

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

3

4 Venu Pullabhatla^{a,1}, Amy L. Roberts^{b,1}, Myles J. Lewis^c, Daniele Mauro^c, David L. Morris^b,
5 Christopher A. Odhams^b, Philip Tomblason^b, Ulrika Liljedahl^d, Simon Vyse^{b,9}, Michael A.
6 Simpson^b, Sascha Sauer^{e,h}, Emanuele de Rinaldis^a, Ann-Christine Syvänen^d, Timothy J.
7 Vyse^{b,f,2}

8

9 ^a NIHR GSTFT/KCL Comprehensive Biomedical Research Centre, Guy's & St. Thomas'
10 NHS Foundation Trust, London, SE1 9RT, UK

11

12 ^b Division of Genetics and Molecular Medicine, King's College London, London, SE1 9RT,
13 UK

14

15 ^c Centre for Experimental Medicine and Rheumatology, William Harvey Research Institute,
16 Queen Mary University of London, London, EC1M 6BQ, UK

17

18 ^d Department of Medical Sciences, Uppsala University, Uppsala, 75144, Sweden

19

20 ^e Max Planck Institute for Molecular Genetics, Berlin, 14195, Germany

21

22 ^f Division of Immunology, Infection and Inflammatory Disease, King's College London,
23 London, SE1 9RT, UK

24

25 ⁹ Current address: Division of Cancer Biology, The Institute of Cancer Research, London,
26 SW3 6JB, UK

27

28 ^h Current address: Max Delbrück Centre for Molecular Medicine (BIMSB/BIH), Berlin, 13092,
29 Germany

30

31 ¹ V.P and A.L.R contributed equally to this work

32

33 ² To whom correspondence should be addressed. Email: timothy.vyse@kcl.ac.uk

34

35

36 **Abstract**

37 The omnigenic model of complex diseases stipulates that the majority of the heritability will
38 be explained by the effects of common variation on genes in the periphery of core disease
39 pathways. Rare variant associations, expected to explain far less of the heritability, may be
40 enriched in core disease genes and thus will be instrumental in the understanding of
41 complex disease pathogenesis and their potential therapeutic targets. Here, using
42 complementary whole-exome sequencing (WES), high-density imputation, and *in vitro*
43 cellular assays, we identify three candidate core genes in the pathogenesis of Systemic
44 Lupus Erythematosus (SLE). Using extreme-phenotype sampling, we sequenced the
45 exomes of 30 SLE parent-affected-offspring trios and identified 14 genes with missense *de*
46 *novo* mutations (DNM), none of which are within the >80 SLE susceptibility loci implicated
47 through genome-wide association studies (GWAS). In a follow-up cohort of 10,995
48 individuals of matched European ancestry, we imputed genotype data to the density of the
49 combined UK10K-1000 genomes Phase III reference panel across the 14 candidate genes.
50 We identify a burden of rare variants across *PRKCD* associated with SLE risk ($P=0.0028$),
51 and across *DNMT3A* associated with two severe disease prognosis sub-phenotypes
52 ($P=0.0005$ and $P=0.0033$). Both genes are functional candidates and significantly
53 constrained against missense mutations in gene-level analyses, along with *C1QTNF4*. We
54 further characterise the TNF-dependent functions of candidate gene *C1QTNF4* on NF- κ B
55 activation and apoptosis, which are inhibited by the p.His198Gln DNM. Our results support
56 extreme-phenotype sampling and DNM gene discovery to aid the search for core disease
57 genes implicated through rare variation.

58 **Significance Statement**

59 Rare variants, present in <1% in population, are expected to explain little of the heritability of
60 complex diseases, such as Systemic Lupus Erythematosus (SLE), yet are likely to identify
61 core genes crucial to disease mechanisms. Their rarity, however, limits the power to show

62 their statistical association with disease. Through sequencing the exomes of SLE patients
63 and their parents, we identified non-inherited *de novo* mutations in 14 genes and
64 hypothesised that these are prime candidates for harbouring additional disease-associated
65 rare variants. We demonstrate that two of these genes also carry a significant excess of
66 rare variants in an independent, large cohort of SLE patients. Our findings will influence
67 future study designs in the search for the 'missing heritability' of complex diseases.

68 **/body**

69 **Introduction**

70 Considerable progress has been made in elucidating the genetic basis of complex diseases.
71 The vast majority of identified disease-associated genetic polymorphisms are common in the
72 population and the risk alleles impart a modest individual increment to the likelihood of
73 developing disease. Although large-scale genome-wide association studies (GWAS) have so
74 far explained less of the heritability than originally predicted (1), much of the 'missing
75 heritability' is expected to be accounted for by common variants with effect sizes below the
76 genome-wide significance threshold (2). However, under the newly proposed omnigenic
77 model of complex traits, the majority of associated common variants – both identified and
78 unidentified - will primarily be found in periphery genes expressed in relevant cell types but
79 not necessarily biologically relevant to disease (3).

80 In contrast, the role of rare variants in complex disease is largely unknown and often
81 dismissed. A recent study, however, with an extremely large sample size, identified rare and
82 low frequency variants contributing to the genetic variance of adult human height (4) – a
83 polygenic trait with a genetic architecture similar to that of complex diseases (5) - suggesting
84 previous complex disease studies with seemingly large sample sizes were perhaps still
85 insufficiently powered to detect rare variant associations (6). Furthermore, studies of rare
86 variants typically find gene sets enriched in biologically relevant functions/pathways (3, 7, 8).
87 Therefore, although estimated to explain less of the heritable disease risk at a population

88 level than common variants, identifying rare and low frequency variants is of paramount
89 importance to understanding disease pathogenesis as they are likely to implicate biologically
90 relevant core genes (3). Supporting the theory that common and rare variant associations
91 will be found in discrete gene sets is the lack of additional rare variant associations in GWAS
92 genes (9).

93 Exome-wide searches, which provides a highly enriched source of potential disease-causing
94 mutations (10), have revealed limited numbers of rare variation associated with complex
95 diseases. Even though greater statistical power is achieved by gene-level analyses whereby
96 aggregated variants are tested for an allelic burden of collective rare variation, widely used
97 gene-based association tests have been shown to lack power at the exome-wide level (11).
98 Coupled with the insufficient sample sizes currently available in the study of most complex
99 diseases, hypothesis-free searches for core genes with rare variant associations are unlikely
100 to be fruitful.

101 Our strategy to address this problem in autoimmune disease Systemic Lupus Erythematosus
102 (SLE), is outlined here and summarised in Fig. 1. Using a discovery cohort of 30 unrelated
103 SLE cases with a severe disease (young age of onset and clinical features associated with
104 poorer outcome), we hypothesized that these individuals would exhibit unique mutation
105 events in their protein-coding DNA that may predisposed to disease risk. We undertook
106 whole exome sequencing (WES) in 30 family trios (both parents and affected offspring) and
107 scrutinized the data for non-inherited *de novo* mutations (DNM) in the individual with SLE to
108 identify a group of candidate genes for an independent follow-up rare variant analysis. This
109 method allowed the identification of novel loci harbouring disease risk through collective rare
110 variation, and emphasises the value of phenotypic extremes in the search for core genes in
111 multifactorial disorders (12).

112

113 **Results**

114 **Identification of DNM in extreme-phenotype SLE cases.** We screened for DNM by WES
115 of 30 family trios with an affected offspring with more severe SLE (Fig. S1). A total of
116 584,798 variants ($\geq 20X$), including single nucleotide variants and indels, were identified in
117 the 30 affected probands. Using three bioinformatic tools and employing conservative
118 parameters, 17 putative missense DNM were identified across 17 genes (Table S1; Fig. S2).
119 We also analysed the SLE proband WES data alone, without the unaffected parents. This
120 revealed 1,194 non-silent, heterozygous, rare variants in 1,067 genes distributed across the
121 genome, which would make prioritisation for downstream analysis a difficult task, highlighting
122 the benefit of parent-offspring trio sequencing (Fig. S3). Sanger sequencing confirmed 14
123 true positive non-silent DNM (Table 1; Table S2), present in the SLE proband but absent in
124 both parents and any unaffected siblings, in 11 of the 30 probands (36.7%) for further
125 analysis. No DNM was found in any of the >80 known SLE-associated genes. Of the three
126 false positive DNM (11.7%; Table S1) one, within *LAMC2*, is likely a result of germline
127 mosaicism because, although not observed in either parent, it is observed in an unaffected
128 sibling in addition to the SLE proband (13), and the other two variants are within *KRTAP10-2*
129 and *KLRC1* - both members of highly homologous gene families. Such sequence identity
130 may have caused false positive identification of DNM in the WES analysis and suggests our
131 NGS error-prone genes (NEPG) filter, which removes loci known to be problematic for
132 genome mapping during NGS analyses, should have been more conservative. Indeed the
133 *KLRC1* p.Ile225Met missense variant appears to be a polymorphic Paralogous Sequence
134 Variant (PSV) – the paralogous variant being Met223Ile in *KLRC2*.

135 **Variant- and gene-level functional characterisation of DNM.** In order to best predict the
136 phenotypic effect of the 14 DNM, we used both variant-level and gene-level metrics (14). We
137 used the ExAC database (15) and Combined Annotation Dependent Depletion (CADD)
138 scores (16) to characterise the frequency and predicted functional effects, respectively, of
139 the variants. Five of the 14 DNM – found in *MICALL1*, *LRP1*, *PNPLA1*, *PLD1*, and *GFTP2* -
140 have been observed, at very rare frequencies, in the ~60,000 exomes documented in ExAC

141 (Table 1). All five mutations are CpG transitions and therefore likely to be identity-by-state,
142 reflecting the higher mutability rate of these sites. Within the mutation set, five (35.7%) –
143 found in *DNMT3A*, *PRKCD*, *MICALL1*, *LRP1*, and *PNPLA1* – have CADD Phred scores >30,
144 placing them in the top 0.1% of possible damaging mutations in the human genome (Table
145 1). We further explored the function, expression (BioGPS), existing autoimmunity
146 associations (ImmunoBase), and gene-level constraint against missense mutations (ExAC),
147 of the DNM genes to build a profile of *a priori* evidence of a role in SLE pathogenesis. None
148 of the candidate genes have been previously associated with SLE through GWAS in any
149 population (17). We also identify candidate genes through known/predicted function and
150 expression profiles (*C1QTNF4*, *SRRM2*, *HMSD*), and four genes (*PRKCD*, *DNMT3A*,
151 *C1QTNF4* and *LRP1*) with a significant ($Z>3.09$) constraint against missense variants (Table
152 2). However, across the entire gene set, there was no difference in the median Z-score
153 (0.50) compared with the median Z-score across all genes in ExAC (0.51).

154 ***PRKCD* and *DNMT3A* as novel SLE genes.** Although the variant- and gene-level metric
155 analyses suggested intriguing functional candidates, we took a comprehensive approach
156 and tested each locus for an allelic burden of rare variation. We hypothesised that, while
157 some observed DNM were random background variation as present in the exome of every
158 individual regardless of disease status (18), others may be reflecting a hitherto unknown
159 gene contributing to SLE risk, and this may be shown through rare variant burden.
160 Therefore, genotype data was imputed (Fig. S6 and S7) to the density of the combined
161 UK10K and 1000 genomes Phase III reference panel (UK10K-1000GP3) across all 14 DNM
162 genes in a follow-up cohort of 10,995 individuals of matched European ancestry previously
163 genotyped on the Illumina HumanOmni1 BeadChip (19). Under the hypothesis that rare
164 variants at these loci would be causal and not protective, we employed a one-tailed
165 collapsing burden test (20) to survey each of the 14 genes for an excess of aggregated rare
166 (MAF<1%) exonic variants in SLE cases compared with healthy controls. We identify an
167 association of *PRKCD* rare variants with SLE (Table S3; $P=0.0028$; $n_{\text{cases}}=4,036$). In sub-

168 phenotype analyses, we identify collective rare exonic variants in *DNMT3A* associated with
169 both anti-dsDNA (Table S3; $P=0.0005$; $n_{\text{cases}}=1,261$) and renal involvement with
170 hypocomplementemia (Table S3; $P=0.0033$; $n_{\text{cases}}=186$), both of which are markers of more
171 severe disease. We also collapsed all exons from the 14 genes together to test for an overall
172 burden of rare variants across these loci. These analyses revealed no excess of rare exonic
173 variants across the grouped genes, reflecting the hypothesis that some/most genes will not
174 be relevant to disease status because the observed DNM are random background variation
175 only. These data reflect the results of our gene-level constraint metric, in which the
176 aggregated gene set do not have a significant mutation constraint. Together these results
177 suggest further prioritisation based on gene-level metrics would not have resulted in true
178 positive associations being excluded from analyses.

179 **Effect of DNM p.His198Gln on C1QTNF4 function.** Although no rare variant association
180 was found at the novel candidate gene *C1QTNF4*, it's potential role in disease is supported
181 by gene-level metrics – it is a compelling functional candidate and one of four genes
182 constrained against missense variants (ExAC gene-level constraints $Z=3.17$, Table 2).
183 Although gene coding length does not correlate with missense constraint scores (15), the
184 small (<1Kb) coding sequence of this candidate gene may have contributed to insufficient
185 power to detect a rare variant association in the burden testing. On the variant-level, the
186 DNM in *C1QTNF4* generates a p.His198Gln sequence change with a modest CADD score
187 of 12.3 (Table 1). Although useful in the absence of suitable functional assays, the sensitivity
188 of bioinformatic prediction tools is known to be suboptimal. Where functional assays are
189 available, previous studies have also demonstrated functional effects of variants predicted to
190 be tolerated/benign (21). We therefore pursued a functional analysis of the p.His198Gln
191 DNM detected in the *C1QTNF4* gene as an alternative method to add support for its
192 potential role in disease. Although its function is rather poorly understood, the protein
193 product, C1QTNF4 (CTRP4) is secreted and may act as a cytokine, as it has homology with
194 TNF and the complement component C1q (Fig. 2). C1QTNF4 has been shown to influence

195 NF- κ B activation (22), a pathway known to be implicated in SLE pathogenesis, therefore we
196 looked for an effect of the p.His198Gln mutation on NF- κ B production. Using a HEK293-NF-
197 κ B reporter cell line, we showed that C1QTNF4 p.His198Gln mutant protein was expressed
198 and that it inhibited the NF- κ B activation generated by exposure to TNF (Fig. 2).
199 Furthermore, we showed that the fibroblast L929 cell line, which is sensitive to TNF-induced
200 cell death, was rescued by exposure to C1QTNF4 p.His198Gln, but not by wild type
201 C1QTNF4. Thus, the mutant form of C1QTNF4 appears to inhibit some of the actions of TNF
202 (23–25).

203 **DNM genes do not harbour common variant associations.** We next tested for additional
204 common variant associations at these 14 loci using the high-density UK10K-1000GP3
205 imputed data. No significant association at any locus was observed with overall risk in a
206 case-control comparison, nor with anti-dsDNA ($n_{\text{cases}}=1,261$) or renal-involvement with
207 hypocomplementemia ($n_{\text{cases}}=186$) sub-phenotypes (Table S4). The lack of an associated
208 common variant within *PRKCD* and *DNMT3A* supports the hypothesis that discrete gene
209 sets will be identified through rare and common variant associations, with the former
210 expecting to be enriched for core disease genes (3).

211 Discussion

212 To fully understand the pathogenesis of complex diseases we must analyse the full
213 frequency spectrum of genetic variants (4). The study of rare variants associated with
214 disease is of paramount importance to the discovery of core genes that have the potential to
215 be therapeutic targets (12). Our data support the omnigenic hypothesis that rare genetic risk
216 may be found in a discrete set of non-canonical susceptibility genes, as we report an
217 association of collective rare variation across *PRKCD* and *DNMT3A*, and found no evidence
218 of an association with common variants across these loci. This, to the best of our knowledge,
219 is the first WES study in polygenic cases of autoimmune disease to use DNM discovery to
220 identify candidate genes for rare variant analyses. Furthermore, our study supports the

221 importance of phenotypic extremes in elucidating the genetic basis of multifactorial disorders
222 (26).

223 Searching GWAS-identified canonical disease susceptibility genes for additional rare variant
224 risk has not been fruitful. Although there are examples – and perhaps more to discover – of
225 canonical disease genes harbouring both common and rare risk alleles (27), the vast
226 majority of such loci do not. Indeed the common variant associated loci which have also
227 been shown to harbor rare coding variant risk are often those distinct minority of loci where
228 the common polymorphisms are non-silent coding variants (e.g. *NCF2* (9)). It is important to
229 note, however, that the separation of periphery and core genes may not necessarily be
230 binary (3).

231 *DNMT3A* and *PRKCD*, although hitherto not associated with polygenic SLE, are known
232 autoimmunity susceptibility loci; *DNMT3A* is associated with Crohn's disease (CD) (28) and
233 *PRKCD* is associated with both CD and ulcerative colitis (UC) (29). The notion that a locus
234 could harbour common variants contributing to one autoimmune disease and rare variants
235 contributing to another is intriguing, and could provide further hypothesis-driven searches in
236 the hunt for disease-specific core genes.

237 A functional missense variant p.G510S (c.G1528A) in *PRKCD* has previously been reported
238 in a consanguineous family with monogenic SLE (30). It was demonstrated that the *PRKCD*-
239 encoded protein, PRC δ , was essential in the regulation of B cell tolerance and affected
240 family members with the homozygous mutation had increased numbers of immature B cells.
241 Our study implicates the role of rare variants in *PRKCD* in the broader context of SLE
242 susceptibility, beyond a monogenic recessive disease model. Indeed the analysis of rare and
243 low frequency variants contributing to human height found significant overlap with genes
244 mutated in monogenic growth disorders (4). Furthermore, *PRKCB*, another member of the
245 protein kinase C gene family, has been implicated in SLE risk in a Chinese study (31).

246 *DNMT3A*, a DNA methyltransferase, is a very intriguing candidate gene for SLE as altered
247 patterns of DNA methylation are reported in autoimmune diseases (32), and
248 hypomethylation of apoptotic DNA has been reported to induce autoantibody production in
249 SLE (33). DNA methylation changes are also associated with monozygotic twin discordance
250 in SLE (34). A candidate gene study previously reported a trend of association between the
251 common *DNMT3A* intronic SNP rs1550117 (MAF~7%) and SLE in a European cohort (35).
252 Our analysis did not replicate this finding ($P=0.23$) and found no evidence of a common
253 variant association at this locus. Instead we find an association of collective rare variants
254 and SLE sub-phenotypes and emphasises the importance of deep phenotyping and the
255 potential role of rare variants in specific sub-phenotype, or indeed autoimmune,
256 manifestations. Despite progress with diagnosis and treatment, particular SLE sub-
257 phenotypes – including those used in this study - are still associated with reduced life
258 expectancy. Therefore, elucidating the specific underlying genetic risk is of paramount
259 importance.

260 Through two in vitro assays, we demonstrated the functional effect of a DNMT, p.His198Gln in
261 *C1QTNF4*, despite this mutation being predicted to be of little functional importance across
262 variant-level prediction tools. We showed the mutated protein product of *C1QTNF4*,
263 *C1QTNF4*, inhibits some TNF-mediated cellular responses, including activation of NF- κ B
264 and TNF-induced apoptosis. The role of TNF in SLE is complex and incompletely
265 understood, although, in this context, it is noteworthy that TNF inhibition may promote
266 antinuclear autoimmunity (24). Gene-level metrics for *C1QTNF4* were supportive of a role in
267 disease and our result support the importance of combined gene- and variant-level metrics,
268 and the dangers of relying heavily on variant-level metrics alone, when interpreting the
269 potential role of mutations (14). *C1QTNF6* is a known susceptibility locus for Type 1
270 Diabetes and is implicated in Rheumatoid Arthritis (36, 37), and an association with SLE has
271 recently been in a transancestral ImmunoChip analysis (38). Together these data suggest a
272 potential role of the hitherto understudied *C1QTNF* superfamily of genes in autoimmunity.

273 Although our study allowed a comprehensive approach to test all DNM genes for allelic
274 burden of rare variants, our results show that filtering based on gene- or variant-level metrics
275 would not have resulted in true associations of *DNMT3A* and *PRKCD* being missed. When
276 larger datasets require further prioritisation of genes, we suggest both variant- and gene-
277 level metrics are used.

278 Each human - regardless of the disease status - is estimated to have one DNM in their
279 exome (18). The simple presence of a provisionally functional DNM in a proband is therefore
280 not sufficient evidence that it contributes to disease risk. A major challenge of WES studies,
281 therefore, is how to differentiate between variants truly important to disease and background
282 variation (39). In light of recent studies which have demonstrated the limitations of large-
283 scale exome-wide case-control studies in detecting rare variant associations (6, 40), despite
284 such associations being found when no limitation on sample size exists (4), our results
285 support extreme-phenotype sampling and DNM discovery to aid a hypothesis-driven search
286 for rare variant associations with complex diseases, in the hunt to determine core disease
287 genes.

288 **Methods**

289 **Selection of trios for sequencing.** SLE patients of European ancestry – as determined by
290 genome-wide genotyping as part of a GWAS (19) - were selected from the UK SLE genetic
291 repository assembled in the Vyse laboratory on the following criteria: age of onset of SLE <
292 25 years (median age 21 years); more marked disease phenotype as shown by either
293 evidence for renal involvement as per standard classification criteria and/or the presence of
294 hypocomplementemia and anti-dsDNA autoantibodies; and DNA available from both
295 unaffected parents. The 30 trios (90 individuals) were exome sequenced, as described in SI
296 Methods. Ethical approval for the research was granted by the NRES Committee London
297 (12/LO/1273 and 06/MRE02/9).

298 **DNM calling.** Three bioinformatics tools with conservative parameters were used for DNM
299 screening: BCFtools (41), DeNovoGear (42) and DeNovoCheck (43). A detailed description
300 of the methods applied can be found in SI Methods. Briefly, 454 variants were identified with
301 BCFtools and DeNovoGear and eight additional variants were identified by DeNovoCheck
302 and validated by IGV, resulting in a total of 462 variants, which map to 257 genes. The
303 variants were next filtered sequentially filtered (Fig. S2): (A) Removal of NGS error prone
304 genes (NEPG); (B) Fulfil a Het:Ref:Ref for Child:Father:Mother *de novo* pattern of
305 inheritance and further selected variants that did not contain any trace of alternate allele in
306 any of the parents; (C) Non-silent variant annotation. This process resulted in a total of 17
307 variants in 17 genes (Table S1).

308 **Analysis of whole exome sequencing (WES) in cases only.** 584,798 variants with $\geq 20X$
309 coverage depth and within Gencode capture regions were identified in the analysis of 30
310 SLE probands only. Stringent filters were applied for variant refinement, described in full in
311 SI Methods, resulting in 1194 variants in 1067 genes (Fig.S3).

312 **Sanger Sequencing confirmation.** Primers were designed using Primer 3. 10ng of DNA
313 from SLE probands, any unaffected siblings and both parents was amplified with Hot Start
314 Taq polymerase. PCR products were first purified with EXO-SAP before BigDye labelling in
315 a linear PCR and sequenced on an ABI 3300XL. Primers and PCR conditions available on
316 request. The reads were analysed using Chromas Lite (v.2.1.1)

317 **Imputation.** Illumina HumanOmni1 BeadChip genotype data from 6,995 controls and 4,036
318 SLE patients of matched European ancestry were used, which had undergone quality control
319 as previously described including Principal Component Analysis (PCA) to account for
320 population structure (19).The UK10K (REL-2012-06-02) plus 1000 Genomes Project Phase3
321 data (release 20131101.v5) merged reference panel (UK10K-1000GP3) was accessed
322 through the European Genome-phenome Archive (EGAD00001000776). The genotype data
323 were imputed using the UK10K-1000GP3 reference panel across the coding regions of the
324 14 DNM genes plus a 2Mb flanking region. To increase the accuracy of imputed genotype

325 calls, a full imputation without pre-phasing was conducted using IMPUTE2 (44, 45). Imputed
326 genotypes were filtered for confidence using an info score (IMPUTE2) threshold of 0.3 (Fig.
327 S6 and S7). The most likely genotype from IMPUTE2 was taken if its probability was > 0.5. If
328 the probability fell below this threshold, it was set as missing. Variants with >10% missing
329 genotype calls were removed for further analysis. All individuals had <8% missing genotype
330 data.

331 **Rare variant burden tests.** Imputed data were filtered, using Plink v1.9, to include only
332 variants mapping to coding exons of hg19 RefSeq transcripts. Plink/SEQv1.0 (20) was used
333 to run gene-wise one-tailed burden testing with a MAF<1% threshold. A 5% false discovery
334 rate was used for multiple testing correction for 14 genes.

335 **Common variant association tests.** SNPTEST 2.5.2 (46) was used to test for associated
336 variants with MAF>1% across the region spanning the encoded gene. The first four
337 covariates from the original GWAS were included (19). Bonferroni correction was used for
338 3,000 tests across the loci ($q=1.66E-5$).

339 **Plasmids.** Myc-Flag-tagged *C1QTNF4* on the pCMV6 vector and the empty pCMV6 vector
340 were used (OriGene). The mutant pCMV6-*C1QTNF4 C594G* (p.His198Gln) was generated
341 by site-directed mutagenesis (Quikchange II XL; Stratagene) according the manufacturer
342 instructions: mutagenic primer: 5'-GCGAGTGGTTGCTGCCGCGGCC-3' (Sigma Aldrich).
343 The plasmids production was carried out in XL10-Gold Ultracompetent cells, isolated and
344 purified using EndoFree Maxi Prep kit (Qiagen) and plasmid ORFs were confirmed by full
345 Sanger sequencing (GATC-Biotech). The expression and secretion of the flagged proteins
346 was confirmed by western blot on cell lysates and supernatants with monoclonal anti-FLAG
347 antibody (clone M2; Sigma-Aldrich).

348 **Luciferase assays and TNF-induced programmed cell death.** GloResponse NF- κ B-RE-
349 luc2P HEK293 cell line (Promega) and TNF-sensitive L929 fibrosarcoma cell line (ATCC)
350 were cultured in Dulbecco's Modified Eagle Medium (DMEM) supplemented with 10% fetal

351 bovine serum (FBS) and 1% Penicillin/Streptomycin at 37°C, 5% CO₂. HEK293 were seeded
352 24 hours before transfection in antibiotic free DMEM in 96 wells plate (2×10⁴ cells/well),
353 transfected with either *C1QTNF4*, *C1QTNF4 C594G* or Empty Vector via Fugene HD
354 (Promega). 48 hours after transfection the cell were left unstimulated or stimulated with
355 TNFα 5 ng/ml (PeproTech) for 4 hours. Luciferase activity was assayed by One-Glo
356 (Promega) on Berthold Orion luminometer, the values were normalized to cell viability
357 measured by CellTiter Glo (Promega). L929 were challenged with TNFα 0.45 ng/ml and
358 Actinomycin D 1 μg/ml (R&D) for 24 hours in presence of C1QTNF4 or C1QTNF4
359 p.His198Gln containing media, cell viability was measured by CellTiter Glo.

360 **Size exclusion chromatography.** Supernatants (750 μl) of HEK293 producing C1QTNF4
361 or C1QTNF4 p.His198Gln were buffer exchanged in PBS on Zeba Spin Desalting Columns
362 (Thermo Fisher) and 0.5 mL loaded on an AKTA FPLC with a Superdex 200 10/300 GL
363 column (GE Healthcare). Absorbance was normalized to the maximum peak of each sample.

364 **Data availability.** WES data on 90 individuals – 30 parent-offspring trios – will be deposited
365 at the European Genome-phenome Archive.

366 **Acknowledgements**

367 The work leading to these results received funding from the European Union FP7
368 programme (grant agreement n° 262055) via the European Sequencing and Genotyping
369 Infrastructure (ESGI). Sequencing was performed by the SNP&SEQ Technology Platform in
370 Uppsala, which is part of the National Genomics Infrastructure (NGI) hosted by Science for
371 Life Laboratory in Sweden. This work was supported in part by the Swedish Research
372 Council for Medicine and Health (grant n° E0226301) and by the Knut and Alice Wallenberg
373 Foundation (KAW 2011.0073). We thank Johanna Lagensjö and Olof Karlberg for assistance
374 with sequencing. The research was funded/supported by the National Institute for Health
375 Research (NIHR) Biomedical Research Centre based at Guy's and St Thomas' NHS
376 Foundation Trust and King's College London.

377

378 **References**

- 379 1. Manolio TA, et al. (2009) Finding the missing heritability of complex diseases. *Nature*
380 461:747–753.
- 381 2. Yang J, et al. (2010) Common SNPs explain a large proportion of the heritability for
382 human height. *Nat Gen* 42(7):565–569.
- 383 3. Boyle EA, Li YI, Pritchard JK (2017) Leading Edge Perspective An Expanded View of
384 Complex Traits: From Polygenic to Omnigenic. *Cell* 169(7):1177–1186.
- 385 4. Marouli E, et al. (2017) Rare and low-frequency coding variants alter human adult
386 height. *Nature* 542(7640):186–190.
- 387 5. Shi H, Kichaev G, Pasaniuc B (2016) Contrasting the Genetic Architecture of 30
388 Complex Traits from Summary Association Data. *Am J Hum Genet* 99(1):139–153.
- 389 6. Fuchsberger C, et al. (2016) The genetic architecture of type 2 diabetes. *Nature*
390 536:41–47.
- 391 7. Purcell SM, et al. (2014) A polygenic burden of rare disruptive mutations in
392 schizophrenia. *Nature* 506(7487):185–190.
- 393 8. Ripke S, et al. (2014) Biological insights from 108 schizophrenia-associated genetic
394 loci. *Nature* 511:421–427.
- 395 9. Hunt K a, et al. (2013) Negligible impact of rare autoimmune-locus coding-region
396 variants on missing heritability. *Nature* 498(7453):232–5.
- 397 10. Bamshad MJ, et al. (2011) Exome sequencing as a tool for Mendelian disease gene
398 discovery. *Nat Rev Genet* 12(11):745–755.
- 399 11. Moutsianas L, et al. (2015) The power of gene-based rare variant methods to detect
400 disease-associated variation and test hypotheses about complex disease. *PLoS*
401 *Genet* 11(4):e1005165.
- 402 12. Chakravarti A, Turner TN (2016) Revealing rate-limiting steps in complex disease
403 biology: The crucial importance of studying rare, extreme-phenotype families.
404 *BioEssays* 38(6):578–586.
- 405 13. Rahbari R, et al. (2015) Timing, rates and spectra of human germline mutation. *Nat*
406 *Genet* 48(December):1–11.
- 407 14. Itan Y, et al. (2015) The human gene damage index as a gene-level approach to
408 prioritizing exome variants. *Proc Natl Acad Sci U S A* 112(44):13615–20.
- 409 15. Lek M, et al. (2016) Analysis of protein-coding genetic variation in 60,706 humans.
410 *Nature* 536(7616):285–291.
- 411 16. Kircher M, et al. (2014) A general framework for estimating the relative pathogenicity
412 of human genetic variants. *Nat Genet* 46(3):310–315.
- 413 17. Chen L, Morris DL, Vyse TJ (2017) Genetic advances in systemic lupus
414 erythematosus. *Curr Opin Rheumatol* 29:423–433.
- 415 18. Veltman J a., Brunner HG (2012) De novo mutations in human genetic disease. *Nat*
416 *Rev Genet* 13(8):565–575.
- 417 19. Bentham J, et al. (2015) Genetic association analyses implicate aberrant regulation of
418 innate and adaptive immunity genes in the pathogenesis of systemic lupus

- 419 erythematosus. *Nat Genet* 47(12):1457–1464.
- 420 20. Purcell SM PLINK/SEQ: A library for the analysis of genetic variation data. Available
421 at: <https://atgu.mgh.harvard.edu/plinkseq/>.
- 422 21. Roberts AL, et al. (2014) Resequencing the susceptibility gene, ITGAM, identifies two
423 functionally deleterious rare variants in systemic lupus erythematosus cases. *Arthritis*
424 *Res Ther* 16(3):R114.
- 425 22. Li Q, et al. (2011) Identification of C1qTNF-related protein 4 as a potential cytokine
426 that stimulates the STAT3 and NF- κ B pathways and promotes cell survival in human
427 cancer cells. *Cancer Lett* 308(2):203–214.
- 428 23. Beigel F, et al. (2011) Formation of antinuclear and double-strand DNA antibodies
429 and frequency of lupus-like syndrome in anti-TNF- α antibody-treated patients with
430 inflammatory bowel disease. *Inflamm Bowel Dis* 17(1):91–98.
- 431 24. Eriksson C, Engstrand S, Sundqvist K-G, Rantapää-Dahlqvist S (2005) Autoantibody
432 formation in patients with rheumatoid arthritis treated with anti-TNF alpha. *Ann Rheum*
433 *Dis* 64(3):403–7.
- 434 25. Pink AE, Fonia A, Allen MH, Smith CH, Barker JNWN (2010) Antinuclear antibodies
435 associate with loss of response to antitumour necrosis factor-alpha therapy in
436 psoriasis: a retrospective, observational study. *Br J Dermatol* 162(4):780–5.
- 437 26. Turner TN, et al. (2015) Loss of δ -catenin function in severe autism. *Nature*
438 520(7545):51–6.
- 439 27. Jordan CT, et al. (2012) Rare and common variants in CARD14, encoding an
440 epidermal regulator of NF-kappaB, in psoriasis. *Am J Hum Genet* 90(5):796–808.
- 441 28. Franke A, et al. (2010) Genome-wide meta-analysis increases to 71 the number of
442 confirmed Crohn's disease susceptibility loci. *Nat Genet* 42(12):1118–25.
- 443 29. Jostins L, et al. (2012) Host-microbe interactions have shaped the genetic
444 architecture of inflammatory bowel disease. *Nature* 491(7422):119–24.
- 445 30. Belot A, et al. (2013) Protein kinase C?? deficiency causes mendelian systemic lupus
446 erythematosus with B cell-defective apoptosis and hyperproliferation. *Arthritis Rheum*
447 65(8):2161–2171.
- 448 31. Sheng Y-J, et al. (2011) Follow-up study identifies two novel susceptibility loci PRKCB
449 and 8p11.21 for systemic lupus erythematosus. *Rheumatology (Oxford)* 50(4):682–
450 688.
- 451 32. Ballestar E (2011) Epigenetic alterations in autoimmune rheumatic diseases. *Nat Rev*
452 *Rheumatol* 7(5):263–71.
- 453 33. Wen ZK, et al. (2007) DNA hypomethylation is crucial for apoptotic DNA to induce
454 systemic lupus erythematosus-like autoimmune disease in SLE-non-susceptible mice.
455 *Rheumatology* 46(12):1796–803.
- 456 34. Javierre BM, et al. (2010) Changes in the pattern of DNA methylation associate with
457 twin discordance in systemic lupus erythematosus. *Genome Res* 20(2):170–9.
- 458 35. Piotrowski P, Grobelna MK, Wudarski M, Olesinska M, Jagodzinski PP (2015)
459 Genetic variants of DNMT3A and systemic lupus erythematosus susceptibility. *Mod*
460 *Rheumatol* 25(1):96–99.
- 461 36. Onengut-Gumuscu S, et al. (2015) Fine mapping of type 1 diabetes susceptibility loci
462 and evidence for colocalization of causal variants with lymphoid gene enhancers. *Nat*

- 463 *Genet* 47(4):381–6.
- 464 37. Murayama MA, et al. (2015) CTRP6 is an endogenous complement regulator that can
465 effectively treat induced arthritis. *Nat Commun* 6:8483.
- 466 38. Langefeld CD, et al. (2017) Transancestral mapping and genetic load in systemic
467 lupus erythematosus. *Nat Commun* In press(May). doi:10.1038/ncomms16021.
- 468 39. Cooper GM, Shendure J (2011) Needles in stacks of needles: finding disease-causal
469 variants in a wealth of genomic data. *Nat Rev Genet* 12(9):628–40.
- 470 40. Luo Y, et al. (2016) Exploring the genetic architecture of inflammatory bowel disease
471 by whole genome sequencing identifies association at ADCY7. *Nat Genet* 49(2):186–
472 192.
- 473 41. Li H, et al. (2009) The Sequence Alignment/Map format and SAMtools. *Bioinformatics*
474 25(16):2078–2079.
- 475 42. Ramu A, et al. (2013) DeNovoGear: de novo indel and point mutation discovery and
476 phasing. *Nat Methods* 10(10):985–987.
- 477 43. de Ligt J, et al. (2012) Diagnostic exome sequencing in persons with severe
478 intellectual disability. *N Engl J Med* 367(20):1921–9.
- 479 44. Howie B, Fuchsberger C, Stephens M, Marchini J, Abecasis GR (2012) Fast and
480 accurate genotype imputation in genome-wide association studies through pre-
481 phasing. *Nat Genet* 44(8):955–959.
- 482 45. Roshyara NR, et al. (2016) Comparing performance of modern genotype imputation
483 methods in different ethnicities. *Sci Rep* 6:34386.
- 484 46. Marchini J, Howie B (2010) Genotype imputation for genome-wide association
485 studies. *Nat Rev Genet* 11(7):499–511.
- 486 47. Kelley LA, Sternberg MJE (2009) Protein structure prediction on the Web: a case
487 study using the Phyre server. *Nat Protoc* 4(3):363–371.

488

489

490

491

492

493

494

495

496

497

498

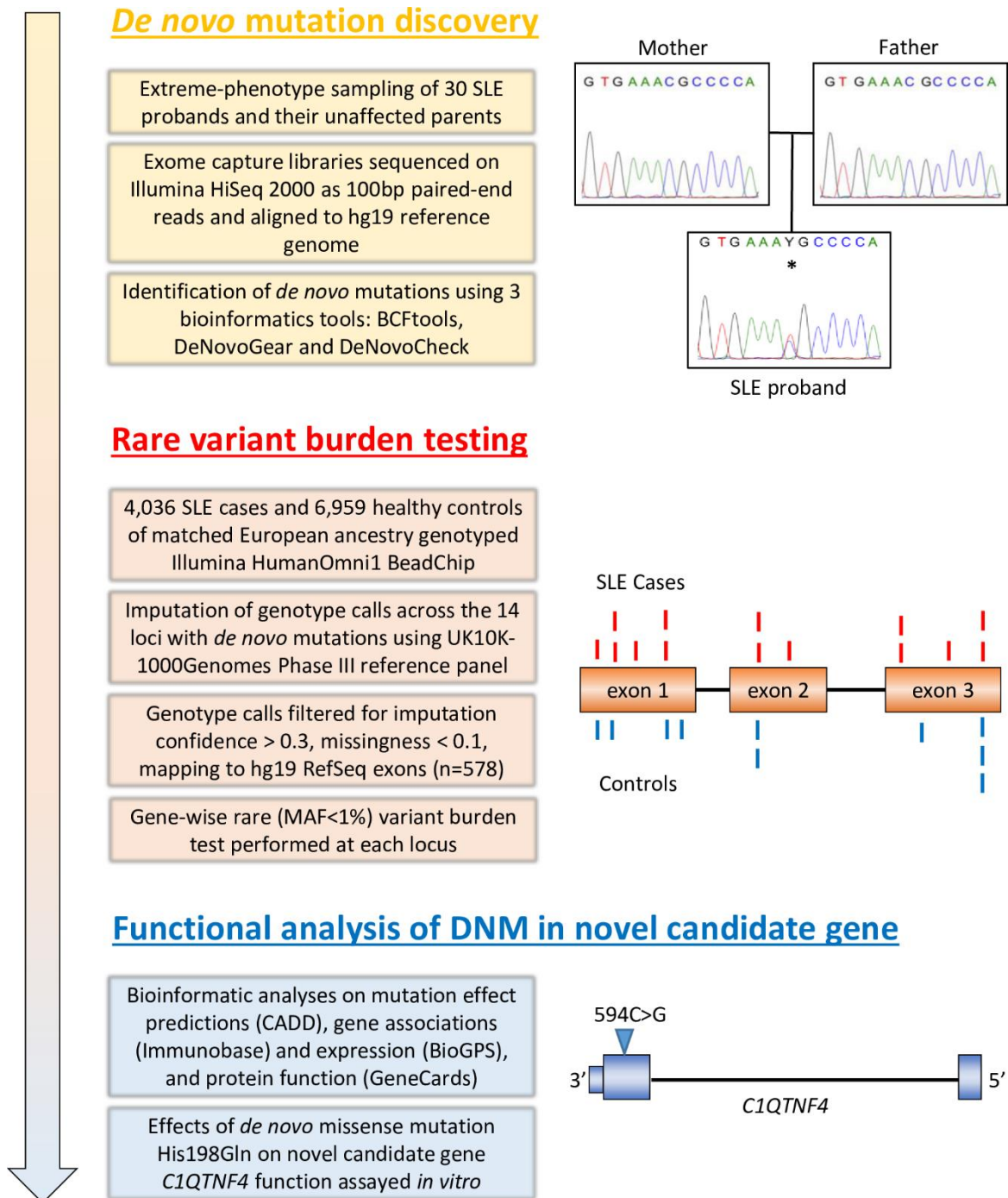
499

500

501 **Figures**

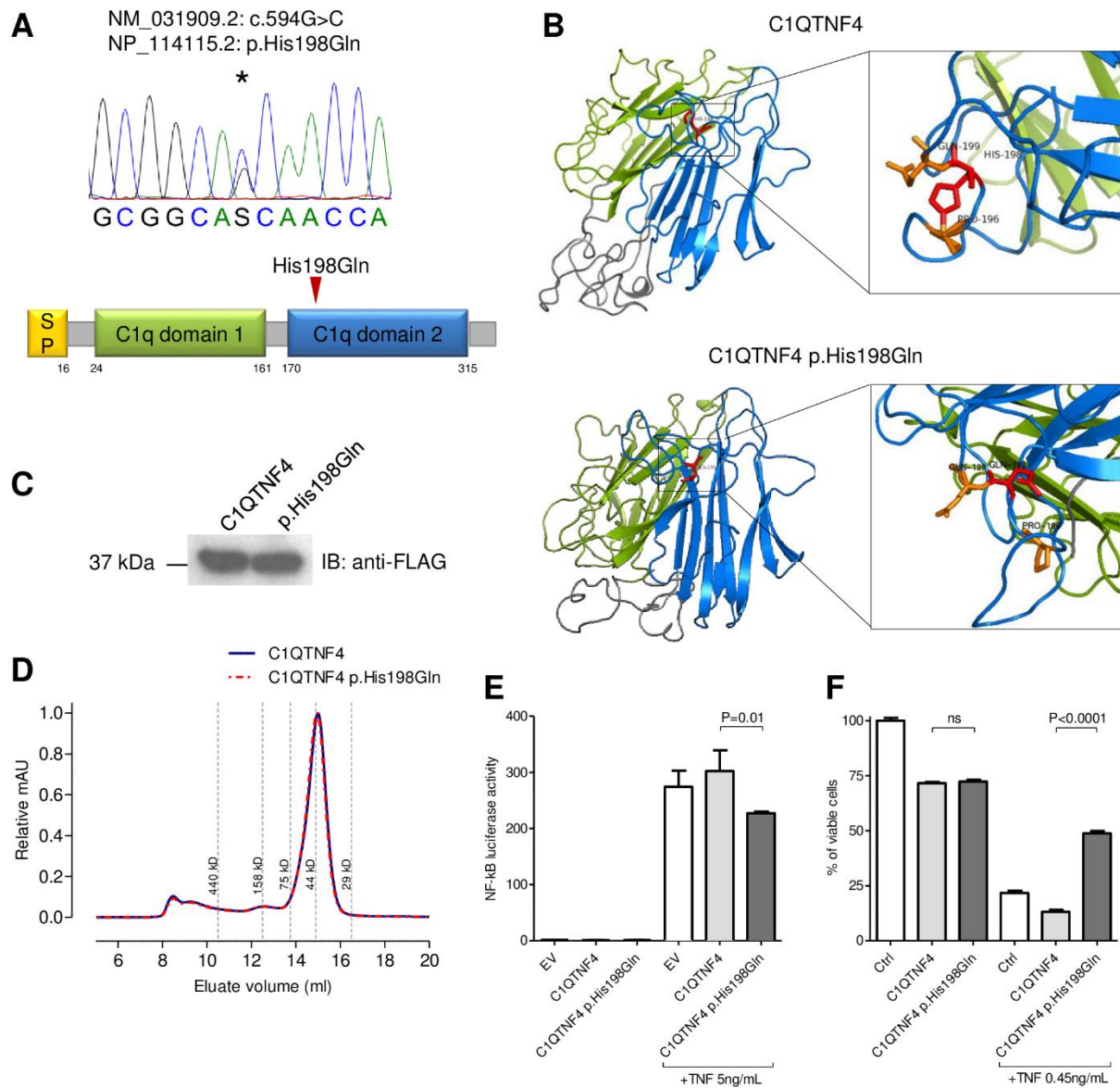
502

503 **Figure 1. Overview of study.** *De novo* mutations (DNM) in a discovery cohort revealed
504 candidate genes for imputation-based rare variant burden testing using a follow-up cohort.
505 Independent functional analyses demonstrate the functional effects of one DNM in a
506 candidate gene.



507

508 **Figure 2. Structural and functional characterization of C1QTNF4 p.His198Gln**
509 **substitution.** (A) Domain organization of human C1QTNF4, showing signal peptide (yellow),
510 first C1q domain (green), second C1q domain (blue) and linker peptides (grey). Arrow
511 highlights substitution site. (B) 3D structure prediction of C1QTNF4 and C1QTNF4
512 p.His198Gln using Phyre2 (47). Ribbons show the interaction between the positively
513 charged Histidine 198 and Proline 196 lost in C1QTNF4 p.His198Gln due to the substitution
514 of Histidine with Glutamine. (C) Immunoblot demonstrating that p.His198Gln does not affect
515 secretion of C1QTNF4 in HEK293 supernatants. (D) Size exclusion chromatography profile
516 showing no difference in oligomerisation between supernatant containing C1QTNF4 (blue)
517 and C1QTNF4 p.His198Gln (red). (E) Luciferase assay in HEK293-NF- κ B reporter cell line
518 showing that C1QTNF4 p.His198Gln inhibits NF- κ B activation in response to 4h stimulation
519 with 5ng/mL TNF α . Error bars represent standard error of the mean. (F) Inhibition of L929
520 induced cell death by C1QTNF4 p.His198Gln after 24h of stimulation with 0.45 ng/mL TNF α
521 in presence of Actinomycin 1 μ g/ml. EV=empty vector.



522

523

524

525

526

527

528

529

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

Family	Mutation (chr:position ref:alt)	Gene	Gene Description	Exon	Amino acid	MAF in ExAC ^a	CADD Phred	Mutation Type ^b
SLE0751	22:38336799 C:T	<i>MICALL1</i>	MICAL-like 1	16	Arg852Cys	1.5 x 10 ⁻⁴	35	Ti CpG
SLE0496	3:53223122 G:A	<i>PRKCD</i>	protein kinase C, delta	16	Gly535Arg	-	34	Ti CpG
SLE0679	12:57588368 C:T	<i>LRP1</i>	Low-density lipoprotein receptor-related protein 1	50	Arg2693Cys	8.3 x 10 ⁻⁴	34	Ti CpG
SLE0592	6:36260896 G:A	<i>PNPLA1</i>	patatin-like phospholipase domain containing 1	3	Arg166His	5.8 x 10 ⁻⁵	33	Ti CpG
SLE0296	2:25457236 G:A	<i>DNMT3A</i>	DNA (cytosine-5-)-methyltransferase 3 alpha	19	Ala695Val	-	32	Ti CpG
SLE0571	4:79512728 G:T	<i>ANXA3</i>	annexin A3	7	Ser145Ile	-	25.2	Tv
SLE0679	3:171431716 G:A	<i>PLD1</i>	phospholipase D1, phosphatidylcholine-specific	9	Thr293Met	5.8 x 10 ⁻⁵	25.1	Ti CpG
SLE0411	5:179743769 C:T	<i>GFPT2</i>	glutamine-fructose-6-phosphate transaminase 2	12	Val383Met	2.6 x 10 ⁻⁵	23.4	Ti CpG
SLE0679	7:138968784 C:A	<i>UBN2</i>	ubiquitin 2	15	Pro1045Thr	-	18.46	Tv
SLE0080	16:2812426 C:T	<i>SRRM2</i>	serine/arginine repetitive matrix 2	11	Arg633Cys	-	14.32	Ti CpG
SLE0852	11:47611769 G:C	<i>C1QTNF4</i>	C1q and tumor necrosis factor related protein 4	2	His198Gln	-	12.29	Tv
SLE0321	18:61621642 G:A	<i>HMSD</i>	histocompatibility (minor) serpin domain containing	3	Ala25Thr	-	9.732	Ti
SLE0390	12:32369376 G:C	<i>BICD1</i>	bicaudal D homolog 1 (Drosophila)	2	Val137Leu	-	8.673	Tv
SLE0321	1:35251125 C:G	<i>GJB3</i>	gap junction protein, beta 3	2	Asp254Glu	-	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 ³⁰	IBD, UC, CD ²⁸	Dendritic	3.75*
<i>DNMT3A</i>	DNA methyltransferase	Candidate gene study ³⁵	CD ²⁹	-	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>GFPT2</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 Table S5

^b See Table S6

^c See SFigure S4. 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 with an asterisk.