1

# Characterising *RAG1* and *RAG2* with predictive genomics

Dylan Lawless, MSc[a,*], Jolan E. Walter, MD, PhD[b,c], Rashida Anwar, PhD[a], Sinisa Savic, MD, PhD[d,e,*]

[a] *Leeds Institute of Biomedical and Clinical Sciences, University of Leeds, Wellcome Trust Brenner Building, St James's University Hospital, Beckett Street, Leeds, UK.*
[b] *University of South Florida and Johns Hopkins All Children's Hospital, Saint Petersburg, FL, USA.*
[c] *Division of Allergy Immunology, Massachusetts General Hospital for Children, Boston, MA, USA.*
[d] *Department of Clinical Immunology and Allergy, St James's University Hospital, Beckett Street, Leeds, UK.*
[e] *National Institute for Health Research Leeds Musculoskeletal Biomedical Research Centre and Leeds Institute of Rheumatic and Musculoskeletal Medicine, Wellcome Trust Brenner Building, St James's University Hospital, Beckett Street, Leeds, UK.*

---

## Abstract

While widespread genome sequencing ushers in a new era of preventive medicine, the tools for predictive genomics are still lacking. The greatest hurdle in diagnosis of rare disease is validation for variants of unknown significance. RAG deficiency presents at an early age with a distinct phenotype of combined immunodeficiency with granuloma and/or autoimmunity. Allele frequency of a SNV in the general population is an indicator of the functional or structural importance of a particular amino acid residue. However, rare diseases are often attributable to variants in genes which are highly conserved. Mutation of a conserved residue does not confirm pathogenicity and functional validation must be confirmed to correctly identify a monogenic disorders such as RAG deficiency. We present protein variants in RAG1 and RAG2 which are most likely to be seen clinically as disease-causing. Our method of mutation rate residue fre-

---

*Addresses for correspondence
Email addresses:* `D.Lawless@leeds.ac.uk` (Dylan Lawless, MSc), `S.Savic@leeds.ac.uk` (Sinisa Savic, MD, PhD)

quency builds a map of most probable mutations allowing pre-emptive functional analysis. We compare the accuracy of our predicted probabilities to previously established functional measurements.

## Acknowledgements

## Key words

RAG1, RAG2, genomics.

## Abbreviations

CADD (combined annotation dependent depletion), GWAS (genome-wide association studies), $M_r$ (mutation rate), MFR (mutation rate residue frequency), (pLI) probability of being loss-of-function intolerant, $R_f$ (residue frequency), RAG1 (recombination activating gene 1).

**Introduction**

Costs associated with genomic investigations continue to reduce [1] while the richness of data generated increases. Globally, the adoption of wide scale genome sequencing implies that all new-born infants may receive screening for pathogenic genetic mutation in an asymptomatic stage, pre-emptively [2]. The one dimensionality of individual genomes is now being expanded by the possibility of massive parallel sequencing for somatic variant analysis and by single-cell or lineage-specific genotyping; culminating in a genotype spectrum. In whole blood, virtually every nucleotide position may be mutated across $10^5$ cells [3]. Mapping one's genotype across multiple cell types and at several periods during a person's life may soon be feasible [4]. Such genotype snapshots might allow for prediction and tracking of somatic, epigenetic, and transcriptomic profiling.

The predictive value of the screening highly depends on the computation tools used for data analysis and its correlation with functional assays or prior clinical experience. Interpretation of that data is especially challenging for variants of unknown significance. There is a need for predictive genomic modelling with aims to provide a reliable guidance for therapeutic intervention for patients harbouring genetic defects for life threatening disease before the illness becomes clinically significant. Although, most genomic investigations currently are not predictive for clinical outcome. The study of predictive genomics is exemplified by consideration of gene essentiality, accomplished by observing intolerance to loss-of-function variants. Several gene essentiality scoring meth-

3

49 ods are available for both the coding and non-coding genome [5]. Approxi-

50 mately 3,000 human genes cannot tolerate the loss of one allele [5]. The great-

51 est hurdle in monogenic disease is the interpretation of variants of unknown

52 significance while functional validation is a major time and cost investment

53 for laboratories investigating rare disease. Severe, life-threatening immune dis-

54 eases are caused by genetic variations in almost 300 genes [6, 7] however, only

55 a small percentage of disease causing variants have been characterised with

56 functional studies. Our investigation aims to apply predictive genomics as a

57 tool to identify pathogenic genetic variants that are most likely to be seen in

58 patient cohorts.

59     We present the first application of our novel approach of predictive genomics

60 using Recombination activating gene 1 (RAG1) and RAG2 deficiency as a model

61 for a rare primary immunodeficiency caused by autosomal recessive variants.

62 *RAG1* and *RAG2* encode lymphoid-specific proteins that are essential for V(D)J

63 recombination. This genetic recombination mechanism is essential for a ro-

64 bust immune response by diversification the T and B cell repertoire in the thy-

65 mus and bone marrow, respectively [8, 9]. RAG deficiency is mesured by in

66 vitro quantification of recombination activity. Hypomorphic RAG1 and RAG2

67 mutations with residual V(D)J recombination activity (in average 5-30%) re-

68 sult in a distinct phenotype of combined immunodeficiency with granuloma

69 and/or autoimmunity (CID-G/A) [2, 10, 11]. *RAG1* and *RAG2* are highly con-

70 served genes but disease is only reported with autosomal recessive inheritance.

71 Only 44% of amino acids in RAG1 and RAG2 are reported as mutated on Gno-

72 mAD and functional validation of clinically relevant variants is difficult. Pre-

4

73 emptive selection of residues for functional validation is a major challenge; a

74 selection based on low allele frequency alone is infeasible. A shortened time be-

75 tween genetic analysis and diagnosis means that treatments may be delivered

76 earlier. With such tools, patients with RAG deficiency may receive hematopoi-

77 etic stem cell transplant [12] or be provided mechanism-based treatment [13].

78 GnomAD was queried to identify conserved residues using a Boolean score

79 $C$ (0 or 1, although allele frequency can be substituted). The gene-specific mu-

80 tation rate $M_r$ of each residue was calculated from allele frequencies. The gene-

81 specific residue frequency $R_f$ was also calculated and together the values calcu-

82 late the most probable disease-causing variants which have not yet been iden-

83 tified in patients. We term the resulting score a mutation rate residue frequency

84 (MRF); where $MRF = C \times M_r \times R_f$. For visualisation, a noise reduction method

85 was also applied where the average MRF per 1% interval is displayed with a

86 cut-off threshold at the 75th percentile.

## Results

87 

*RAG1 and RAG2 conservation and mutation rate residue frequency.*

89 Fig 1 presents the most probable unidentified disease-causing variants in RAG1/2.

90 Phenotypic, epigenetic, or other such weighting data may also be applied to

91 this model. Variants with a low MRF may still be damaging but resources for

92 functional validation are best spent on gene regions with high MRF. Clusters

93 of conserved residues are shown in Fig 1(i) however; these clusters do not pre-

94 dict the likelihood of clinical presentation. Raw MRF scores are presented in

95 Fig 1(ii). A histogram illustrates the MRF without Boolean scoring applied and

5

96   Fig 1(iii) presents a clearer visualisation. Table S1 provides all MRF scores for

97   both proteins as well as raw data used for calculations and the list of validated

98   residues of RAG1 and RAG2.

99   *MRF score versus known variant pathogenicity measure*

100  The functional validation of these predictions is presented in Fig 1(v). We have

101  previously measured the recombination activity of RAG1 and RAG2 disease-

102  causing variants in several patients [14]. We have combined the known func-

103  tional activity from other extensive reports [15], to compare a total of 44 vari-

104  ants. RAG deficiency is measured by the level of recombination potential. We

105  expected that damaging mutations (resulting in low recombination activity in

106  vitro) would be identified with high MRF scores. MRF pathogenicity prediction

107  correctly identified damaging mutations in RAG1 and RAG2 (Fig 1(v)). Variants

108  reported on GnomAD which are clinically found to cause disease have signifi-

109  cantly higher MRF scores than variants which have not been reported to cause

110  disease (Fig 1(v)). Table S1 provides all MRF scores for both proteins as well as

111  raw data used for calculations and the list of validated residues of RAG1 and

112  RAG2.

113      Allele frequency is generally the single most important filtering method for

114  rare disease in whole genome (and exome) sequencing experiments. *RAG1* and

115  *RAG2* have probability of being loss-of-function intolerant (pLI) scores of 0.00

116  and 0.01, respectively. Mutations under pressure from purifying selection are

117  more likely to cause disease than common variants. However, allele frequen-

118  cies of rare variants reported on GnomAD cannot differentially predict likeli-

6

119  hood of causing disease. This is particularly important for recessive diseases

120  such as RAG deficiency. As such we find no significant difference between clini-

121  cally damaging variants and those which have not been reported yet as disease-

122  causing, illustrating the reasoning for our method design (Fig 1 (vi)). Many

123  non-clinically-reported rare variants may cause disease; the MRF score identi-

124  fies the top clinically-relevant candidates. Conserved residues with the highest

125  MRF for both RAG1 and RAG2 are mapped onto the protein structure in Fig 3

126  and frequently show high MRF at DNA contact points. The accuracy for cor-

127  rectly identifying all disease-causing variants reported to date is shown in (Fig

128  1(vii). We found >80% accuracy for 21 known variants tested, >50% accuracy

129  for 48 tested and <50% accuracy for only 23 tested. The raw values comparing

130  functional pathogenicity and MRF scores are illustrated in Fig 2.

131  *False positives in Transib domains do not worsen probability prediction*

132  A set of conserved motifs in core *RAG1* are shared with the *Transib* transposase,

133  including the critical DDE residue catalytic triad [16]. Ten *RAG1* core motifs are

134  conserved amongst a set of diverse species including human [16]. To assess the

135  influence of false positive effect on MRF prediction the conserved residues in

136  this dataset are compared to GnomAD allele frequencies and MRF score. Fig 4

137  (i) plots the MRF (lacking the Boolean component *C*) for conserved *Transib* mo-

138  tif residues, non-conserved *Transib* motif residues, and non-*Transib* residues.

139  Fig 4 (ii) shows the percentage of these which are reported as mutated on Gno-

140  mAD. Removing reported variants by applying *C*, the resulting effect on incor-

141  rectly scoring MRF in the conserved *Transib* motifs remains neutral. Com-

7

142 bined Annotation Dependent Depletion (CADD) scoring [17] is an important

143 bioinformatics filtering method. We compare MRF to the PHRED-scaled *RAG1*

144 CADD scores for all possible SNVs (Fig 5). While CADD is a valuable scoring

145 method its purpose is not to predict likelihood of variation.

**Discussion**

147 Determining disease-causing variants for functional analysis typically aims to

148 target conserved gene regions. On GnomAD 55.99% of *RAG1* (approx. 246,000

149 alleles) has no reported variants. Functionally validating unknown variants in

150 genes with this level purifying selection is generally infeasible. Conserved re-

151 gions are likely high importance regions, yet determining the likelihood of pa-

152 tients presenting with mutations in these clusters requires a scoring mecha-

153 nism. An example of such clustering of highly scoring MRFs occured in the

154 RAG1 catalytic RNase H (RNH) domain at p.Ser638-Leu658 which is also con-

155 sidered a conserved *Transib* motif. Targeting clearly defined regions with high

156 MRF scores allows for functional validation studies tailored to the most clinically-

157 relevant protein regions. Phenotypic, epigenetic, or other such weighting data

158 may also be applied to this model. Variants with a low MRF may still be damag-

159 ing but resources for functional validation are best spent on gene regions with

160 high MRF. Table S1 lists the values for calculated MRFs for RAG1 and RAG2.

161 We have presented a basic application of MRF scoring for RAG deficiency.

162 Furthermore, we have suggested its genome-wide application with to the infor-

163 mation retrieval method; term frequency, inverse document frequency ($tf -$

164 $idf$). In this case the"term" will represent an amino acid residue $r$ while the

8

165 "document" represents a gene $g$ such that,

$$rf - igf_{r,g} = rf_{r,g} \times igf_r \tag{1}$$

166 We may view each gene as a vector with one component corresponding to each
167 residue mutation in the gene, together with a weight for each component that is
168 given by (1). Therefore, we can find the overlap score measure with the $rf - igf$
169 weight of each term in $g$.

$$\text{Score}(q, g) = \sum_{r \in q} \text{rf-igf}_{r,g}.$$

170 We expand here briefly on the technical description of this method. Log weight-
171 ing may offer clearer disease-causing variant discovery depending on the scor-
172 ing method. In respect to MRF scoring, this information retrieval method might
173 be applied as follows; the $rf - igf$ weight of a term is the product of its $rf$
174 weight and its $igf$ weight ($W_{r,g} = rf_{r,g} \times \log \frac{N}{gf_r}$) or ($W_{r,g} = (1 + \log rf_{r,g}) \times$
175 $\log \frac{N}{gf_r}$). That is, firstly, the number of times a residue mutates in a gene ($rf =$
176 $rf_{r,g}$). Secondly, the rarity of the mutation genome-wide in $N$ number of genes
177 ($igf = N/gf_r$). Finally, ranking the score of genes for a mutation query $q$ by;

$$\text{Score}(q, g) = \sum_{r \in q \cap g} \text{rf-igf}_{r,g}$$

178 The score of the query (Score($q, g$)) equals the mutations (terms) that appear
179 in both the query and the gene ($r \in q \cap g$). Working out the $rf - igf$ weight for
180 each of those variants ($rf.igf_{r,g}$) and then summing them ($\sum$) to give the score
181 for the specific gene with respect to the query.

182     During clinical investigations using personalised analysis of patient data,
183 further scoring methods may be applied based on disease features. A patient

9

184 with autoinflammatory features may require weighting for genes such as *MEFV*

185 and *TNFAIP3*, whereas a patient with mainly immunodeficiency may have weighted

186 scoring for genes such as *BTK* and *DOCK8*. A patient phenotype can contribute

187 a weight based on known genotype correlations separating primary immunod-

188 eficiencies or autoinflammatory diseases [6]. However, validation of these ex-

189 panded implementations requires a deeper consolidation of functional stud-

190 ies than is currently available. A method with similar possible applications for

191 human health mapping constrained coding regions has been recently released

192 [18]. This study employed a method which included weighting by sequencing

193 depth. We have not included this method as our analysis was gene-specific but

194 implementation is advised when calculating genome-wide MRF scores.

195 Predicting the likelihood of discovering novel mutations has implications

196 in genome-wide association studies (GWAS). Variants with low minor allele fre-

197 quencies have a low discovery rate and low probability of disease association

198 [19]; an important consideration for rare diseases such as RAG deficiency. An

199 analysis of the NHGRI-EBI catalogue data highlighted diseases whose average

200 risk allele frequency was low. Autoimmune diseases had risk allele frequen-

201 cies considered low at approximately 0.4 [19]. Without a method to rank most

202 probable novel disease-causing variants, it is unlikely that GWAS will identify

203 very rare disease alleles (with frequencies <0.001). It is conceivable that a num-

204 ber of rare immune diseases are attributable to polygenic rare variants. How-

205 ever, evidence for low-frequency polygenic compounding mutations will not be

206 available until large, accessible genetics databases are available, exemplified by

207 the NIHR BioResource Rare Diseases study [14]. An interesting consideration

10

208 when predicting probabilities of variant frequency, is that of protective muta-

209 tions. Disease risk variants are quelled at low frequency by negative selection,

210 while protective variants may drift at higher allele frequencies [20].

211     The cost-effectiveness of genomic diagnostic tests is already outperforming

212 traditional, targeted sequencing [1]. Even with substantial increases in data

213 sharing capabilities and adoption of clinical genomics, rare diseases due to

214 variants of unknown significance and low allele frequencies (<0.0001) will re-

215 main non-actionable until reliable predictive genomics practices are developed.

216 Bioinformatics a a whole has made staggering advances in the field of genet-

217 ics [21]. Challenges which remain unsolved, hindering the benefit of national

218 or global genomics databases, include DNA data storage and random access

219 retrieval [22], data privacy management [23], and predictive genomics analy-

220 sis methods. Variant filtration in rare disease is based on reference allele fre-

221 quency, yet the result is not clinically actionable in most cases. Development of

222 predictive genomics tools may provide a critical role for single patient studies

223 and timely diagnosis [13].

224 **Conclusion**

225 We provide the amino acid residue list for RAG1 and RAG2 which have not been

226 reported to date but are most likely to present clinically as RAG deficiency. This

227 method may be applied to other diseases with hopes of improving prepared-

228 ness for clinical diagnosis.

11

**Methods**

*Population genetics*

GnomAD (version r2.0.2) [24] was queried for the canonical transcripts of *RAG1* and *RAG2* from population genetics data of approximately 146,000 alleles; ENST00000299440 (*RAG1*) 1495 vaiants (including filtered: 1586), GRCh37 11:36532259-36614706 and ENST00000311485 (*RAG2*) 786 varaitns (including filtered: 831), GRCh37 11:36597124 - 36619829. Data was filtered to contain the identifiers: frameshift, inframe deletion, inframe insertion, missense, stop lost, or stop gained. Reference transcripts were sourced from Ensembl in the FASTA format amino acid sequence; transcript: RAG1-201 ENST00000299440.5 [HGNC:9831] and transcript: RAG2-201 ENST00000311485.7 [HGNC:9832]. These sequences were converted to their three-letter code format using One to Three from the Sequence Manipulation Suite (http://bioinformatics.org/sms2/mirror.html).

Input sets used GnomAD variant allele frequencies and reference sequences processed as cvs files, cleaned and sorted to contain only coding amino acid residues, amino acid code, residue number, alternate variants, allele frequencies of variants, and a score (*C*) of 0 or 1 where 1 represented no reported variants. A score was also given where multiple alternate variants existed. A separate statistics report was generated from this processed input data. The percentage of conserved residues was calculated (55.99% of amino acids contained no reported variants in RAG1, 55.98% in RAG2). The count of variants per residue was found for both proteins. The ratio was also found per residue

12

conservation rate / mutation rate. Basic protein statistics were generated using reference canonical transcript sequences of RAG1 and RAG2 with the Sequence Manipulation Suite. The residue frequency was calculated based on the respective polypeptide chain length.

The calculated mutation rate value and residue frequency score together produce the mutation rate residue frequency as shown in Table S1. Our investigation used the Boolean $C$ score of 0 or 1 to weight mutation rate residue frequencies. An important consideration for future application is whether to use this Boolean score or a frequency score. In the clinical setting, the likelihood of *de novo* mutations versus inherited mutations have different impact on recessive and dominant diseases. The likelihood of a patient presenting with a particular (predicted) variant is more likely if the variant exists even at a very low frequency in the patients ancestral population. Therefore, an allele frequency may be used to replace $C$ in many investigations.

*Data visualisation*

For our visualisation of MRF scores, small clusters of high MRF were of more significance than individual highly conserved residues. Therefore, we applied a 1% average filter where values were averaged over a sliding window of N number of residues (10 in the case of RAG1, 6 in the case of RAG2). However, when using Boolean scoring $C$, this method should be applied before $C$. Alternatively, if using allele frequency scoring, this visualisation method can be applied subsequently. Lastly, for a clear distinction of MRF clusters a cut-off threshold was applied at the 75th percentile (0.0168 in RAG1).

13

A gene map for coding regions in RAG1 and RAG2 were populated with (1) Boolean *C* score from population genetics data, (2) raw MRF scores, and (3) MRF clusters with 1% average and cutoff threshold. GraphPad Prism was used for heatmaps and Adobe Adobe Illustrator and Photoshop were used for protein domain illustrations.

*Validation of MRF against functional data*

The recombination activity of RAG1 and RAG2 was previously measured on 44 known pathogenic variants [14, 15]. Briefly, the pathogenicity of variants in RAG1 and RAG2 are measured functionally *in vitro* by expression of RAG1 and RAG2 in combination with a recombination substrate plasmid containing re-combination signal sequences which are targeted by RAG complex during nor-mal V(D)J recombination. Recombination events are assessed by quantitative real-time PCR using comparative CT. The inverse score of recombination activ-ity (0-100%) is used to quantify pathogenicity of variants in our study. Compar-ison between known pathogenicity scores and MFR was done by scaling MRF scores from 0-100% (100% being highest probability of occurring as damaging).

**References**

[1] Katherine Payne, Sean P Gavan, Stuart J Wright, and Alexander J Thomp-son. Cost-effectiveness analyses of genetic and genomic diagnostic tests. *Nature Reviews Genetics*, 2018.

[2] Antonia Kwan, Roshini S Abraham, Robert Currier, Amy Brower, Karen An-druszewski, Jordan K Abbott, Mei Baker, Mark Ballow, Louis E Bartoshesky,

14

Vincent R Bonagura, et al. Newborn screening for severe combined immunodeficiency in 11 screening programs in the united states. *Jama*, 312 (7):729–738, 2014.

[3] L. Alexander Liggett, Anchal Sharma, Subhajyoti De, and James DeGregori. Conserved patterns of somatic mutations in human peripheral blood cells. *bioRxiv*, 2017. doi: 10.1101/208066.

[4] Kapourani Chantriolnt-Andreas Stubbs Thomas M. Lee Heather J. Alda-Catalinas Celia Krueger Felix Sanguinetti Guido Kelsey Gavin Marioni John C. Stegle Oliver Reik Wolf Clark Stephen J., Argelaguet Ricard. scnmt-seq enables joint profiling of chromatin accessibility dna methylation and transcription in single cells. *Nature Communications*, 9(1):781, 2018. ISSN 2041-1723. doi: 10.1038/s41467-018-03149-4.

[5] István Bartha, Julia di Iulio, J Craig Venter, and Amalio Telenti. Human gene essentiality. *Nature Reviews Genetics*, pages nrg–2017, 2017.

[6] Capucine Picard, Waleed Al-Herz, Aziz Bousfiha, Jean-Laurent Casanova, Talal Chatila, Mary Ellen Conley, Charlotte Cunningham-Rundles, Amos Etzioni, Steven M Holland, Christoph Klein, et al. Primary immunodeficiency diseases: an update on the classification from the international union of immunological societies expert committee for primary immunodeficiency 2015. *Journal of clinical immunology*, 35(8):696–726, 2015.

[7] Mary Ellen Conley and Jean-Laurent Casanova. Discovery of single-gene

15

318    inborn errors of immunity by next generation sequencing. *Current opin-*

319    *ion in immunology*, 30:17–23, 2014.

320  [8] David G Schatz, Marjorie A Oettinger, and David Baltimore. The v (d) j

321    recombination activating gene, rag-1. *Cell*, 59(6):1035–1048, 1989.

322  [9] Marjorie A Oettinger, David G Schatz, Carolyn Gorka, and David Balti-

323    more. Rag-1 and rag-2, adjacent genes that synergistically activate v (d)

324    j recombination. *Science*, 248(4962):1517–1523, 1990.

325  [10] Jolan E Walter, Lindsey B Rosen, Krisztian Csomos, Jacob M Rosenberg,

326    Divij Mathew, Marton Keszei, Boglarka Ujhazi, Karin Chen, Yu Nee Lee,

327    Irit Tirosh, et al. Broad-spectrum antibodies against self-antigens and cy-

328    tokines in rag deficiency. *The Journal of clinical investigation*, 125(11):

329    4135–4148, 2015.

330  [11] Catharina Schuetz, Kirsten Huck, Sonja Gudowius, Mosaad Megahed,

331    Oliver Feyen, Bernd Hubner, Dominik T Schneider, Burkhard Manfras, Ul-

332    rich Pannicke, Rein Willemze, et al. An immunodeficiency disease with

333    rag mutations and granulomas. *New England Journal of Medicine*, 358(19):

334    2030–2038, 2008.

335  [12] Tami John, Jolan E Walter, Catherina Schuetz, Karin Chen, Roshini S Abra-

336    ham, Carmem Bonfim, Thomas G Boyce, Avni Y Joshi, Elizabeth Kang,

337    Beatriz Tavares Costa Carvalho, et al. Unrelated hematopoietic cell trans-

338    plantation in a patient with combined immunodeficiency with granulo-

16

339    matous disease and autoimmunity secondary to rag deficiency. *Journal of*
340    *clinical immunology*, 36(7):725–732, 2016.

341    [13] Jean-Laurent Casanova, Mary Ellen Conley, Stephen J Seligman, Laurent
342    Abel, and Luigi D Notarangelo. Guidelines for genetic studies in single pa-
343    tients: lessons from primary immunodeficiencies. *Journal of Experimen-*
344    *tal Medicine*, pages jem–20140520, 2014.

345    [14] Dylan Lawless, Christoph B Geier, Jocelyn R Farmer, Hana Allen Lango,
346    Daniel Thwaites, Faranaz Atschekzei, Matthew Brown, David Buchbinder,
347    Siobhan O Burns, Manish J Butte, et al. Prevalence and clinical chal-
348    lenges among adult primary immunodeficiency patients with rag defi-
349    ciency. *Journal of Allergy and Clinical Immunology*.

350    [15] Yu Nee Lee, Francesco Frugoni, Kerry Dobbs, Irit Tirosh, Likun Du,
351    Francesca A Ververs, Heng Ru, Lisa Ott de Bruin, Mehdi Adeli, Jacob H
352    Bleesing, et al. Characterization of t and b cell repertoire diversity in pa-
353    tients with rag deficiency. *Science immunology*, 1(6), 2016.

354    [16] Vladimir V Kapitonov and Jerzy Jurka. Rag1 core and v (d) j recombination
355    signal sequences were derived from transib transposons. *PLoS biology*, 3
356    (6):e181, 2005.

357    [17] Martin Kircher, Daniela M Witten, Preti Jain, Brian J O'roak, Gregory M
358    Cooper, and Jay Shendure. A general framework for estimating the relative
359    pathogenicity of human genetic variants. *Nature genetics*, 46(3):310, 2014.

17

[18] James M Havrilla, Brent S Pedersen, Ryan M Layer, and Aaron R Quinlan. A map of constrained coding regions in the human genome. *bioRxiv*, 2017. doi: 10.1101/220814.

[19] Takashi Kido, Weronika Sikora-Wohlfeld, Minae Kawashima, Shinichi Kikuchi, Naoyuki Kamatani, Anil Patwardhan, Richard Chen, Marina Sirota, Keiichi Kodama, Dexter Hadley, et al. Are minor alleles more likely to be risk alleles? *BMC medical genomics*, 11(1):3, 2018.

[20] Yingleong Chan, Elaine T Lim, Niina Sandholm, Sophie R Wang, Amy Jayne McKnight, Stephan Ripke, Mark J Daly, Benjamin M Neale, Rany M Salem, Joel N Hirschhorn, et al. An excess of risk-increasing low-frequency variants can be a signal of polygenic inheritance in complex diseases. *The American Journal of Human Genetics*, 94(3):437–452, 2014.

[21] Maxwell W Libbrecht and William Stafford Noble. Machine learning applications in genetics and genomics. *Nature Reviews Genetics*, 16(6):321, 2015.

[22] Lee Organick, Siena Dumas Ang, Yuan-Jyue Chen, Randolph Lopez, Sergey Yekhanin, Konstantin Makarychev, Miklos Z. Racz, Govinda Kamath, Parikshit Gopalan, Bichlien Nguyen, Christopher Takahashi, Sharon Newman, Hsing-Yeh Parker, Cyrus Rashtchian, Kendall Stewart, Gagan Gupta, Robert Carlson, John Mulligan, Douglas Carmean, Georg Seelig, Luis Ceze, and Karin Strauss. Scaling up dna data storage and random access retrieval. *bioRxiv*, 2017. doi: 10.1101/114553.

18

[23] Zhicong Huang, Erman Ayday, Huang Lin, Raeka S. Aiyar, Adam Molyneaux, Zhenyu Xu, Jacques Fellay, Lars M. Steinmetz, and Jean-Pierre Hubaux. A privacy-preserving solution for compressed storage and selective retrieval of genomic data. *Genome Research*, 26(12):10. 1687–1696, 2016.

[24] Monkol Lek, Konrad J Karczewski, Eric V Minikel, Kaitlin E Samocha, Eric Banks, Timothy Fennell, Anne H OâĂŹDonnell-Luria, James S Ware, Andrew J Hill, Beryl B Cummings, et al. Analysis of protein-coding genetic variation in 60,706 humans. *Nature*, 536(7616):285, 2016.

19

Table S1: MRF data tables. The complete RAG1 and RAG2 amino acid residue MRF scores are provided along with known clinically pathogenic variant residues and raw data used for calculations.
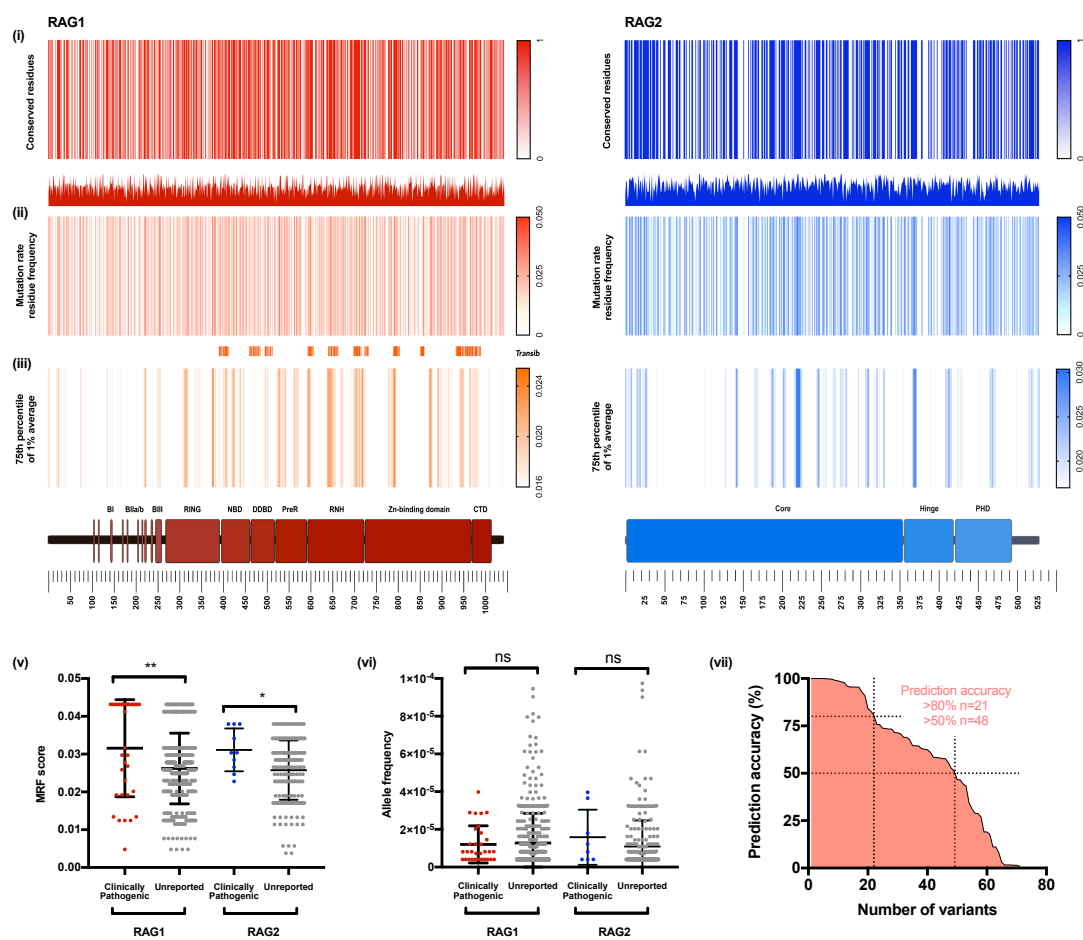
Figure 1: RAG1 and RAG2 conservation and mutation rate residue frequency. (i) Gene conservation score, non-conserved 0 and conserved 1. (ii) Histogram; raw MRF score. Heatmap; MRF prediction for conserved residues, graded 0 to 0.05. (iii) MRF score averaged with 1% intervals for each respective gene and cut-off below 75th percentile, graded 0 to 0.03 (Noise reduction method). (iv) Gene structure with functional domains. (v) Clinically damaging variants reported on GnomAD have significantly higher MRF scores than non-pathogenic variants. (Unpaired t test. RAG1 P value 0.002** RAG2 P value 0.0339*). (vi) GnomAD allele frequency <0.0001. No significant difference in allele frequency is found between clinically damaging variants and non-clinically reported. (vi) Accuracy of MRF scoring compared to functionally validated pathogenicity.
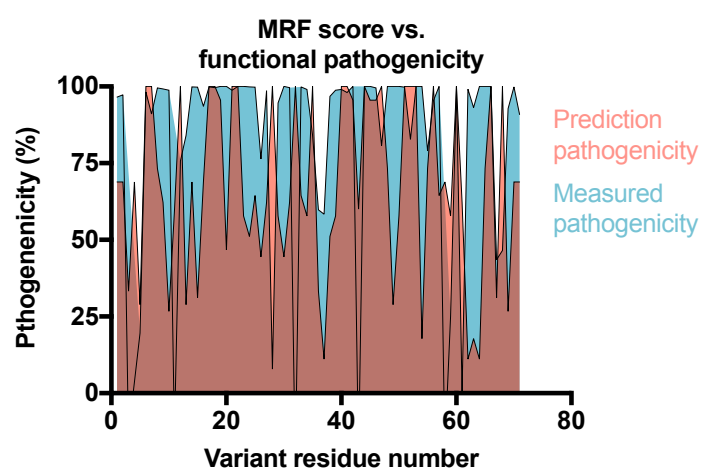
21

Figure 2: MRF score vs. known variant pathogenicity measure. Predicted pathogenicity likelihood (based on maximum and minimum MRF score as a percentage) is shown in red. In blue, the functionally measured recombination activity of each variant where complete loss of protein activity is measured as 100% pathogenicity. These values are summarised in Fig 1v(ii).
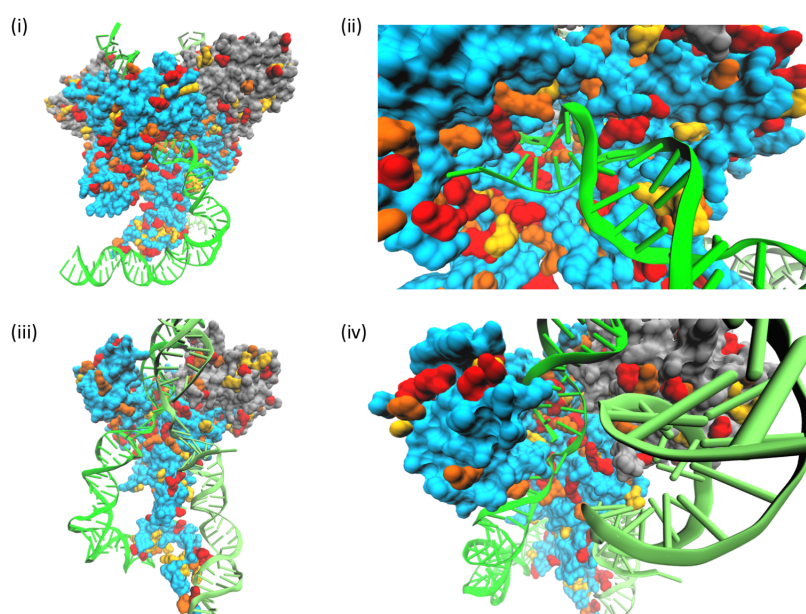
Figure 3: The RAG1 (blue) and RAG2 (grey) protein structure with MRF scores. (i) Protein dimers and (ii=iv) monomers illustrating the three highest category MRF scores for predicted clinically-relevant variants. Increasing in MRF score; yellow, orange, red. DNA contact points are integral to protein function and generally score as high MRF residues. (PDB:3jbw)
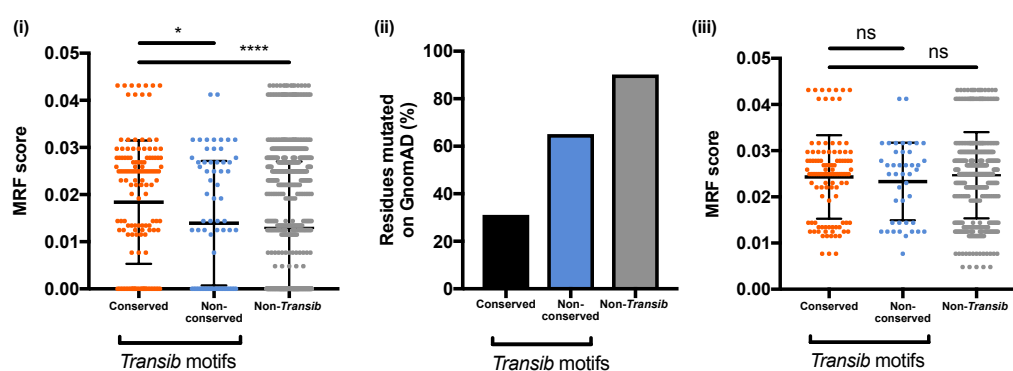
Figure 4: False positives in *Transib* domains do not worsen probability prediction. The *Transib* domains contain critical conserved protein residues. (i) False positives are simulated by scoring *Transib* domains MRF without their Boolean conservation weight $C$. (ii) Allele frequencies on GnomAD have inversely proportional conservation to simulated false-positive MRF scoring. (iii) When the Boolean component $C$ is applied in MRF calculation the effect of false positives remains non-significant, illustrating the non-negative impact of MRF for pathogenicity rate prediction.
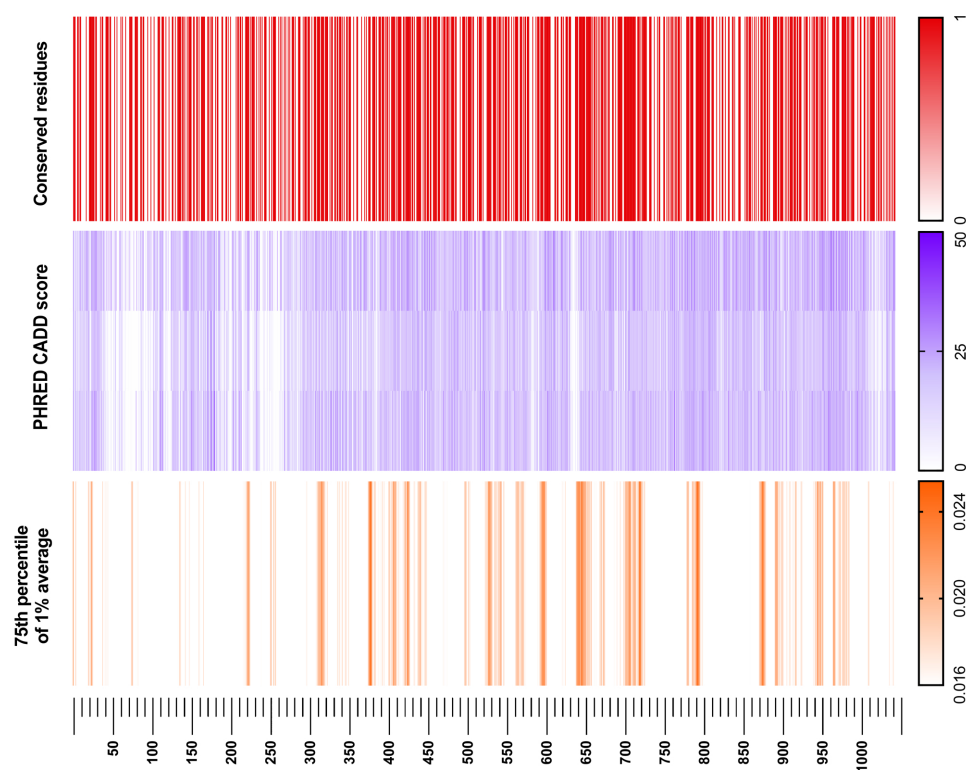
Figure 5: *RAG1* PHRED-scaled CADD score versus GnomAD conservation rate and MRF score. Allele frequency conservation rate (top) is vastly important for identifying critical structural and functional protein regions. The impact of mutation in one of these conserved regions is often estimated using CADD scoring (middle). The MRF score (bottom)(visualised using the 75th percentile with 1% averaging) highlights protein regions which are most likely to present clinically and may require pre-emptive functional investigation.