

1 **De novo mutations implicate novel genes with burden of rare variants in Systemic**
2 **Lupus Erythematosus**

3

4 Venu Pullabhatla^{1*}, Amy L. Roberts^{2*}, Myles J. Lewis³, Daniele Mauro³, David L. Morris²,
5 Christopher A. Odhams², Philip Tomblinson², Ulrika Liljedahl⁴, Simon Vyse^{2,7}, Michael A.
6 Simpson², Sascha Sauer^{5,8}, Emanuele de Rinaldis¹, Ann-Christine Syvänen⁴, Timothy J.
7 Vyse^{2,6,9}

8 * Contributed equally to this work

9

10 ¹ NIHR GSTFT/KCL Comprehensive Biomedical Research Centre, Guy's & St. Thomas'
11 NHS Foundation Trust, London, UK

12

13 ² Division of Genetics and Molecular Medicine, King's College London, London, UK

14

15 ³ Centre for Experimental Medicine and Rheumatology, William Harvey Research Institute,
16 Barts and The London School of Medicine and Dentistry, Queen Mary University of London,
17 London, UK

18

19 ⁴ Department of Medical Sciences, Uppsala University, Uppsala, Sweden

20

21 ⁵ Max Planck Institute for Molecular Genetics, Berlin, Germany

22

23 ⁶ Division of Immunology, Infection and Inflammatory Disease, King's College London,
24 London, UK

25

26 ⁷ Current address: Division of Cancer Biology, The Institute of Cancer Research, London,
27 UK

28

29 ⁸ Current address: Max Delbrück Centre for Molecular Medicine (BIMSB/BIH), Berlin,
30 Germany

31

32 ⁹ Corresponding author: timothy.vyse@kcl.ac.uk

33

34

35

36 **Abstract**

37 The heritability of most complex diseases, including autoimmune disease Systemic Lupus
38 Erythematosus (SLE), remains largely unexplained by common variation, and few examples
39 of rare variant associations have been identified. Here, using complementary whole-exome
40 sequencing (WES) and high-density imputation, we identify candidate genes through *de*
41 *novo* mutation discovery and demonstrate collective rare variant associations at novel SLE-
42 susceptibility genes. Using extreme-phenotype sampling, we sequenced the exomes of 30
43 SLE parent-affected-offspring trios and identified 14 genes with missense *de novo*
44 mutations, none of which are within the >70 SLE susceptibility loci implicated through
45 genome-wide association studies (GWAS). In a follow-up cohort of 10,995 individuals of
46 matched European ancestry, including 4,036 SLE cases, we imputed genotype data to the
47 density of the combined UK10K-1000 genomes Phase III reference panel across the 14
48 candidate genes. We identify a burden of rare exonic variants across *PRKCD* associated
49 with SLE risk ($P=0.0028$), and across *DNMT3A* associated with two severe disease
50 prognosis sub-phenotypes ($P=0.0005$ and $P=0.0033$). Additionally, we show the
51 p.His198Gln *de novo* mutation within the candidate gene *C1QTNF4* inhibits NF- κ B activation
52 following TNF exposure. Exome sequencing studies typically lack power to detect rare
53 variant associations for complex traits. Our results support extreme-phenotype sampling and
54 using *de novo* mutation gene discovery to aid the search for rare variation contributing to the
55 heritability of complex diseases.

56

57

58

59

60

61

62 **Introduction**

63 Considerable progress has been made in elucidating the genetic basis of complex disease
64 traits. The associated genetic polymorphisms identified are usually relatively common in the
65 population and the risk alleles impart a modest individual increment to the likelihood of
66 developing disease. The allelic identities of such genetic factors are established using large-
67 scale genotyping chips with a predetermined composition of tagging variants. Advances in
68 DNA sequencing technology, such as next-generation sequencing (NGS), enable the
69 characterization of rare and unique genetic variants. Targeting the exome, which comprises
70 approximately 1% of the genome, facilitates (by scale) analysis of NGS data and provides a
71 highly enriched source of highly penetrant disease causing mutations¹. The impact of NGS
72 has been dramatic in monogenic disease; in a recent review it was stated that the genes
73 underlying approximately half of known Mendelian disorders had been discovered². In
74 contrast, the role of rare genetic variants in complex diseases is unknown. It had been
75 proposed that rare variants might explain, at least in part, the missing heritability that is
76 frequently described in complex disease traits³. Yet very few examples of individual rare
77 variants contributing to complex disease risk have been identified⁴. However, rare variants
78 are collectively very common; indeed the Exome Aggregate Consortium (ExAC) have
79 demonstrated that the vast majority of genetic variation is extremely rare⁵. Therefore, as
80 opposed to case-control analyses on the *variant*-level – as routinely employed for common
81 polymorphisms – the contribution of rare variation can be assessed on the *gene*-level by
82 aggregating all observed rare variants in a defined region and performing burden tests.

83 Although there are examples of disease predisposition genes, typically identified through
84 genome-wide association studies (GWAS), harbouring associated variants across a
85 spectrum of allele frequencies^{4,6,7}, studies focusing on canonical disease-associated loci
86 have been far from fruitful, suggesting these loci by and large do not harbour additional risk
87 through rare variation⁸. Exome-wide searches have similarly revealed limited numbers of
88 rare variation associated with complex diseases; a recent large-scale whole-exome

89 sequencing case-control study in Type 2 Diabetes concluded that rare variants do not play a
90 major role in disease predisposition⁹. Furthermore, widely used gene-based association
91 tests have been shown to lack power at the exome-wide level, even with sample sizes up to
92 10,000, suggesting the need for reduced number of tested regions¹⁰.

93 Our strategy to address this problem is outlined here and summarised in **Fig. 1**. We selected
94 SLE cases with a severe phenotype (young age of onset and clinical features associated
95 with poorer outcome) and hypothesized that these individuals would exhibit unique mutation
96 events in their protein coding DNA that predisposed to disease risk. Therefore, we undertook
97 whole exome sequencing (WES) in 30 family trios (both parents and affected offspring) and
98 scrutinized the data for *de novo* mutations in the individual with SLE to identify a group of
99 candidate genes for an independent follow-up rare variant analysis. This method allowed the
100 identification of novel loci harbouring disease risk through collective rare variation.

101

102

103

104

105

106

107

108

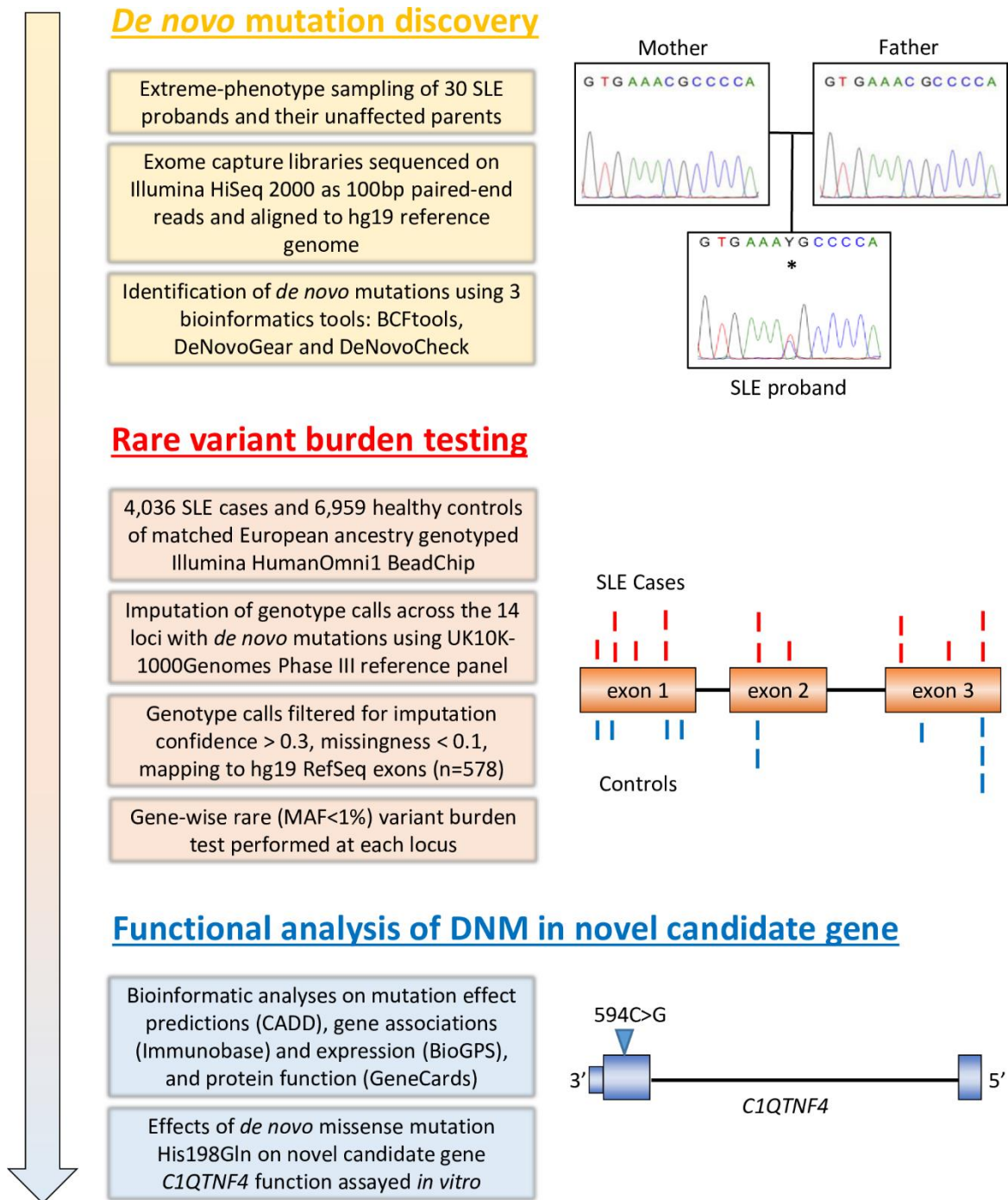
109

110

111

112 **Figure 1 Overview of study**

113 *De novo* mutations (DNM) in a discovery cohort revealed candidate genes for imputation-
114 based rare variant burden testing using a follow-up cohort. Independent functional analyses
115 demonstrate the functional effects of one DNM in a candidate gene.



116

117

118 **Results**

119 ***Identification of de novo mutations in SLE cases***

120 We screened for *de novo* mutations by WES of 30 family trios with an affected offspring with
121 more severe SLE (**Supplementary Fig. 1**). A total of 584,798 variants ($\geq 20X$), including
122 single nucleotide variants and indels, were identified in the 30 affected probands. Using
123 three bioinformatic tools and employing conservative parameters, 26 putative exonic *de*
124 *novo* mutations were identified across the 30 trios, at an average of 0.87 variants per
125 proband, including 17 missense mutations across 17 genes (**Supplementary Table 1**;
126 **Supplementary Fig. 2**). These data fit with the estimation that each exome contains one *de*
127 *novo* mutation¹¹.

128 We also analysed the SLE proband WES data alone, without the unaffected parents. This
129 revealed 1,194 non-silent, heterozygous, rare variants in 1,067 genes, which would make
130 prioritisation for downstream analysis a difficult task, highlighting the benefit of parent-
131 offspring trio sequencing (**Supplementary Fig. 3**).

132 Through Sanger sequencing all three members of the parent-offspring trio, plus any
133 additional unaffected siblings (**Supplementary Table 2**), we confirmed 14 true non-silent *de*
134 *novo* mutations (**Table 1**) in 11 of the 30 probands (36.7%) for further analysis. These true
135 *de novo* mutations were absent in both parents and any unaffected siblings, but present in
136 the SLE proband.

137 Of the three mutations that did not pass Sanger verification (**Supplementary Table 1**), one,
138 within *LAMC2*, is likely a result of germline mosaicism because, although not observed in
139 either parent, it is observed in an unaffected sibling in addition to the SLE proband¹². It is
140 therefore perhaps unlikely to be pathogenic and was not taken forward in downstream
141 analyses. The two remaining putative mutations in *KLRC1* and *KRTAP10-2* are both
142 member of highly homologous gene families. Such sequence identity may have caused false
143 positive identification of *de novo* mutations in the WES analysis. Indeed the *KLRC1*

144 p.Ile225Met missense variant appears to be a polymorphic Paralogous Sequence Variant
145 (PSV) – the paralogous variant being Met223Ile in *KLRC2* - once more making this variant
146 unlikely to be pathogenic.

147

148 ***Supporting evidence of a role of candidate genes in SLE***

149 We explored the function, expression (BioGPS), known associations with autoimmunity
150 (ImmunoBase), and gene-level constraint against missense mutations (ExAC), of the genes
151 with *de novo* mutations to build a profile of *a priori* evidence of a role in SLE pathogenesis
152 (**Table 2**). The candidate genes include autoimmune susceptibility genes (*PRKCD*, *DNMT3A*
153 and *ANXA3*), although none have been previously associated with SLE through GWAS in
154 any population^{13,14}. We also identify candidate genes through known/predicted function and
155 expression profiles (*C1QTNF4*, *SRRM2*, *HMSD*), and four genes (*PRKCD*, *DNMT3A*,
156 *C1QTNF4* and *LRP1*) with a significant ($Z > 3.09$) constraint against missense variants (**Table**
157 **2**).

158

159 ***Functional characterisation of de novo mutations***

160 To analyse the potential impact of the *de novo* mutations, we used the ExAC database⁵ and
161 Combined Annotation Dependent Depletion (CADD) scores¹⁵ to characterise their frequency
162 and predicted functional effects, respectively. Five of the 14 *de novo* mutations – found in
163 *MICALL1*, *LRP1*, *PNPLA1*, *PLD1*, and *GFTP2* - have been observed, at very rare
164 frequencies, in the ~60,000 exomes documented in ExAC (**Table 1**). All five mutations are
165 CpG transitions and therefore likely to be identity-by-state, reflecting the higher mutability
166 rate of these sites. Within the mutation set, five (35.7%) – found in *DNMT3A*, *PRKCD*,
167 *MICALL1*, *LRP1*, and *PNPLA1* – have CADD Phred scores > 30 , placing them in the top
168 0.1% of possible damaging mutations in the human genome (**Table 1**).

169 ***Independent gene-based analysis of rare variants in de novo mutated genes***

170 We hypothesised that, while some observed *de novo* mutations were random background
171 variation as present in the exome of every individual regardless of disease status¹¹, others
172 could be contributing to disease through their deleterious effects on the encoded protein,
173 and would thus be indicative of a hitherto unknown gene which contributes to SLE risk. Such
174 putative genes could also harbour a burden of additional rare variation. Therefore, in a
175 follow-up cohort of 10,995 individuals of matched European ancestry previously genotyped
176 on the Illumina HumanOmni1 BeadChip¹³, including 4,036 SLE cases, we imputed genotype
177 data to the density of the combined UK10K and 1000 genomes Phase III reference panel
178 (UK10K-1000GP3) across the 14 genes with *de novo* mutations to assay rare variation risk.
179 Using a collapsing burden test¹⁶, we surveyed each of the 14 genes for an excess of
180 aggregated rare (minor allele frequencies (MAF)<1%) exonic variants in SLE cases
181 compared to healthy controls, and we identify an association of rare exonic variation in
182 *PRKCD* with SLE (**Table 3**; $P=0.0028$). In sub-phenotype analyses, using healthy controls
183 and only SLE cases with anti-dsDNA ($n_{\text{cases}}=1261$) or renal-involvement with
184 hypocomplementemia ($n_{\text{cases}}=186$), both of which are markers of more severe disease, we
185 identify collective rare exonic variants in *DNMT3A* associated with both anti-dsDNA (**Table**
186 **3**; $P=0.0005$) and renal involvement with hypocomplementemia (**Table 3**; $P=0.0033$). We
187 also collapsed all exons from the 14 genes together to test for an overall burden of rare
188 variants across these loci. These analyses revealed no excess of rare exonic variants across
189 the grouped genes, reflecting the hypothesis that some/most genes will not be relevant to
190 disease status because the observed *de novo* mutations are random background variation
191 only.

192 Using gene-level constraint metric data from ExAC⁵, *DNMT3A* and *PRKCD* are two of the
193 four genes with *de novo* mutations with a significant constraint against missense variants (Z
194 >3.09 ; **Table 2**). However, across the entire gene set, there was no difference in the median
195 Z-score (0.50) compared with the median Z-score across all genes in ExAC (0.51). These

196 data reflect the results of our rare variant burden tests, in which the aggregated gene set do
197 not contribute to disease risk.

198

199 ***Common variant analysis across de novo mutation genes***

200 Loci harbouring autoimmune risk through both common and rare variants have been
201 reported^{6,17}. Therefore, using the high-density UK10K-1000GP3 imputed data, we
202 reassessed the contribution of common variation to SLE risk across these loci. No significant
203 association at any locus was observed with overall risk in a case-control comparison, as
204 previously reported¹³, nor with anti-dsDNA ($n_{\text{cases}}=1261$) or renal-involvement with
205 hypocomplementemia ($n_{\text{cases}}=186$) sub-phenotypes (**Supplementary Table 3**). A candidate
206 gene study previously reported a trend of association between the common *DNMT3A*
207 intronic SNP rs1550117 (MAF~7%) and SLE in a European cohort¹⁸. Our analysis did not
208 replicate this finding ($P=0.23$).

209

210 ***Underfunctioning of C1QTNF4 p.His198Gln***

211 The candidate gene *C1QTNF4* is a very small gene with <1Kb coding sequence over two
212 exons and is one of four genes constrained against missense variants (ExAC gene-level
213 constraints $Z=3.17$, **Table 2**). Although gene coding length does not correlate with missense
214 constraint scores⁵, it may contribute to insufficient power to detect rare variant associations
215 when using imputed data derived from reference haplotypes from healthy individuals alone.

216 The *de novo* mutation in *C1QTNF4* generates a His198Gln protein sequence change with a
217 modest CADD score of 12.3 (**Table 1**). Although they are useful in the absence of suitable
218 functional assays, the sensitivity of bioinformatic prediction tools is known to be suboptimal.
219 Where functional assays are available, previous studies have also demonstrated functional
220 effects of variants predicted to be tolerated/benign¹⁹.

221 We therefore pursued a functional analysis of the His198Gln *de novo* mutation detected in
222 the *C1QTNF4* gene. Although its function is rather poorly understood, the protein product,
223 C1QTNF4 (CTRP4) is secreted and may act as a cytokine, as it has homology with TNF and
224 the complement component C1q (**Fig. 2**). C1QTNF4 has been shown to influence NF- κ B
225 activation²⁰, a pathway known to be implicated in SLE pathogenesis. In order to study the
226 effect of the *C1QTNF4* mutation, we looked for an effect on NF- κ B production. Using a
227 HEK293-NF- κ B reporter cell line, we showed that C1QTNF4 p.His198Gln mutant protein
228 was expressed and that it inhibited the NF- κ B activation generated by exposure to TNF α
229 (**Fig. 2**). Furthermore, we showed that the fibroblast L929 cell line, which is sensitive to TNF-
230 induced cell death, was rescued by exposure to C1QTNF4 p.His198Gln, but not by wild type
231 C1QTNF4. Thus, the mutant form of C1QTNF4 appears to inhibit some of the actions of
232 TNF, which may promote antinuclear autoimmunity²¹⁻²³.

233

234

235

236

237

238

239

240

241

242

243

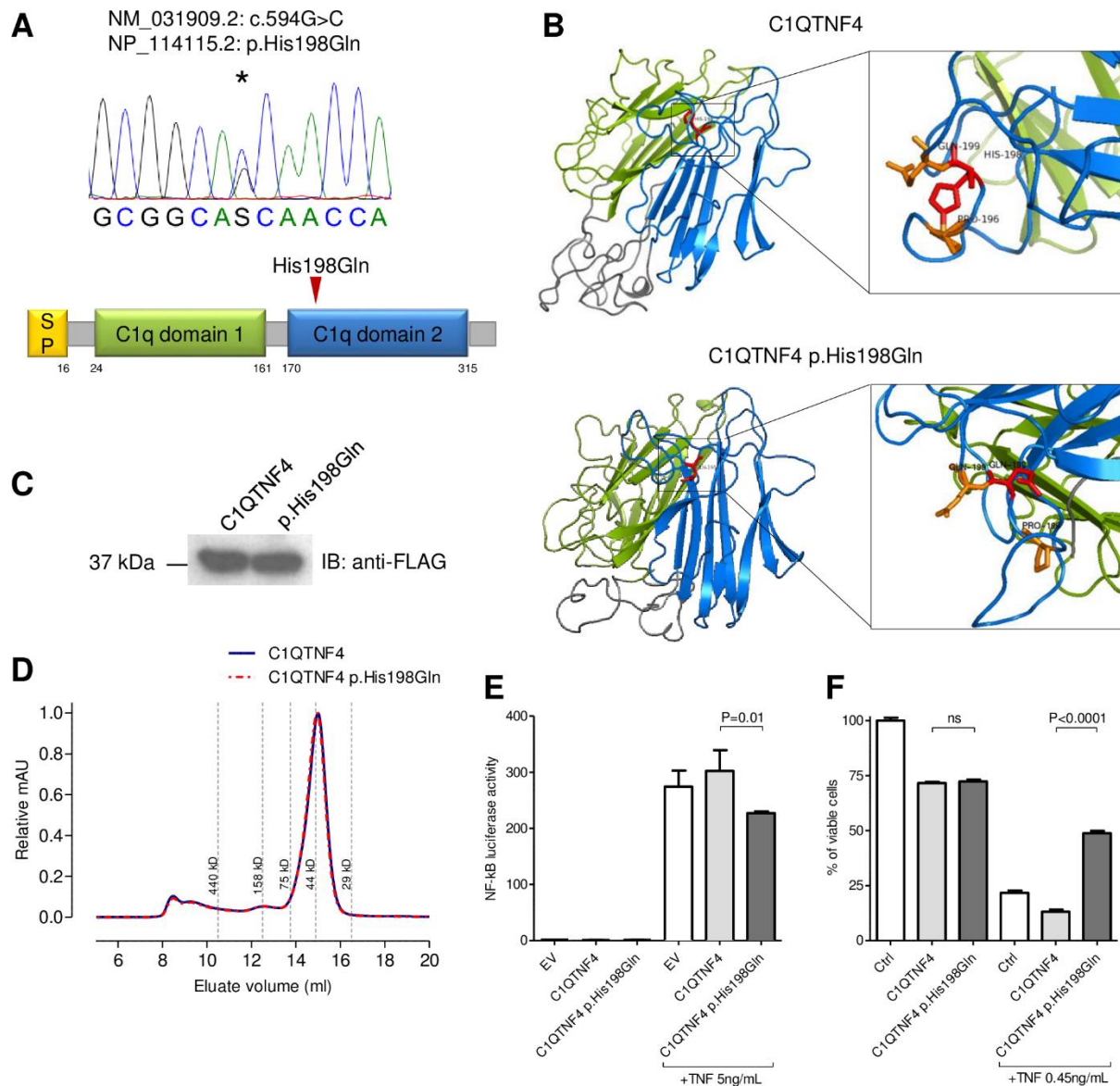
244

245

246

247 **Figure 2. Structural and functional characterization of C1QTNF4 p.His198Gln**
248 **substitution**

249 (A) Domain organization of human C1QTNF4, showing signal peptide (yellow), first C1q
250 domain (green), second C1q domain (blue) and linker peptides (grey). Arrow highlights
251 substitution site. (B) 3D structure prediction of C1QTNF4 and C1QTNF4 p.His198Gln.
252 Ribbons show the interaction between the positively charged Histidine 198 and Proline 196
253 lost in C1QTNF4 p.His198Gln due to the substitution of Histidine with Glutamine. (C)
254 Immunoblot demonstrating that p.His198Gln does not affect secretion of C1QTNF4 in
255 HEK293 supernatants. (D) Size exclusion chromatography profile showing no difference in
256 oligomerisation between supernatant containing C1QTNF4 (blue) and C1QTNF4
257 p.His198Gln (red). (E) Luciferase assay in HEK293-NF- κ B reporter cell line showing that
258 C1QTNF4 p.His198Gln inhibits NF- κ B activation in response to 4h stimulation with 5ng/mL
259 TNF α . Error bars represent standard error of the mean. (F) Inhibition of L929 induced cell
260 death by C1QTNF4 p.His198Gln after 24h of stimulation with 0.45 ng/mL TNF α in presence
261 of Actinomycin 1 μ g/ml.



262

263

264

265 Discussion

266 Following the unexplained heritability left in the wake of massive GWAS in most complex
 267 diseases, searching the GWAS-identified canonical disease susceptibility genes for rare
 268 variants has added little to the heritability explained. Although there are examples – and
 269 perhaps more to discover – of canonical disease genes harbouring both common and rare
 270 risk alleles, the vast majority of such loci do not⁶. Indeed the common-variant associated loci
 271 which have also been shown to harbor rare coding variant risk are often those loci where the

272 common polymorphisms are non-silent coding variants (e.g. *NCF2*⁸). Our data suggest that
273 rare genetic risk may be found in a discrete set of non-canonical susceptibility genes, as we
274 report an association of collective rare variation at *PRKCD* and *DNMT3A*, and found no
275 evidence of an association with common variants across these loci. This, to the best of our
276 knowledge, is the first WES study in polygenic cases of autoimmune disease to use *de novo*
277 mutation discovery to identify candidate genes for rare variant analyses.

278 We saw a relatively high drop-out rate (11.7%) from our NGS to Sanger sequencing
279 confirmed mutations due to the two variants – *KRTAP10-2* and *KLRC1* - found in members
280 of highly homologous gene families. This suggests our NGS error-prone genes (NEPG) filter,
281 which removes loci known to be problematic for genome mapping during NGS analyses,
282 should have been more conservative.

283 *DNMT3A* and *PRKCD*, although hitherto not associated with polygenic SLE, are known
284 autoimmunity susceptibility loci; *DNMT3A* is associated with Crohn's disease (CD)²⁴ and
285 *PRKCD* is associated with both CD and ulcerative colitis (UC)²⁵. The notion that a locus
286 could harbour common variants contributing to one autoimmune disease and rare variants
287 contributing to another is intriguing, and could provide hypothesis-driven searches in the
288 hunt for 'missing heritability'.

289 A study by Berlot et al. identified a functional missense variant p.G510S (c.G1528A) in
290 *PRKCD* in a consanguineous family with monogenic SLE²⁶. It was demonstrated that the
291 *PRKCD*-encoded protein, PRC δ , was essential in the regulation of B cell tolerance and
292 affected family members with the homozygous mutation had increased numbers of immature
293 B cells. Our study implicates the role of rare variants in *PRKCD* in the broader context of
294 SLE susceptibility, beyond a monogenic recessive disease model. Furthermore, *PRKCB*,
295 another member of the protein kinase C gene family, has been implicated in SLE risk in a
296 Chinese study²⁷.

297 *DNMT3A*, a DNA methyltransferase, is a very intriguing candidate gene for SLE as altered
298 patterns of DNA methylation are reported in autoimmune diseases²⁸, and hypomethylation of
299 apoptotic DNA has been reported to induce autoantibody production in SLE²⁹. DNA
300 methylation changes are also associated with monozygotic twin discordance in SLE³⁰.
301 Although previously implicated through a candidate gene study¹⁸, we found no evidence of a
302 common variant association at this locus. Instead we find an association of collective rare
303 variants and SLE sub-phenotypes.

304 Together with the results of Hunt et al., in which autoimmune cases were aggregated and
305 rare variants were found to play a negligible role⁸, our study suggests the importance of
306 deep phenotyping and the potential role of rare variants in specific sub-phenotype, or indeed
307 autoimmune, manifestations. Despite progress with diagnosis and treatment, particular SLE
308 sub-phenotypes – including those used in this study - are still associated with reduced life
309 expectancy. Therefore, elucidating the specific underlying genetic risk is of paramount
310 importance.

311 Our rare variant study is likely to be underpowered given the burden testing was conducted
312 on imputation data from an, albeit large, reference panel of healthy individuals and therefore
313 will not include putative SLE-specific variants. Therefore, other loci with suggestive p-values
314 at candidate genes (e.g. *HMSD*) warrant follow-up analyses through sequencing and/or
315 SLE-weighted imputation reference panels.

316 Through two in vitro assays, we demonstrated the functional effect of a *de novo* variant,
317 His198Gln in *C1QTNF4*, despite this variant being predicted to be of little functional
318 importance across multiple prediction tools. We showed the mutated protein product of
319 *C1QTNF4*, C1QTNF4, inhibits some TNF-mediated cellular responses, including activation
320 of NF-κB and TNF-induced apoptosis. The role of TNF in SLE is complex and incompletely
321 understood, although, in this context, it is noteworthy that TNF inhibition may promote
322 antinuclear autoimmunity²². A recent transancestral ImmunoChip analysis found an
323 association at the chromosome 5-encoded *C1QTNF2* with SLE (Nature Communications, in

324 press), and chromosome 22-encoded *C1QTNF6* is a known susceptibility locus for Type 1
325 Diabetes and is implicated in Rheumatoid Arthritis^{31,32}. Together these data suggest a
326 potential role of the hitherto understudied *C1QTNF* superfamily of genes in autoimmunity.
327 Each human - regardless of the disease status - is estimated to have one *de novo* mutation
328 in their exome¹¹. The simple presence of a provisionally functional *de novo* mutation in a
329 proband is therefore not sufficient evidence that it contributes to disease risk. A major
330 challenge of WES studies, therefore, is how to differentiate between variants truly important
331 to disease and background variation³³. In light of recent studies which have demonstrated
332 the limitations of large-scale exome-wide case-control studies in detecting rare variant
333 contributions^{9,34}, our results support extreme-phenotype sampling and *de novo* mutation
334 discovery to aid a hypothesis-driven search for rare variation contributing to the heritability of
335 complex diseases.

336

337 **Methods**

338 *Selection of trios for sequencing*

339 SLE patients of European ancestry were selected from the UK SLE genetic repository
340 assembled in the Vyse laboratory. The study cases have been subject to genome-wide
341 genotyping as part of a GWAS¹³. The criteria for inclusion were as follows: age of onset of
342 SLE < 25 years (median age 21 years); more marked disease phenotype as shown by either
343 evidence for renal involvement as per standard classification criteria and/or the presence of
344 hypocomplementemia and anti-dsDNA autoantibodies; and DNA available from both
345 unaffected parents. Thirty trios (90 individuals) were studied by WES. Ethical approval for
346 the research was granted by the NRES Committee London (12/LO/1273 and 06/MRE02/9).

347

348

349 *Sequencing and alignment*

350 Sequencing libraries were prepared from 1ug of DNA using the SureSelect XT Human All
351 exon v4 +UTR kit (Agilent Technologies). The libraries were prepared according to the
352 manufacturer's manual (SureSelect XT target enrichment system for Illumina paired end
353 sequencing library, v1.4.1 Sept2012).The exome capture libraries from each individual in the
354 trio were sequenced on Illumina HiSeq 2000 as 100bp paired-end reads. The resulting BCL
355 files from the sequencer were processed with Illumina Casava software v1.8 to obtain paired
356 end reads in FastQ format. The paired end FastQ reads were aligned to the human
357 reference genome hg19 (GCRh37) using Novoalign v2.07.11 with the following parameters
358 (-i 200 30 -o SAM -o SoftClip -k -a -g 65 -x 7). The resulting BAM file was processed to sort
359 and remove PCR duplicates using Picard tools. Only reads uniquely aligned to the reference
360 genome were considered for further analysis. At the end of this process one BAM file per
361 individual was obtained totaling 90 BAM files (30 trios).

362

363 *Quality control (QC) and variant annotation*

364 Variants were retained if they passed the following criteria: (i) read depth $\geq 20x$, (ii) located
365 within exome-captured regions as annotated in Gencode (on-target) and NCBI RefSeq
366 annotation³⁵. Read depth estimation was performed using DepthOfCoverage from GATK
367 tool³⁶, as evaluated using bedtools³⁷. All but one family had 75% representation of the
368 exome at 20X (**Supplementary Fig. 1**). Variants passing these QC criteria were annotated
369 using ANNOVAR³⁸. Any variant observed in either ExAC, 1000 genomes, ESP6500, or in-
370 house databases was considered polymorphic.

371

372

373

374 *De novo variant calling*

375 To screen for *de novo* genetic variants in the affected offspring, three different bioinformatics
376 tools were used: BCFtools³⁹, DeNovoGear⁴⁰ and DeNovoCheck⁴¹. They are based on the
377 SAMtools algorithm for genotype calling and call *de novo* variants in the data from parents
378 and offspring in each family trio, and have previously been used in *de novo* variant
379 identification studies^{41,42}. For each identified *de novo* variant, BCFtools assigns a combined
380 likelihood ratio score (CLR, range 1-255) and DeNovoGear assigns a posterior probability
381 score (PP_dnm, range 0.0 – 1.0), while DeNovoCheck flags a variant as ‘denovo’ without
382 providing a score. Conservative threshold scores for ascertainment of *de novo* variation
383 were applied: CLR \geq 80 for BCFtools and PP_dnm \geq 0.8 for DeNovoGear (**Supplementary**
384 **Fig. 5**). 454 variants were identified at these thresholds. Eight additional variants that were
385 identified by DeNovoCheck and validated by IGV, resulting in a total of 462 variants, which
386 map to 257 genes. The variants were next filtered in the following sequential steps
387 (**Supplementary Figure 2**):

- 388 1. Removal of NGS error prone genes (NEPG): genes previously reported as
389 probable false positive signals from NGS studies due to high frequency of
390 rearrangements, polymorphisms or present in multiple copies⁴³
- 391 2. Fulfil a *de novo* pattern of inheritance: any variant that supports Het:Ref:Ref for
392 Child:Father:Mother, respectively, was considered a potential *de novo* variant.
393 We also further selected variants that did not contain any trace of alternate allele
394 in any of the parents. IGVtools was used to count the number of non-reference
395 bases at each identified variant position from both father and mother. Only
396 variants with a zero count of non-reference bases in both father and mother were
397 considered as very high quality variants and retained for further analysis.
- 398 3. Variant annotation: only non-silent variants (Missense, Nonsense, splicing and
399 insertions/deletions) was retained for further analysis.

400 This process resulted in a total of 17 variants in 17 genes (**Supplementary Table 1**).

401 *Analysis of whole exome sequencing (WES) in cases only*

402 To quantify the advantage of using parent-offspring trios, the whole exome samples of the
403 30 probands were analysed alone. 584,798 variants with $\geq 20X$ coverage depth and within
404 Gencode capture regions were identified. All stringent filters to help in refinement of variants
405 that could be potentially causal were applied. These include two filters previously described
406 in the *de novo* filtering approach (Non-NEPG and variant annotation) along with selection of
407 heterozygous variants (based on the pattern of zygosity expected for *de novo* mutations)
408 and non-polymorphic filters (variants not observed in control datasets). To note, this is unlike
409 the trio analysis where variants were not filtered for non-polymorphic. Filters were applied
410 sequentially 1) Non-NEPG 2) Variant annotation 3) Heterozygosity 4) Non-polymorphic and
411 resulted in 1194 variants in 1067 genes (**Supplementary Figure 3**).

412

413 *Sanger Sequencing confirmation*

414 Primers were designed using Primer 3 to target the exon containing the *de novo* mutation.
415 Primers and PCR conditions available on request. 10ng of DNA from SLE probands, any
416 unaffected siblings and both parents was amplified with Hot Start Taq polymerase on a G-
417 Storm Thermocycler. PCR products were first cleaned with EXO-SAP before BigDye
418 labelling in a linear PCR. Samples were sequenced on an ABI 3300XL.

419

420 *Imputation*

421 Genotype data from 10,995 individuals of matched European ancestry, including 4,036 SLE
422 cases, genotyped on the Illumina Chip Illumina HumanOmni1 BeadChip for a previous
423 study¹³ were used. These data had undergone quality control as previously described¹³,
424 including Principal Component Analysis (PCA) to account for population structure. The
425 UK10K (REL-2012-06-02) plus 1000 Genomes Project Phase3 data (release 20131101.v5)

426 merged reference panel (UK10K-1000GP3) was accessed through the European Genome-
427 phenome Archive (EGAD00001000776). The genotype data were imputed using the UK10K-
428 1000GP3 reference panel across the coding regions of the 14 genes with *de novo* mutations
429 plus a 2Mb flanking region. To increase the accuracy of imputed genotype calls, a full
430 imputation without pre-phasing was conducted using IMPUTE2^{44,45}. Imputed genotypes were
431 filtered for confidence using an info score (IMPUTE2) threshold of 0.3 (**Supplementary**
432 **Figure 6**). The most likely genotype from IMPUTE2 was taken if its probability was > 0.5. If
433 the probability fell below this threshold, it was set as missing. Variants with >10% missing
434 genotype calls were removed for further analysis. All individuals had <8% missing genotype
435 data.

436

437 *Rare Variant Burden Tests*

438 Imputed data were filtered, using Plink v1.9, to include only variants mapping to coding
439 exons of hg19 RefSeq transcripts. Plink/SEQv1.0¹⁶ was used to run gene-wise burden
440 testing with a MAF<1% threshold. A 5% false discovery rate was used for multiple testing
441 correction.

442

443 *Common variant association tests*

444 Imputed data were filtered to include variants with MAF>1% and SNPTEST 2.5.2⁴⁶ was
445 used to test for associations across the region spanning the encoded gene. Bonferroni
446 correction was used for 3,000 tests across the loci ($q=1.66E-5$).

447

448 *Plasmids*

449 Myc-Flag-tagged *C1QTNF4* on the pCMV6 vector and the empty pCMV6 vector were
450 purchased from OriGene. The mutant pCMV6-*C1QTNF4 C594G* (p.His198Gln) was

451 generated by site-directed mutagenesis (Quikchange II XL; Stratagene) according the
452 manufacturer instructions: mutagenic primer: 5'-GCGAGTGGTTGCTGCCGCGGCC-3'
453 (Sigma Aldrich). The plasmids production was carried on in XL10-Gold Ultracompetent cells,
454 isolated and purified using EndoFree Maxi Prep kit (Qiagen). All the plasmid ORFs were
455 confirmed by full Sanger sequencing (GATC-Biotech). The expression and secretion of the
456 flagged proteins was confirmed by western blot on cell lysates and supernatants with
457 monoclonal anti-FLAG antibody (clone M2; Sigma-Aldrich).

458

459 *Luciferase assays and TNF-induced programmed cell death*

460 GloResponse NF- κ B-RE-luc2P HEK293 cell line (Promega) and TNF-sensitive L929
461 fibrosarcoma cell line (ATCC) were cultured in Dulbecco's Modified Eagle Medium (DMEM)
462 enriched with 10% fetal bovine serum (FBS) and 1% Penicillin/Streptomycin (complete
463 DMEM) at 37°C, 5% CO₂. HEK293 were seeded 24 hours before transfection in antibiotic
464 free DMEM in 96 wells plate (2×10^4 cells/well), transfected with either *C1QTNF4*, *C1QTNF4*
465 *C594G* or Empty Vector via Fugene HD (Promega). 48 hours after transfection the cell were
466 left unstimulated or stimulated with TNF α 5 ng/ml (PeproTech) for 4 hours. Luciferase
467 activity was assayed by One-Glo (Promega) on Berthold Orion luminometer, the values were
468 normalized to cell viability measured by CellTiter Glo (Promega). L929 were challenged with
469 TNF α 0.45 ng/ml and Actinomycin D 1 μ g/ml (R&D) for 24 hours in presence of *C1QTNF4* or
470 *C1QTNF4* p.His198Gln containing media, cell viability was measured by CellTiter Glo.

471

472 *Size exclusion chromatography*

473 Supernatants (750 μ l) of HEK293 producing *C1QTNF4* or *C1QTNF4* p.His198Gln were
474 buffer exchanged in PBS on Zeba Spin Desalting Columns (Thermo Fisher) and 0.5 mL
475 loaded on an AKTA FPLC with a Superdex 200 10/300 GL column (GE Healthcare).
476 Absorbance was normalized to the maximum peak of each sample.

477 *In silico protein structure prediction*

478 The web-based service Protein Homology/AnalogY Recognition Engine (Phyre2)⁴⁷ was
479 used for the protein structure prediction of C1QTNF4 and C1QTNF4 p.His198Gln. The PDB
480 file produced was loaded on Pymol software for visualization (Schrödinger, LLC).

481

482 *Data availability*

483 WES data on 90 individuals – 30 parent-offspring trios – will be deposited at the European
484 Genome-phenome Archive.

485

486 **Acknowledgements**

487 The work leading to these results received funding from the European Union FP7
488 programme (grant agreement n° 262055) via the European Sequencing and Genotyping
489 Infrastructure (ESGI). Sequencing was performed by the SNP&SEQ Technology Platform in
490 Uppsala, which is part of the National Genomics Infrastructure (NGI) hosted by Science for
491 Life Laboratory in Sweden. This work was supported in part by the Swedish Research
492 Council for Medicine and Health (grant n° E0226301) and by the Knut and Alice Wallenberg
493 Foundation (KAW 2011.0073). We thank Johanna Lagensjö and Olof Karlberg for assistance
494 with sequencing.

495 The research was funded/supported by the National Institute for Health Research (NIHR)
496 Biomedical Research Centre based at Guy's and St Thomas' NHS Foundation Trust and
497 King's College London

1. Bamshad, M. J. *et al.* Exome sequencing as a tool for Mendelian disease gene discovery. *Nat. Rev. Genet.* **12**, 745–755 (2011).
2. Chong, J. X. *et al.* The Genetic Basis of Mendelian Phenotypes: Discoveries, Challenges, and Opportunities. *Am. J. Hum. Genet.* **97**, 199–215 (2015).
3. Manolio, T. A. *et al.* Finding the missing heritability of complex diseases. *Nature* **461**, 747–753 (2009).
4. Rivas, M. a *et al.* Deep resequencing of GWAS loci identifies independent rare variants associated with inflammatory bowel disease. *Nat. Genet.* **43**, 1066–73 (2011).
5. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291 (2016).
6. Jordan, C. T. *et al.* Rare and common variants in CARD14, encoding an epidermal regulator of NF-kappaB, in psoriasis. *Am. J. Hum. Genet.* **90**, 796–808 (2012).
7. Diogo, D. *et al.* Rare, low-frequency, and common variants in the protein-coding sequence of biological candidate genes from GWASs contribute to risk of rheumatoid arthritis. *Am. J. Hum. Genet.* **92**, 15–27 (2013).
8. Hunt, K. a *et al.* Negligible impact of rare autoimmune-locus coding-region variants on missing heritability. *Nature* **498**, 232–5 (2013).
9. Fuchsberger, C. *et al.* The genetic architecture of type 2 diabetes. *Nature* **536**, 41–47 (2016).
10. Moutsianas, L. *et al.* The power of gene-based rare variant methods to detect disease-associated variation and test hypotheses about complex disease. *PLoS Genet.* **11**, e1005165 (2015).
11. Veltman, J. a. & Brunner, H. G. De novo mutations in human genetic disease. *Nat. Rev. Genet.* **13**, 565–575 (2012).
12. Rahbari, R. *et al.* Timing, rates and spectra of human germline mutation. *Nat. Genet.* **48**, 1–11 (2015).
13. Bentham, J. *et al.* Genetic association analyses implicate aberrant regulation of innate and adaptive immunity genes in the pathogenesis of systemic lupus erythematosus. *Nat. Genet.* **47**, 1457–1464 (2015).
14. Morris, D. L. *et al.* Genome-wide association meta-analysis in Chinese and European individuals identifies ten new loci associated with systemic lupus erythematosus. *Nat. Genet.* **48**, (2016).
15. Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of human genetic variants. *Nat. Genet.* **46**, 310–315 (2014).
16. Purcell, S. M. PLINK/SEQ: A library for the analysis of genetic variation data. at <<https://atgu.mgh.harvard.edu/plinkseq/>>
17. Okada, Y. *et al.* Integration of sequence data from a consanguineous family with genetic data from an outbred population identifies PLB1 as a candidate rheumatoid arthritis risk gene. *PLoS One* **9**, 1–12 (2014).
18. Piotrowski, P., Grobelna, M. K., Wudarski, M., Olesinska, M. & Jagodzinski, P. P. Genetic variants of DNMT3A and systemic lupus erythematosus susceptibility. *Mod Rheumatol* **25**, 96–99 (2015).

19. Roberts, A. L. *et al.* Resequencing the susceptibility gene, ITGAM, identifies two functionally deleterious rare variants in systemic lupus erythematosus cases. *Arthritis Res. Ther.* **16**, R114 (2014).
20. Li, Q. *et al.* Identification of C1qTNF-related protein 4 as a potential cytokine that stimulates the STAT3 and NF- κ B pathways and promotes cell survival in human cancer cells. *Cancer Lett.* **308**, 203–214 (2011).
21. Beigel, F. *et al.* Formation of antinuclear and double-strand DNA antibodies and frequency of lupus-like syndrome in anti-TNF- α antibody-treated patients with inflammatory bowel disease. *Inflamm. Bowel Dis.* **17**, 91–98 (2011).
22. Eriksson, C., Engstrand, S., Sundqvist, K.-G. & Rantapää-Dahlqvist, S. Autoantibody formation in patients with rheumatoid arthritis treated with anti-TNF alpha. *Ann. Rheum. Dis.* **64**, 403–7 (2005).
23. Pink, A. E., Fonia, A., Allen, M. H., Smith, C. H. & Barker, J. N. W. N. Antinuclear antibodies associate with loss of response to antitumour necrosis factor-alpha therapy in psoriasis: a retrospective, observational study. *Br. J. Dermatol.* **162**, 780–5 (2010).
24. Franke, A. *et al.* Genome-wide meta-analysis increases to 71 the number of confirmed Crohn's disease susceptibility loci. *Nat. Genet.* **42**, 1118–25 (2010).
25. Jostins, L. *et al.* Host-microbe interactions have shaped the genetic architecture of inflammatory bowel disease. *Nature* **491**, 119–24 (2012).
26. Belot, A. *et al.* Protein kinase C δ deficiency causes mendelian systemic lupus erythematosus with B cell-defective apoptosis and hyperproliferation. *Arthritis Rheum.* **65**, 2161–2171 (2013).
27. Sheng, Y.-J. *et al.* Follow-up study identifies two novel susceptibility loci PRKCB and 8p11.21 for systemic lupus erythematosus. *Rheumatology (Oxford)*. **50**, 682–688 (2011).
28. Ballestar, E. Epigenetic alterations in autoimmune rheumatic diseases. *Nat. Rev. Rheumatol.* **7**, 263–71 (2011).
29. Wen, Z. K. *et al.* DNA hypomethylation is crucial for apoptotic DNA to induce systemic lupus erythematosus-like autoimmune disease in SLE-non-susceptible mice. *Rheumatology* **46**, 1796–803 (2007).
30. Javierre, B. M. *et al.* Changes in the pattern of DNA methylation associate with twin discordance in systemic lupus erythematosus. *Genome Res.* **20**, 170–9 (2010).
31. Onengut-Gumuscu, S. *et al.* Fine mapping of type 1 diabetes susceptibility loci and evidence for colocalization of causal variants with lymphoid gene enhancers. *Nat. Genet.* **47**, 381–6 (2015).
32. Murayama, M. A. *et al.* CTRP6 is an endogenous complement regulator that can effectively treat induced arthritis. *Nat Commun* **6**, 8483 (2015).
33. Cooper, G. M. & Shendure, J. Needles in stacks of needles: finding disease-causal variants in a wealth of genomic data. *Nat. Rev. Genet.* **12**, 628–40 (2011).
34. Luo, Y. *et al.* Exploring the genetic architecture of inflammatory bowel disease by whole genome sequencing identifies association at ADCY7. *Nat. Genet.* **49**, 186–192 (2016).
35. Frankish, A. *et al.* Comparison of GENCODE and RefSeq gene annotation and the impact of reference geneset on variant effect prediction. *BMC Genomics* **16(Suppl 8)**, S2 (2015).

36. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491–8 (2011).
37. Quinlan, A. R. & Hall, I. M. BEDTools: A flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).
38. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic Acids Res.* **38**, e164 (2010).
39. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
40. Ramu, A. *et al.* DeNovoGear: de novo indel and point mutation discovery and phasing. *Nat. Methods* **10**, 985–987 (2013).
41. de Ligt, J. *et al.* Diagnostic exome sequencing in persons with severe intellectual disability. *N. Engl. J. Med.* **367**, 1921–9 (2012).
42. Nava, C. *et al.* De novo mutations in HCN1 cause early infantile epileptic encephalopathy. *Nat. Genet.* **46**, 640–5 (2014).
43. Fuentes Fajardo, K. V. *et al.* Detecting false-positive signals in exome sequencing. *Hum. Mutat.* **33**, 609–613 (2012).
44. Howie, B., Fuchsberger, C., Stephens, M., Marchini, J. & Abecasis, G. R. Fast and accurate genotype imputation in genome-wide association studies through pre-phasing. *Nat. Genet.* **44**, 955–959 (2012).
45. Roshyara, N. R. *et al.* Comparing performance of modern genotype imputation methods in different ethnicities. *Sci. Rep.* **6**, 34386 (2016).
46. Marchini, J. & Howie, B. Genotype imputation for genome-wide association studies. *Nat. Rev. Genet.* **11**, 499–511 (2010).
47. Kelley, L. A. & Sternberg, M. J. E. Protein structure prediction on the Web: a case study using the Phyre server. *Nat. Protoc.* **4**, 363–371 (2009).

Table 1: *De novo* mutations in SLE probands with extreme phenotypes

Family	Chr	Position (hg19)	Ref	Alt	Gene	Gene Description	Exon	Amino acid	MAF in ExAC ^a	PhyloP	CADD Phred	Mutation Type ^b
SLE0751	22	38336799	C	T	<i>MICALL1</i>	MICAL-like 1	16	Arg852Cys	1.5E-04	2.31	35	Ti CpG
SLE0496	3	53223122	G	A	<i>PRKCD</i>	protein kinase C, delta	16	Gly535Arg	-	2.85	34	Ti CpG
SLE0679	12	57588368	C	T	<i>LRP1</i>	Low-density lipoprotein receptor-related protein 1	50	Arg2693Cys	8.3E-06	1.34	34	Ti CpG
SLE0592	6	36260896	G	A	<i>PNPLA1</i>	patatin-like phospholipase domain containing 1	3	Arg166His	5.8E-05	2.56	33	Ti CpG
SLE0296	2	25457236	G	A	<i>DNMT3A</i>	DNA (cytosine-5-)-methyltransferase 3 alpha	19	Ala695Val	-	2.75	32	Ti CpG
SLE0571	4	79512728	G	T	<i>ANXA3</i>	annexin A3	7	Ser145Ile	-	2.56	25.2	Tv
SLE0679	3	171431716	G	A	<i>PLD1</i>	phospholipase D1, phosphatidylcholine-specific	9	Thr293Met	5.8E-05	2.68	25.1	Ti CpG
SLE0411	5	179743769	C	T	<i>GFPT2</i>	glutamine-fructose-6-phosphate transaminase 2	12	Val383Met	2.6E-05	1.48	23.4	Ti CpG
SLE0679	7	138968784	C	A	<i>UBN2</i>	ubiquitin 2	15	Pro1045Thr	-	2.75	18.46	Tv
SLE0080	16	2812426	C	T	<i>SRRM2</i>	serine/arginine repetitive matrix 2	11	Arg633Cys	-	0.77	14.32	Ti CpG
SLE0852	11	47611769	G	C	<i>C1QTNF4</i>	C1q and tumor necrosis factor related protein 4	2	His198Gln	-	2.09	12.29	Tv
SLE0321	18	61621642	G	A	<i>HMSD</i>	histocompatibility (minor) serpin domain containing	3	Ala25Thr	-	1.26	9.732	Ti
SLE0390	12	32369376	G	C	<i>BICD1</i>	bicaudal D homolog 1 (Drosophila)	2	Val137Leu	-	0.72	8.673	Tv
SLE0321	1	35251125	C	G	<i>GJB3</i>	gap junction protein, beta 3	2	Asp254Glu	-	-0.25	0.002	Tv

The mutations are ordered by level of severity, from most to least, predicted by CADD score

^aFrequencies are presented from all 61,468 multiethnic individuals in ExAC because the *de novo* mutations observed in ExAC are likely to be identity-by-state not identity-by-descent.

^bTv = Transversion; Ti = Transition; Ti CpG = Transition within a CpG dinucleotide

Table 2: Evidence for role of *de novo* mutation gene in autoimmunity

Gene	Functional Candidate ^a	Association with SLE ^b	Associations with other AID ^b	Immune cell type with highest expression ^c	Missense Constraint ^d
<i>PRKCD</i>	B cell signaling and self-antigen induced B cell tolerance induction	Monogenic forms ²⁶	IBD, UC, CD ²⁴	Dendritic	3.75
<i>DNMT3A</i>	DNA methyltransferase	Candidate gene study ¹⁸	CD ²⁵	-	4.31
<i>C1QTNF4</i>	Pro-inflammatory cytokine	-	-	CD34+	3.17
<i>SRRM2</i>	Spliceosome-associated pre-mRNA splicing	-	-	CD8+	No data
<i>LRP1</i>	Endo/Phagocytosis of apoptotic cells	-	-	-	10.60
<i>HMSD</i>	Minor histocompatibility antigen	-	-	n/a	0.25
<i>UBN2</i>	DNA binding	-	-	-	0.01
<i>ANXA3</i>	-	-	RA ¹⁷	-	-0.37
<i>PLD1</i>	-	-	-	Lymphoblasts	-0.73
<i>PNPLA1</i>	-	-	-	-	0.27
<i>GFTP2</i>	-	-	-	-	1.59
<i>BICD1</i>	-	-	-	-	2.12
<i>GJB3</i>	-	-	-	-	-0.81
<i>MICALL1</i>	-	-	-	-	0.50

Genes appear in descending order of supporting evidence. UC=ulcerative colitis, CD=Crohn's Disease, IBD=inflammatory bowel disease, RA=Rheumatoid Arthritis

^a See Supplementary Table 4

^b See Supplementary Table 5

^c See Supplementary Figure 4. Data from BioGPS. If gene expression is highest in immune cells compared to all other cells, the immune cell type with highest expression is listed.

^d Gene-wise ExAC Constraint Z-scores. Genes with significant restraint against missense variants are highlighted in bold.

Table 3: Gene-based rare variant burden analyses

Locus	# variants	# minor alleles controls	SLE		Anti-dsDNA		Renal with hypocomplementemia	
			# minor alleles cases	p-value	# minor alleles cases	p-value	# minor alleles cases	p-value
<i>LRP1</i>	84	927	514	1.00	143	1.00	25	0.57
<i>BICD1</i>	68	673	397	0.38	147	0.50	22	0.45
<i>UBN2</i>	63	338	214	0.26	77	0.07	17	0.11
<i>PLD1</i>	55	910	530	0.60	153	0.24	23	1.00
<i>SRRM2</i>	37	380	188	1.00	57	1.00	12	0.29
<i>MICALL1</i>	29	146	75	1.00	16	1.00	6	0.22
<i>GFTP2</i>	25	350	212	0.35	72	0.26	13	0.14
<i>DNMT3A</i>	24	110	91	0.0075	38	0.0005	9	0.0033
<i>PRKCD</i>	13	69	69	0.0028	20	0.06	2	0.71
<i>ANXA3</i>	12	311	155	1.00	47	1.00	4	1.00
<i>GJB3</i>	11	145	71	1.00	19	1.00	1	1.00
<i>PNPLA1</i>	11	194	139	0.04	51	0.02	3	1.00
<i>C1QTNF4</i>	9	186	100	1.00	31	1.00	4	1.00
<i>HMSD</i>	5	10	8	0.71	6	0.03	0	1.00
Grouped	446	4749	2763	0.57	877	0.31	141	0.23

Individual genes are ordered by descending number of observed rare variants. Significant p-values (burden test) at 5% FDR are highlighted in bold ($q=0.003$)