processes (the “disease genome”). However, recognition of causal variants amongst many thousands of mostly neutral variants is a huge challenge and a pressing problem. For example, Chong et al. [1] state that the genes underlying ~50% of all Mendelian phenotypes remain unknown and many more Mendelian conditions are still to be described.

1
2
3 Alongside methods for predicting the potential pathogenicity of individual DNA variants at
4 least 20 gene-specific metrics (scores) have been developed in recent years which may help
5 facilitate recognition of disease causing variation. An example of one of these methods is
6 RVIS (residual variation intolerance score) which ranks genes by whether they have more or
7 less common functional genetic variation relative to the genome wide expectation [2]. A
8 candidate pathogenic variant found in a gene classed as intolerant of common functional
9 variation might be worthy of follow-up as a potential causal variant.
10
11

12
13 Understanding the properties of the disease genome and integrating existing gene-specific
14 predictors may help in classifying genes based on their specific features to refine molecular
15 diagnosis. Pathogenicity scores for individual DNA variants are often inconsistent in that
16 different methods can provide conflicting evidence on potential pathogenicity. The degree of
17 redundancy in the genome makes the task of picking out causal variation particularly
18 challenging. We recognize that variant prediction tools alone are currently not conclusive and
19 that evidence at the gene-specific level has the potential to enhance the recognition of variant
20 pathogenicity [3].
21
22

23
24 This systematic review considers the literature related to gene-specific scores and their
25 applicability to improve filtering of genome sequence data. We set out to achieve a
26 satisfactory answer to the following research question: “Can the use of gene-specific metrics
27 facilitate the identification of disease genes in patient genomes?”
28
29

30 Gene-specific metrics are frequently based on properties of genic coding regions. The extent
31 to which they provide information on the tendency of a gene to have associated disease
32 causal variation outside the coding region is limited. Most of the tools analyzed in this
33 review, with a few exceptions, are concerned with genomic coding variation.
34
35

36 Details of the methodology used in this systematic review are given in the Supplementary
37 methods, Supplementary Figures 1 and 2 and Supplementary Table 1.
38
39

40 41 42 43 **3 Findings: Key models**

44
45 Each of the twenty gene-specific approaches identified by the systematic review were classified
46 into one of three groups according to the main focus of each method. We consider below each
47 of the three groups: (i) Essentiality and conservation (ii) Haploinsufficiency (iii) Selection.
48 Supplementary Tables 2-4 give details of the main methods and scores allocated into each
49 category.
50
51

52 53 **3.1 Characteristics of essential and conserved genes.**

54
55 Essential and conserved genes encode proteins which have core biological functions that are
56 essential for an organism’s viability. Genes vary in their degree of essentiality and a number
57
58
59

of different quantitative scores provide approximations to essentiality. These include predictions of the extent to which a gene is tolerant or intolerant of loss of function (LoF) mutations and estimation of the expected rate of *de novo* mutations [14]. Supplementary Table 2 outlines the key approaches in this category. The Residual Variation Intolerance Score (RVIS) ranks genes by probability of carrying more, or less, functional genetic variation than expected highlighting genes intolerant to common functional variation [2]. Genes with positive scores have more common functional variation, while negative scoring genes are less tolerant having reduced associated common functional variation. Genes containing variation involved in monogenic diseases have lower RVIS scores than other genes.

By examining the evolutionary conservation of protein sequences, Rackham et al. developed the Evolutionary Intolerance score (EvoTol) to identify genes which are intolerant to mutation [15] [16]. Because only small areas of a gene may be intolerant, for example protein-coding domains, these sub-regions might be particularly important domains of essentiality [16]. EvoTol allows identification of intolerant protein sub-domains alongside the identification of intolerant genes more generally.

The development of NGS makes possible the identification of newly arising (*de novo*) mutations (DNMs) and their potential roles in rare disease. Recognition of these variants is not without difficulty because of errors in alignment and poorly supported variant calls. Validation by re-sequencing and, in particular, sequencing of additional family members (often the parents of a patient) can help correctly resolve *de novo* variation which might be of disease significance. Such mutations are not considered to play a significant role in the pathogenesis of complex diseases [17]. To accurately estimate the expected rate of *de novo* mutations in a given gene, careful assessment of gene mutability is required. Gene length and local sequence context are essential factors underlying mutation rate differences [17]. Samocha et al. calculated per-gene probabilities of mutation which are correlated with observed counts of rare missense variants in the Exome Sequencing Project (ESP) data set. The Samocha et al. study extends a model which investigated *de novo* mutations in epileptic encephalopathy patients (Epi4K consortium) by considering depth of coverage (i.e., how many sequence reads were present on average per base) and the regional divergence in genes between humans and Macaques. Significant numbers of genes with missense variant deficits were observed, compared to expectation from predicted mutation rates, suggesting strong evolutionary constraint removing variants by negative selection [17] [18]. The Samocha et al. model utilizes exome sequence data to evaluate the DNM rate by gene set and on a single gene basis [17], this score is referred to as *de novo* excess (DNE). The metric is predictive of selective constraint in the human genome and identifies 1,003 constrained genes known to cause severe human disease [17]. It was found that constrained genes contain higher *de novo* LoF mutation rate than expected by chance [17].

The LoFtool measures the ratio of LoF mutations to synonymous mutations for every gene. The performance of the LoFtool, compared to RVIS, DNE Z-score, and EvoTol, suggests enhanced prediction of *de novo* haploinsufficient disease-causing genes. The LoFtool represents values as intolerance percentiles: genes that are intolerant to LoF variation have low LoFtool percentiles [15]. The four measures of genic intolerance outlined so far were included by Bartha et al. who described them as essentiality scores [19].

1
2
3
4 In early 2016, using data from 1000 Genomes Project, Aggarwala et al. proposed the
5 Substitution Intolerance Score (SIS) as a gene-level measurement of essentiality. Genes with
6 high SIS scores are functionally constrained, while genes which score low are tolerant of
7 functional changes in the protein which might arise through mutations in the DNA sequence
8 [20].
9

10
11 Another scoring system by Gussow et al. evaluates intolerance in genic sub-regions
12 proposing that more conserved regions within a gene are expected to contain more variants
13 which are pathogenic [21]. Genes are divided into sub-regions and tiered by intolerance to
14 functional variation. This 'subRVIS' score ranks regions using RVIS but with the addition of
15 information on conservation. Regions intolerant to functional variation are scored low by the
16 subRVIS scoring system. The method utilizes the GERP++ [22] score to evaluate
17 evolutionary constraint for bases in each sub-region [21].
18
19

20
21 The Loss Intolerance probability (pLI) score quantifies the likelihood that a gene is
22 intolerant to a mutation which produces LoF in the protein product [23]. The score is derived
23 using the Exome Aggregation Consortium (ExAC) database which is an extensive catalogue
24 of human genetic diversity. This catalogue identifies one variant every eight bases on average
25 in the exome providing a powerful filter for analysis of candidate deleterious variants in
26 severe Mendelian diseases [23]. Lek et al. proposed that genes with high pLI score ($pLI \geq$
27 0.9) are most intolerant of LoF variation. Genes in this category are the most evolutionarily
28 constrained. The least constrained genes (LoF tolerant) have low pLI scores ($pLI < 0.1$) and
29 typically contribute to the least constrained biological pathways, such as sensory perception,
30 where high haplotype diversity is potentially advantageous [23].
31
32

33
34 It is challenging to assess the relationship between the DNM rate and genes involved in
35 disease. In 2017, Jiang et al. utilized available DNM data to correct for the background
36 mutation rate seen as one of the main limitations of the Samocha et al. [17] model. The
37 problem arises because by sequencing more individuals, more DNMs are inevitably observed
38 in the same gene by chance. Therefore, in a given disease, if a *de novo* mutation is related to
39 pathogenesis, disease-genes might be expected to contain more DNMs than predicted from
40 background rates. This work includes the development of a database which describes the
41 background DNM rate (DNMR), acquired from population variation data [24].
42
43
44
45

46 **3.2 Characteristics of Haploinsufficient genes**

47

48 Haploinsufficiency (HI) occurs whenever there is a missing or damaged copy of a gene
49 leaving a single copy which is insufficient to maintain normal function [3].
50 Haploinsufficiency is mostly caused by LoF mutations and results in dominant diseases.
51 Recognition and prediction of genes which are haploinsufficient can facilitate the filtering of
52 disease genome data wherever the phenotype is likely to have arisen through reduced levels
53 of gene product.
54
55
56
57
58
59
60

1
2
3 In 2010, Haung et al. proposed a deletion-based HI score by identifying differences between
4 HI and haplosufficient (HS) genes, aiming to better distinguish pathogenic from benign
5 deletions which helps in variant prioritization [3]. The analysis develops a logarithm-of-odds
6 (LOD) score to estimate the probability of a deletion causing a HI phenotype. A high LOD
7 score suggests deletions are likely to be deleterious through HI and therefore potential
8 candidates for causing dominant traits. The score assumes there are no statistical interactions
9 between the genes [3]. Previously, and to try to assess the pathogenicity of a deletion,
10 clinicians considered the length of a deletion or the number of genes deleted. The Haung et
11 al. score provides a rational basis to classify pathogenic deletions by comparing deletions
12 seen in patients with deletions in controls and calculating the fraction of controls with a
13 deletion at least as deleterious as that seen in the patient [3].
14
15
16
17

18 Distinguishing false-positive disease variants from the genuinely causal variants is crucial for
19 accurate molecular diagnoses. MacArthur et al. developed the REcessive (REC) score for
20 distinguishing genes involved in recessive diseases from genes which are LoF- variation
21 tolerant [25]. A “healthy” genome might contain 100 true LoF variants, the majority in a
22 heterozygous state. Evidence suggests that the average human carries five recessive lethal
23 alleles in single copy in their genome. Consequently, the majority of LoF variants are
24 considered common variants. However, these variants might still have a phenotypic effect
25 [25]. MacArthur et al. demonstrated differences in functional and evolutionary features
26 between recessive disease and LoF-tolerant genes, allowing for the development of a
27 predictive model to predict recessive disease variants [25].
28
29
30
31

32 Khurana et al. developed the “gene position in NETworks” (NET) indispensability score to
33 investigate relationships between degree of network centrality of a gene and selection within
34 biological networks [26]. They consider a range of biological networks relating to
35 phosphorylation, signaling, protein-protein interaction and regulatory and genetic networks.
36 Genes which are highly connected to many biological networks are the most functionally
37 significant, therefore, mutations in those genes might have serious consequences[26].
38 However, genes connected to metabolic networks were found to have an excess of
39 duplicated copies through more paralogs with LoF mutations[26]. This score was included as
40 a predictor of haploinsufficient genes in the Hsu et al. study [5]
41
42
43
44

45 Ge et al. consider gene-specific pathogenicity using the ratio of non-synonymous to
46 synonymous substitution rates (dN/dS) for X-chromosome genes [27]. Genes with unusually
47 low ratios suggest intolerance to non-synonymous variation, indicating they may be
48 susceptible to disease-related variation. The authors found correlation between genomic
49 regions depleted for missense variation with disease-causal variants [27].
50
51

52 Steinberg et al. proposed that study biases existing in many biological networks might affect
53 the ability of previous HI prediction scores to recognize the genuinely haploinsufficient
54 genes. For that reason they constructed a new, unbiased, HI score, the Genome-wide
55 HaploInsufficiency Score (GHIS) which replaces biological networks with co-expression
56
57
58
59
60

1
2
3 networks [28] [29]. They compared their model with the three pre-existing methods (i.e., HI
4 [3], NET [26] and RVIS [2]) and demonstrated that GHIS provides a score for many genes
5 not scored by other methods [28] with enhanced performance at classifying less well
6 studied genes [28].
7

8
9 Scores have been developed to recognize Mendelian genes with different modes of
10 inheritance. Hsu et al. considered Mendelian disease gene characteristics according to their
11 mode of inheritance. Haploinsufficiency is an essential characteristic of Mendelian disease
12 genes with an autosomal dominant (AD) mode of inheritance and sensitivity to *de novo*
13 mutations was recognized for this group of genes [5]. In contrast disease genes with
14 autosomal recessive (AR) modes of inheritance tend to have more non-synonymous variants
15 and regulatory transcript isoforms [5]. However, the X-linked (XL) pattern of inheritance is
16 associated with fewer non-synonymous and synonymous variants [5]. Based on these
17 findings they create a new approach to prioritize Mendelian disease genes based on their
18 mode of inheritance (AD, AR, and XL) termed Inheritance-mode Specific Pathogenicity
19 Prioritization (ISPP) [5]. This score integrates pre-existing gene-specific prediction methods
20 namely: HI [3], REC [25], RVIS [2], NET [26], DNE [17] and GDI [30] along with numerous
21 genetic properties including global expression from RNA-Seq data, DNA replication time
22 and the noncoding (intronic region) mutation rate [5].
23
24
25
26

27 Because the human genome contains an abundance of non-deleterious heterozygous variants,
28 the identification of dominant mutations for monogenic disorders is challenging. Quinodoz et
29 al. created DOMINO a method using machine learning to identify whether a given gene is
30 liable to carry dominant changes [31].
31
32

33 Inevitably, well-studied genes are over-represented in most biological networks used to
34 create scores that predict HI compared to less-studied genes, hence most biological networks
35 are affected by study bias. Therefore the creation of unbiased HI score becomes particularly
36 important [29]. Recently, Shihab et al. produced an integrated machine learning approach
37 called (HIPred) merging functional annotations with genomic and evolutionary features to
38 predict HI genes without study bias using data from NIH Roadmap Epigenomics [32] and the
39 ENCODE [33] project. The performance of this approach is considered to exceed the pre-
40 existing HI predictors [29]. Supplementary Table 3 outlines the key approaches in this
41 category.
42
43
44
45
46

47 **3.3 Characteristics of genes under selection.**

48

49 Genetic variants may be subject to positive selection whereby, if they are advantageous,
50 they may increase in frequency. Negative selection, in contrast, acts to remove deleterious
51 alleles. Scores which quantify the intensity of negative selection acting on genes provide
52 insights into which genes are more likely to have variation which may have damaging
53 consequences. The pattern is complex because some essential genes are not known to have
54
55
56
57
58
59
60

1
2
3 any associated disease variation and are perhaps subject to negative selection at
4 particularly high intensity [34].

5 Bustamante et al. calculate the extent and directionality of Selection operating on a given
6 gene, this score referred to here as “Sel”. They first compared fixed sequence differences,
7 both synonymous and non-synonymous, between humans in the sample and Chimpanzees
8 over 11.81 Mb region of aligned coding DNA. The ratio of non-synonymous to
9 synonymous differences (divergence) was 23.76%. In contrast the ratio of non-
10 synonymous to synonymous polymorphisms in the human subjects was 38.42%. This
11 shows a significant excess of amino acid variation, relative to divergence, consistent with
12 previous work stating that much amino acid variation in the human genome is slightly to
13 moderately damaging [35].
14
15
16

17 Eilertson et al. create a model to identify genes under natural selection with a non-parametric
18 approach (with no assumption of a specific population genetic model) which is robust to
19 demography [36]. This approach, called Selection Inference using Poisson Random Effects
20 (SnIPRE), utilizes polymorphism and divergence data from synonymous and non-
21 synonymous sites within genes.
22
23
24

25 The Gene-level Integrated Metric of negative Selection (GIMS) was created by combining
26 two meta-analyses into a single meta-analysis. The first meta-analysis combines comparative
27 genomic metrics (GERP++) [22] and functional genomic metrics (Poly-phen2) [37], and the
28 second meta-analysis combines mutation rates (as SNPs/kb) and allele frequencies (as
29 percentage rare) from the 1000 Genomes Project. Meta-analysis was achieved by combining
30 those metrics into GIMS scores for 20,079 genes [38]. Because the majority of genes are
31 under purifying selection, the aim was to quantify the degree of negative selection applied to
32 genes. Conservation and functional scores were initially combined as ‘functional genomic
33 metrics’ integrated with mutation rates and fraction of rare variants as ‘population genetic
34 metrics’. The GIMS score combines these two metrics and provides a unified score per-
35 gene. GIMS gives a probability distribution across the entire genome in quantiles. Genes
36 under negative selection are scored low by GIMS [38].
37
38
39
40

41 The Gene Damage Index (GDI) is a gene-specific score which predicts the liability of a
42 human protein-coding gene to contain disease-causing mutations considering the influences
43 of selection and genetic drift. In GDI, Combined Annotation Dependent Depletion (CADD)
44 [39] scores are used as the variant-level damage prediction method because this method is
45 efficient at distinguishing between benign and deleterious variants and is strongly dependent
46 on evolutionary conservation [30]. Moreover, CADD scores can assess most types of variants
47 while other methods, like Poly-Phen-2 [37] and SIFT [40], can only predict missense
48 variants. To construct the GDI score the cumulative predicted damage in exonic regions of
49 the gene is calculated using the CADD score for each allele compared to the expected score
50 for variants with similar allele frequencies. The homogenized Phred I-score is calculated for
51 each metric to indicate the ranking of the targeted gene relative to all other genes. A low
52 Phred score: indicates a human gene with a low GDI and high Phred score indicates a gene
53 susceptible to contain damaging variation. Genes with high GDI tend to be under less
54
55
56
57
58
59
60

intense purifying selective pressure. A low GDI score is associated with highly conserved genes (including genes enriched for ribosome, chemokine signaling proteasome and spliceosome functions) reflecting essentiality. Such genes tend to be under stronger purifying selection than the median selective pressure acting on human genes [30]. Supplementary Table 4 outlines the key approaches in this category.

4 Discussion

Considering approaches which score genes according to essentiality and conservation the DNE score offers some advantages. The main limitation of DNE is its validity only for interpretation of *de novo* mutations [5] but it considers more variables related to mutation rate going beyond sequence context compared to other methods like RVIS and Sel. These additional variables include consideration of sequence depth of coverage and regional divergence in genes between humans and Macaques independently, which improve the predictive value of this model [17]. The DNE score has been compared to the RVIS and negative selection score Sel. The comparison showed that DNE and RVIS were equally effective emphasizing the benefits predicted from combining the two scores [17].

The strength of Samocha et al. model is enhanced by incorporation of the depth of coverage (i.e., how many sequence reads were present on average per base) and the regional divergence in genes between humans and Macaques independently. These strengths play a significant role in the improvement of their predictive model. The number of rare synonymous variants in the Exome Sequencing Project (ESP) which comprises a relatively small sample of 6700 exomes [41] is shown to be highly correlated with the probability of a synonymous mutation determined by their model. As rare variant allele frequencies are impacted by sample size evaluation in larger databases such as ExAC would be of interest [41].

EvoTol was compared to the RVIS and the DNE scores and shown to have increased performance at classifying intolerant genes compared to RVIS. EvoTol was shown to be highly sensitive and more powerful to characterize genes with high pathogenicity [16]. Although there was no significant correlation between RVIS and EvoTol, the application of the two scores simultaneously will likely be advantageous [16].

Considering approaches for scoring genes for potential roles in haploinsufficiency phenotypes the HIPred approach has been evaluated against five predictors (HI Score, NET, RVIS, EvoTol and GHIS, Supplementary Tables 2 and 3). HIPred was found to outperform all in predicting HI genes [29]. Using different perspectives across the 26 disease-associated gene lists, Hsu et al. estimates the power of several methods that predict gene pathogenicity showing a substantial positive correlation between HI and REC (correlation $r=0.77$) while the six scores have a moderate relationship with each other ($r=0.46$) [5]. Among these gene scores (DNE, GDI, HI, NET, RVIS, and REC) the best predictor of disease-predisposing genes was the REC score [5]. The performance of the ISPP score was significantly superior

1
2
3 for prioritizing AR and X-linked disease-associated genes [5]. The REC score is effective at
4 predicting disease-associated genes generally but less successful in discriminating recessive
5 and dominant disease genes [5].
6

7
8 DNE measures the rate of per-gene *de novo* mutation while RVIS ranks human genes based
9 on the strength and consistency of the purifying selection acting against functional variation.
10 Analysis has shown that GDI and RVIS capture unique sets of reciprocal information from
11 population genetic data [30]. In essence, RVIS reflects selective pressure while DNE is based
12 on *de novo* mutation rate estimates; both methods do not quantitatively estimate the
13 mutational load for a gene in a healthy human population. For this reason, these methods are
14 not optimal for filtering genes with high mutation rates and many residual false positives
15 might be expected. GDI has proved to be the most efficient approach for filtering out false
16 positive variants in genes known to contain damaging variation [30].
17
18

19
20 The Ge et al. X-linked scoring system is not limited by previous gene annotation and the
21 dN/dS ratio can be calculated for any protein-coding gene. This score applies to all X-
22 chromosome protein-coding genes and therefore can assess genes for multiple disease
23 phenotypes [27]. Because the intra-human dN/dS ratio is not specific to the X-chromosome
24 the analysis of more genomic data using dN/dS ratio is recommended for future studies to
25 identify genes which may have disease variation [27].
26
27

28
29 The effort to improve the predictive ability of variant-level scores now includes combination
30 of evidence from multiple pathogenicity scores and other data. An example is the
31 “Mendelian Clinically Applicable Pathogenicity” (M-CAP) score [42] which uses machine
32 learning classification based on existing pathogenicity scores and measures of
33 evolutionary conservation. Such a combinatorial approach might usefully integrate
34 evidence from both variant-level and gene-level metrics to improve predictive abilities
35 overall [42].
36
37

38
39 This work aims to bring together the growing evidence that gene properties, alongside variant
40 scoring systems, can play an important role in filtering disease sequence data. As healthy
41 individuals can have genetic variants that lead to disruption of protein-coding genes (with no
42 clinical phenotype) [25,28,29,43], challenges remain to distinguish which loss of function
43 variants are associated with disease phenotypes from those that do not cause any functional
44 disturbance [28]. Data from the 1000 Genomes Project show that on average a healthy person
45 might carry 250-300 LoF SNVs (1000 Genomes Project Consortium et al., 2010; The 1000
46 Genomes Project Consortium, 2012) [5].
47
48

49
50 The ACMG guidelines consider *in silico* predictions of whether a variant is involved in
51 disease, but without specifying which or how many variant interpretation algorithms to use.
52 These data can be used only as ‘supporting’ evidence for variant interpretation. There are
53 difficulties with respect to validation of these methods and there is a relatively high error rate
54 with many pathogenic variants assessed as benign by some methods and many benign
55 variants assessed as pathogenic [44]. The guidelines do not currently consider gene-specific
56
57
58
59
60

1
2
3 metrics which are the subject of this review but presumably could similarly constitute
4 supporting evidence given alongside stronger independent evidence suggesting role or lack of
5 role in disease. Ultimately, functional validation is optimal although is frequently not
6 timely, practical or reimbursable [44,45].
7

8
9 The understanding of human genomes is advanced through the accumulation of sequence
10 data in publically available databases. The ExAC resource provides a potent filter to aid
11 recognition of pathogenic variants in severe Mendelian diseases. Using ExAC for filtering to
12 remove false positive, but plausibly pathogenic, variants decreases the number of candidate
13 protein-altering variants by 7-fold compared to the smaller Exome Sequencing Project
14 database (ESP) which has fewer exome sequences [23].
15
16

17
18 Coupled with the previous evidence, another study suggests that the missense Z score which
19 represents genes rather than variants adds more information than variant-specific Poly-phen2
20 and CADD classifications signifying that gene-level scores of constraints provide additional
21 information for evaluating pathogenicity [23]. Furthermore, Huang et al. contend that variant
22 level scores (e.g., SIFT [40] and poly-phen 2 [37]) are limited by lacking the capability to
23 determine , from cross-species alignments, whether negative selection at a given site is
24 acting in a recessive, additive or dominant mode [3].
25
26

27
28 The work proposed by Gussow et al. was based on dividing the genes into sub-regions to
29 identify exactly where the pathogenic mutations are likely to present [21]. This study
30 identified an important question: is the whole gene the correct unit by which to judge patterns
31 of intolerance? Future analyses may consider refinements to gene-specific scores which
32 consider within-gene regional patterns of intolerance in more detail.
33

34
35 Another controversial issue is the difficulty in interpretation of benign LoF variants for which
36 the nomenclature is still not unified. It is important to realize that there are overlaps in the
37 interpretation of LoF variants in healthy people. In the literature, all the following categories
38 are represent LoF variants in healthy individuals: true variants that do not seriously disrupt
39 gene function, benign LoF variation in redundant genes, non-deleterious or less-deleterious
40 variants that have an impact on risk of phenotype or disease [25].
41
42

43
44 Because each genic scoring approach considers only a specific property of genetic
45 architecture, each individual score has limitations. For example: (i) the REC score does not
46 consider dominant disease-predisposing genes (ii) Non-CNV (Non-Copy Number Variation)
47 genetic variants were not included in HI prediction score. (iii) the NET score lacks the
48 systematic comparison of different known disease-associated genes (iv) the RVIS score does
49 not consider variations in allele frequencies across different populations (v) the DNE score
50 has limited applicability for testing *de novo* mutations. (vi) the GDI score only considers
51 mutation profiles [5]. Furthermore, a major limitation of the GHIS score is that the genetic
52 background in individuals is not considered, which is an important issue since genetic
53 variants do not act in isolation and disturbance of individual genes within a single biological
54 pathway might affect the risk of a disease [28]. Accordingly, this analysis which provides a
55
56
57
58
59
60

comprehensive review of each prediction scheme, may help establish new routes for prioritizing disease-causal variants.

Many advances have been developed to assess whether a gene is tolerant or intolerant to common functional variation. Initially, scores were developed per gene then studies were published showing that dividing the gene into sub-regions might help in allocating the mutation accurately. At that time all scores that measure genic intolerance required disease knowledge, this limitation was addressed by developing a tool with no prior disease knowledge required, an essential step to better predict genic intolerance.

Reviewed here are a range of well-studied gene-specific predictors with various independent genetic properties. It is hoped that recognizing some of the limitations of each score and perhaps combining evidence from both variant-specific scores and gene-wise evidence might enable better prediction since there is currently no single method that is reliably predictive of gene pathogenicity. Therefore this hopefully might help to overcome one of the main challenges of 100,000 genome project which is variant annotation to prioritize important variants from harmless neutral variants. This review is intended to highlight existing work to identify and explain different gene-specific pathogenicity predictors, while pointing to the gaps in disease-gene prioritization and annotation issues to facilitate new scores and better prioritization of disease-causal genes.

Key points

1. A wide range of well-established models exist that prioritize genes based on their associated disease variation potential.
2. Integration of these strategies to represent individual genes could have a significant impact on our understanding of genic properties and the recognition of disease-related functional variation.
3. Evaluation and comparison of these individual scores and the development of integrated models to enhance NGS filtering strategies in disease genomes is a fertile area for future studies.

Funding

DA is funded by the Saudi Arabia Cultural Bureau, London, UK.

DB is funded through a NIHR Research Professorship.

References

- [1] J. X. Chong *et al.*, "The Genetic Basis of Mendelian Phenotypes: Discoveries, Challenges, and Opportunities," *Am. J. Hum. Genet.*, vol. 97, no. 2, pp. 199–215, 2015.
- [2] S. Petrovski, Q. Wang, E. L. Heinzen, A. S. Allen, and D. B. Goldstein, "Genic

- Intolerance to Functional Variation and the Interpretation of Personal Genomes,” *PLoS Genet.*, vol. 9, no. 8, 2013.
- [3] N. Huang, I. Lee, E. M. Marcotte, and M. E. Hurles, “Characterising and predicting haploinsufficiency in the human genome,” *PLoS Genet.*, vol. 6, no. 10, pp. 1–11, 2010.
- [4] X. Huang, J. Lin, and D. Demner-Fushman, “Evaluation of PICO as a knowledge representation for clinical questions,” *AMIA Annu. Symp. Proc.*, pp. 359–63, 2006.
- [5] I. T. Newsweekly, “Mendelian Disease ; Findings from University of Hong Kong Yields New Data on Mendelian Disease [(ISPP) for human protein coding genes],” vol. 32, no. 20, pp. 2016–2018, 2016.
- [6] B. Kitchenham and S. Charters, “Guidelines for performing Systematic Literature Reviews in Software Engineering,” *Engineering*, vol. 2, p. 1051, 2007.
- [7] S. U. Khan, M. Niazi, and R. Ahmad, “Barriers in the selection of offshore software development outsourcing vendors: An exploratory study using a systematic literature review,” *Inf. Softw. Technol.*, vol. 53, no. 7, pp. 693–706, 2011.
- [8] C. Jalal, Samireh and Wohlin, “Systematic Literature Studies: Database Searches vs. Backward Snowballing Samireh.” .
- [9] D. Badampudi, “Experiences from using snowballing and database searches in systematic literature studies.” .
- [10] J. Gehanno, L. Rollin, and S. Darmoni, “Is the coverage of google scholar enough to be used alone for systematic reviews,” no. December 2009, pp. 0–4, 2013.
- [11] S. Becker, A. Bryman, and H. (Thomas H. Ferguson, *Understanding research for social policy and practice : themes, methods and approaches*. Policy, 2012.
- [12] G. Craswell and M. Poore, *Writing for academic success*. SAGE, 2012.
- [13] C. Thermes, “Ten years of next-generation sequencing technology,” *Trends Genet.*, vol. 30, no. 9, pp. 418–426, 2014.
- [14] R. J. Pengelly, A. Vergara-Lope, D. Alyousfi, M. R. Jabalameli, and A. Collins, “Understanding the disease genome: gene essentiality and the interplay of selection, recombination and mutation,” *Brief. Bioinform.*, no. June, pp. 1–7, 2017.
- [15] J. Fadista, N. Oskolkov, O. Hansson, and L. Groop, “LoFtool: A gene intolerance score based on loss-of-function variants in 60 706 individuals,” *Bioinformatics*, vol. 33, no. 4, pp. 471–474, 2017.
- [16] O. J. L. Rackham, H. A. Shihab, M. R. Johnson, and E. Petretto, “EvoTol : a protein-sequence based evolutionary intolerance framework for disease-gene prioritization,” vol. 43, no. 5, 2014.
- [17] K. E. Samocha *et al.*, “A framework for the interpretation of de novo mutation in human disease,” vol. 46, no. 9, pp. 944–950, 2015.
- [18] A. S. Allen *et al.*, “De novo mutations in epileptic encephalopathies,” *Nature*, vol. 501, no. 7466, pp. 217–221, 2013.
- [19] A. Bartha, István & di Iulio, Julia & Venter, J & Telenti, “Human gene essentiality. Nature Reviews Genetics.” 2017.
- [20] V. Aggarwala and B. F. Voight, “An expanded sequence context model broadly explains variability in polymorphism levels across the human genome,” *Nat. Genet.*, vol. 48, no. 4, pp. 349–355, 2016.
- [21] A. B. Gussow, S. Petrovski, Q. Wang, A. S. Allen, and D. B. Goldstein, “The intolerance to functional genetic variation of protein domains predicts the localization of pathogenic mutations within genes,” *Genome Biol.*, vol. 17, no. 1, pp. 1–11, 2016.
- [22] E. V. Davydov, D. L. Goode, M. Sirota, G. M. Cooper, A. Sidow, and S. Batzoglou, “Identifying a high fraction of the human genome to be under selective constraint using GERP++,” *PLoS Comput. Biol.*, vol. 6, no. 12, 2010.
- [23] Lek, “Analysis of protein-coding genetic variation in 60,706 humans,” *Nature*, vol.

- 536, no. 7616, pp. 285–291, 2017.
- [24] Y. Jiang *et al.*, “MirDNMR: A gene-centered database of background de novo mutation rates in human,” *Nucleic Acids Res.*, vol. 45, no. D1, pp. D796–D803, 2017.
- [25] D. MacArthur, S. Balasubramanian, and A. Frankish, “A Systematic Survey of Loss-of-Function Variants in Human Protein-Coding Genes,” *Science (80-.)*, vol. 335, no. 6070, pp. 1–14, 2012.
- [26] E. Khurana, Y. Fu, J. Chen, and M. Gerstein, “Interpretation of Genomic Variants Using a Unified Biological Network Approach,” *PLoS Comput. Biol.*, vol. 9, no. 3, 2013.
- [27] X. Ge, P. Y. Kwok, and J. T. C. Shieh, “Prioritizing genes for X-linked diseases using population exome data,” *Hum. Mol. Genet.*, vol. 24, no. 3, pp. 599–608, 2015.
- [28] J. Steinberg, F. Honti, S. Meader, and C. Webber, “Haploinsufficiency predictions without study bias,” *Nucleic Acids Res.*, vol. 43, no. 15, pp. 1–9, 2015.
- [29] H. A. Shihab, M. F. Rogers, C. Campbell, and T. R. Gaunt, “HIPred: An integrative approach to predicting haploinsufficient genes,” *Bioinformatics*, vol. 33, no. 12, pp. 1751–1757, 2017.
- [30] Y. Itan *et al.*, “The human gene damage index as a gene-level approach to prioritizing exome variants,” *Proc. Natl. Acad. Sci.*, vol. 112, no. 44, pp. 13615–13620, 2015.
- [31] M. Quinodoz *et al.*, “REPORT DOMINO : Using Machine Learning to Predict Genes Associated with Dominant Disorders,” *Am. J. Hum. Genet.*, vol. 101, no. 4, pp. 623–629, 2017.
- [32] Roadmap Epigenomics Consortium *et al.*, “Integrative analysis of 111 reference human epigenomes,” *Nature*, vol. 518, no. 7539, pp. 317–329, 2015.
- [33] Encode Consortium, N. Carolina, and C. Hill, “For Junk DNA,” *Nature*, vol. 489, no. 7414, pp. 57–74, 2013.
- [34] N. Spataro, J. A. Rodríguez, A. Navarro, and E. Bosch, “Properties of human disease genes and the role of genes linked to Mendelian disorders in complex disease aetiology,” *Hum. Mol. Genet.*, vol. 26, no. 3, pp. 489–500, 2017.
- [35] C. D. Bustamante *et al.*, “Natural selection on protein-coding genes in the human genome,” *Nature*, vol. 437, no. 7062, pp. 1153–1157, 2005.
- [36] K. E. Eilertson, J. G. Booth, and C. D. Bustamante, “SnIPRE: Selection Inference Using a Poisson Random Effects Model,” *PLoS Comput. Biol.*, vol. 8, no. 12, 2012.
- [37] I. A. Adzhubei1, “A method and server for predicting damaging missense mutations,” *October*, vol. 7, no. 4, pp. 248–249, 2010.
- [38] M. G. Sampson, C. E. Gillies, W. Ju, M. Kretzler, and H. M. Kang, “Gene-level integrated metric of negative selection (GIMS) prioritizes candidate genes for nephrotic syndrome,” *PLoS One*, vol. 8, no. 11, pp. 1–9, 2013.
- [39] M. Kircher, “A general framework for estimating the relative pathogenicity of human genetic variants,” *Nat. g.*, vol. 46, no. 3, pp. 310–315, 2014.
- [40] P. Kumar, S. Henikoff, and P. C. Ng, “Predicting the effects of coding non-synonymous variants on protein function using the SIFT algorithm,” *Nat. Protoc.*, vol. 4, no. 7, pp. 1073–1082, 2009.
- [41] P. L. Auer *et al.*, “Guidelines for Large-Scale Sequence-Based Complex Trait Association Studies: Lessons Learned from the NHLBI Exome Sequencing Project,” *Am. J. Hum. Genet.*, vol. 99, no. 4, pp. 791–801, 2016.
- [42] K. A. Jagadeesh *et al.*, “M-CAP eliminates a majority of variants of uncertain significance in clinical exomes at high sensitivity,” *Nat. Genet.*, vol. 48, no. 12, pp. 1581–1586, 2016.
- [43] P. C. Ng *et al.*, “Genetic variation in an individual human exome,” *PLoS Genet.*, vol. 4, no. 8, 2008.

- 1
2
3 [44] L. J. H. Bean and M. R. Hegde, "Clinical implications and considerations for
4 evaluation of in silico algorithms for use with ACMG/AMP clinical variant
5 interpretation guidelines," *Genome Med.*, vol. 9, no. 1, pp. 9–11, 2017.
6 [45] K. D. Laboratories *et al.*, "HHS Public Access," *CA Cancer J Clin*, vol. 17, no. 5, pp.
7 405–424, 2015.
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review