

1 **Genome-wide analysis of genetic risk factors for rheumatic heart disease in** 2 **Aboriginal Australians provides support for pathogenic molecular mimicry**

3 Short title: A GWAS of RHD in Aboriginal Australians

4 Lesley-Ann Gray^{†, 1,2} Heather A D'Antoine^{†, 3} Steven Y.C. Tong^{†, 3,4} Melita McKinnon,³ Dawn
5 Bessarab,⁵ Ngiare Brown,⁶ Bo Reményi,³ Andrew Steer,⁷ Genevieve Syn,⁸ Jenefer M
6 Blackwell*,⁸ Michael Inouye*,^{1,2} Jonathan R Carapetis*,^{3,8}

7 [†] Equal first authors; * Equal senior authors

8 ¹ School of BioSciences, The University of Melbourne, Parkville 3010, Victoria, Australia (L.G.,
9 M.I.)

10 ² Department of Pathology, The University of Melbourne, Parkville, Victoria, Australia (L.G., M.I.)

11 ³ Menzies School of Health Research, Charles Darwin University, Darwin, Northern Territory,
12 Australia (H.A.D'A., S.Y.C.T., M.McK., B.R., J.C.)

13 ⁴ Victorian Infectious Disease Service, The Royal Melbourne Hospital; and Peter Doherty Institute
14 for Infection and Immunity, The University of Melbourne, Victoria, Australia (S.Y.C.T.)

15 ⁵ Centre for Aboriginal Medical and Dental Health, The University of Western Australia, Crawley,
16 Western Australia, Australia (D.B.)

17 ⁶ School of Education, The University of Wollongong, Wollongong, New South Wales, Australia
18 (N.B.)

19 ⁷ Group A Streptococcal Research Group, Murdoch Childrens Research Institute, Melbourne,
20 Victoria, Australia; and Centre for International Child Health, Department of Paediatrics, Royal
21 Children's Hospital, Melbourne, Victoria, Australia (A.S.)

22 ⁸ Telethon Kids Institute, The University of Western Australia, Roberts Road, Subiaco, Western
23 Australia, Australia (J.B., J.C.)

24 #Current address: Baker Heart & Diabetes Institute, Melbourne, Australia
25

26
27 **Corresponding authors:** Jenefer M. Blackwell (jenefer.blackwell@telethonkids.org.au); Michael
28 Inouye (minouye@baker.edu.au); Jonathan Carapetis (Jonathan.Carapetis@telethonkids.org.au)
29

30 **Correspondence to:** Jonathan Carapetis, PO Box 855, West Perth, Western Australia 6872; 100
31 Roberts Road, Subiaco, Western Australia 6009; Phone: +61 8 9489 7967.
32

33 **Funding** National Health and Medical Research Council (NHMRC) Grant APP1023462, NHMRC and
34 National Heart Foundation of Australia Career Development Fellow (#1061435), NHMRC Career
35 Development Fellow (#1065736).
36

37 **Competing interests**

38 The authors disclose no conflicts of interest. The study sponsor had no role in study design, data
39 collection, data analysis, data interpretation, or writing of the report. The corresponding author had full
40 access to all study data and had final responsibility for the decision to submit for publication.
41

42 **Abstract**

43

44 **Background.** Rheumatic heart disease (RHD) following Group A Streptococcus
45 (GAS) infections is heritable and prevalent in Indigenous populations. Molecular
46 mimicry between human and GAS proteins triggers pro-inflammatory cardiac
47 valve-reactive T-cells.

48

49 **Methods.** Genome-wide genetic analysis was undertaken in 1263 Aboriginal
50 Australians (398 RHD cases; 865 controls). Single nucleotide polymorphisms
51 (SNPs) were genotyped using Illumina HumanCoreExome BeadChips. Direct
52 typing and imputation was used to fine-map the human leukocyte antigen (HLA)
53 region. Epitope binding affinities were mapped for human cross-reactive GAS
54 proteins, including M5 and M6.

55

56 **Results.** The strongest genetic association was intronic to HLA-DQA1
57 (rs9272622; $P=1.86 \times 10^{-7}$). Conditional analyses showed rs9272622 and/or
58 DQA1*AA16 account for the HLA signal. HLA-DQA1*0101_DQB1*0503 (OR 1.44,
59 95%CI 1.09-1.90, $P=9.56 \times 10^{-3}$) and HLA-DQA1*0103_DQB1*0601 (OR 1.27,
60 95%CI 1.07-1.52, $P=7.15 \times 10^{-3}$) were risk haplotypes; HLA-DQA1*0301-
61 DQB1*0402 (OR 0.30, 95%CI 0.14-0.65, $P=2.36 \times 10^{-3}$) was protective. Human
62 myosin cross-reactive N-terminal and B repeat epitopes of GAS M5/M6 bind with
63 higher affinity to DQA1/DQB1 alpha/beta dimers for the two risk haplotypes
64 than the protective haplotype.

65

66 **Conclusions.** Variation at HLA_DQA1-DQB1 is the major genetic risk factor for
67 RHD in Aboriginal Australians studied here. Cross-reactive epitopes bind with
68 higher affinity to alpha/beta dimers formed by risk haplotypes, supporting
69 molecular mimicry as the key mechanism of RHD pathogenesis.

70

71

72 Abstract: 197 Words.

73

74 **Keywords.** GWAS; HLA; Acute rheumatic fever; rheumatic heart disease;
75 epitope mapping; M proteins; molecular mimicry.

76

77

78

79

80 INTRODUCTION

81

82 Acute rheumatic fever (ARF) results from an autoimmune response to infections
83 due to group A *Streptococcus* (GAS), *Streptococcus pyogenes*. Recurrences of ARF
84 and its associated cardiac valvular inflammation lead to chronic valvular damage
85 and rheumatic heart disease (RHD). RHD causes an estimated 275,000 deaths
86 annually with an estimated 33 million prevalent cases globally (reviewed [1]). In
87 Australia, RHD is most prevalent in the Indigenous population, affecting 2-6 per
88 1000 individuals (and as high as 15/1000 school aged children in the northern
89 tropical regions[2])[3, 4].

90

91 The precise pathological mechanisms underlying RHD remain unclear. One
92 hypothesis to explain inflammation of valvular tissue is molecular mimicry
93 (reviewed [1, 5-8]). Accordingly, peptides from GAS proteins are processed by
94 antigen presenting cells in the throat and heart tissue and presented on Human
95 Leukocyte Antigen (HLA) class II molecules to CD4+ T lymphocytes that elicit
96 pro-inflammatory cytokine responses and/or provide help to B lymphocytes for
97 antibody secretion. In RHD patients, the CD4 T cell epitopes and antigenic
98 specificities of antibodies show cross-reactivity to proteins in heart tissue,
99 specifically targeting cardiac valves [5, 9]. This cross-reactivity is thought to be
100 due to sequence similarities between heart tissues and GAS proteins, amongst
101 which GAS M-proteins feature prominently [10], which is supported by studies of
102 HLA-DQ-restricted T cell clones that recognise the M protein and myosin
103 peptides in the blood and hearts of RHD patients [11, 12] as well as studies in
104 animal models of disease [13] . The precise mechanism by which these cross-

105 reactive antibodies target the valve is unclear, and cross-reactive antibodies
106 have been observed in streptococcal pharyngitis without complications [14]. An
107 alternative hypothesis (reviewed [7]) is that a streptococcal M protein N-
108 terminus domain binds to the CB3 region in collagen type IV, initiating an
109 antibody response to the collagen which results in inflammation. These
110 antibodies do not cross-react with M proteins, and hence do not involve
111 molecular mimicry.

112

113 Key aspects of molecular mimicry are the relevant proteins/peptides in GAS
114 strains and host susceptibility. In the Northern Territory of Australia there is
115 high genetic diversity amongst GAS strains which reflect global-scale
116 transmission rather than localised diversification [15, 16]. Despite ubiquitous
117 exposure to GAS, only 1-2% of Indigenous Australians living in this region
118 develop RHD, and the cumulative incidence of ARF only reaches 5-6% in
119 communities with the most complete case ascertainment [17]. ARF is a
120 precursor to RHD, and in a meta-analysis of 435 twin pairs susceptibility to
121 rheumatic fever was estimated to be 60% heritable [18]. For RHD, a number of
122 candidate gene studies have variably reported associations with genes
123 controlling innate and adaptive immune responses (reviewed [6]). Among these
124 candidates, HLA Class I and II genes feature most prominently, but with little
125 consistency in risk and protective genes/alleles reported [6, 19, 20]. Recently, a
126 GWAS of RHD was performed in Oceania populations but did not report an HLA
127 signal [21]. This variability in reported associations likely reflects differing study
128 designs, population-related genetic heterogeneity, failure to control for
129 confounding factors, and the vagaries of small samples sizes and candidate gene

130 approaches. Here we undertake an unbiased genome-wide approach to identify
131 genetic risk factors for RHD in echocardiogram-confirmed cases from the
132 Northern Territory of Australia. The HLA-DQA1-DQB1 locus was the only region
133 to show strong association in this population. We show that differential binding
134 of GAS/human cross-reactive epitopes to MHC Class II dimers for specific HLA
135 DQA1_DQB1 risk and protective haplotypes may underpin the molecular
136 mimicry hypothesis for RHD pathogenesis.

137

138 **METHODS**

139

140 **Ethical Considerations, Sampling and Clinical Data Collection**

141 This study was undertaken with ethical approval from the Human Research
142 Ethics Committee (HREC) of the Northern Territory Department of Health and
143 Menzies School of Health Research (ID HREC-2010-1484) and the Central
144 Australian HREC (ID HREC-2014-241). The study was overseen by a project
145 steering committee and three sub-committees – Aboriginal governance, clinical
146 and scientific. The protocol and any key changes required agreement from the
147 Aboriginal governance committee. Stage 1 of the project involved community
148 engagement and consent, development of culturally appropriate consent
149 materials, and establishment of appropriate governance for collection and
150 subsequent storage of samples. Stage 2 involved identifying individual
151 participants, obtaining informed consent, and collection of samples and
152 associated meta-data. The individual consent incorporated an “opt-in” design
153 where participants selected which components of the study they were
154 comfortable to participate in, and were able to withdraw from the study at any

155 stage [22]. This included an option to accept or refuse continued use of their
156 genetic or clinical data in further studies. De-identified post QC (cf. below)
157 genotype data for individuals who consented to continued use of their data have
158 been lodged in the European Genome-phenome Archive (accession number
159 EGAS00000000000) with access controlled through a study-specific Data Access
160 Committee.

161

162 Participants were recruited from 19 communities in the Northern Territory of
163 Australia (Figure 1). Case participants were defined a priori as having had, at
164 some stage, echocardiographically confirmed evidence of RHD and/or ARF with
165 carditis. For each of the 19 communities, we obtained a list of individuals on the
166 Northern Territory Rheumatic Heart Disease register. These lists were further
167 screened for patients with a history of ARF and associated carditis (defined using
168 the 2015 revised Jones Criteria [23]) or RHD confirmed on echocardiogram
169 (defined using the 2012 World Heart Federation criteria [24]). We aimed for a
170 1:2 ratio of cases to controls. Controls were selected from the same communities
171 (range 4-215 participants/community) to ensure similar likelihood of exposure
172 to GAS among cases and controls, and included a selection of family members as
173 well as unrelated community-based controls. Medical records of potential
174 control participants were checked to exclude a prior history of rheumatic fever.
175 We did not perform echocardiograms on control participants. Both cases and
176 controls had to be aged ≥ 18 years, to minimise the likelihood of enrolling
177 controls that might subsequently become cases (given that ARF is largely a
178 disease of school aged children and most RHD cases are diagnosed before the age

179 of 30). Data were collected for age, gender, community location and RHD
180 case/control status.

181

182 We collected clinical data and saliva from 1382 individuals. Of these, 11 later
183 withdrew consent for the study, and an additional 71 individuals were deemed
184 ineligible for case or control status following detailed medical record review,
185 leaving 1291 eligible to include in the study. Demographic details (age, sex,
186 case/control status) for the 1263 (of 1291) study participants who also passed
187 QC following genotyping (cf. below) are summarised in Table S1.

188

189 Array Genotyping and Marker QC

190 Saliva was collected using Oragene OG-500 saliva kits (DNA Genotek Inc.,
191 Ontario, Canada) and DNA extracted according to manufacturer protocols. DNAs
192 were genotyped on the Illumina Infinium® HumanCoreExome Beadchip
193 (Illumina Inc., San Diego, CA, USA), which includes probes for 547,644 single
194 nucleotide polymorphisms (SNPs), 281,725 of which are genome-wide tag SNPs
195 that represent core content and are highly informative across ancestries, and
196 265,919 SNPs that are exome-focused markers. All genotyping data and
197 reference panels were analysed using human genome build 37 (hg19).
198 Individuals were excluded if they had a missing data rate >5%. SNP variants
199 were excluded if they had genotype missingness >5%, minor allele frequency
200 (MAF) <0.01, or if they deviated from Hardy-Weinberg equilibrium (HWE;
201 threshold of $P < 1.0 \times 10^{-6}$). This provided a post-QC dataset of 1263 individuals
202 genotyped for 239,536 markers. This sample comprised 398 cases and 865
203 controls (Table S1), providing 68% power to detect genome wide significance

204 ($P < 5 \times 10^{-8}$) for genetic effects with a disease allele frequency of 0.25, effect size
205 (genotype relative risk) of 2, and assuming a disease prevalence of 2%. Overt
206 non-Aboriginal population stratification was assessed using the top 10 principal
207 components (PCs) from FlashPCA [25].

208

209 **SNP Imputation and GWAS**

210 Imputation of missing and unassayed genetic variants was performed using the
211 1000 Genomes Project phase 3 reference panel [26], which contains 88 million
212 variants for 2502 samples from 26 populations throughout Africa, America, East
213 Asia, Europe and South-East Asia. Array variants were phased using SHAPEIT v2
214 (r644) [27] and imputed with IMPUTE v2.3.2 [28]. We excluded imputed SNPs
215 with an information metric < 0.4 or genotype probability < 0.9 , and the remaining
216 variants were converted to genotype calls and filtered for $< 10\%$ missingness and
217 $MAF > 0.01$. Imputation accuracy was assessed using the r^2 metric ($r^2 > 0.8$), which
218 represents the squared Pearson correlation between the imputed SNP dosage
219 and the known allele dosage.

220

221 Genome-wide association analysis for the RHD phenotype was performed
222 using a linear mixed model as implemented in FaST-LMM v2.07, which takes
223 account of both relatedness and population substructure [29]. Age and gender
224 were included as fixed effects in the model. Population structure and relatedness
225 were controlled using the genetic similarity matrix, computed from 41,926 LD-
226 pruned array variants, and any systematic confounding assessed using QQ plots
227 and a test statistic inflation factor (λ). Genome-wide significance was set at
228 $P \leq 5 \times 10^{-8}$ [30].

229

230 **Fine-Mapping Associations in the HLA Region**

231 Conditional association analyses in the HLA region also utilised FaST-LMM.
 232 Univariate conditional analysis can fail to uncover residual signals due to the
 233 long-distance haplotypes observed in the HLA region [31], therefore we used a
 234 step-wise conditional analysis of classical HLA alleles and amino acids to scan for
 235 independent signals in HLA. First, we typed exons of 10 classical HLA alleles for
 236 716 samples using the TruSight HLA sequencing panel and produced 4-digit
 237 phase-resolved genotype calls against the IMGT v3210 database (Murdoch
 238 University Centre for Clinical Immunology and Biomedical Statistics, Perth,
 239 Western Australia). We generated an Aboriginal reference panel of typed HLA
 240 variants from these individuals and imputed the HLA region for the untyped
 241 individuals using HIBAG [32]. Phased genotype calls with *prob* >0.8 (i.e.
 242 conditional probability of pairs of haplotypes consistent with observed
 243 genotypes) were converted to amino acid variants and merged with the SNP
 244 variants for association analysis in FaST-LMM, as described above. Haplotype
 245 analyses were performed in PLINK [33] on phased haplotype data using logistic
 246 regression under an additive model with gender, age and 10 PCs as covariates.

247

248 **Functional Predictions for Candidate Loci**

249 We assessed the functional role for the candidate causal HLA variants *in silico*
 250 using NetMHCIIpan 3.1 [34] to map epitopes and their binding affinities to two
 251 risk and one protective HLA-DQA1_HLA-DQB1 haplotypes across GAS proteins
 252 known to contain human cross-reactive epitopes. A literature review of the GAS
 253 proteins reported to show cross-reactivity with host tissue proteins was

254 undertaken (Table S2). Full-length amino acid sequences of all GAS proteins,
255 including M5 and M6 proteins, shown to have cross-reactive epitopes were
256 converted to a series of 20-mer sequences with a 1-mer sliding window and
257 assessed for binding to each significantly associated DQA1_DQB1 haplotype.
258 Cross-reactive epitopes from human proteins were mapped onto the epitope
259 binding maps of M5 and M6, as indicated. Binding affinities were compared
260 (GraphPad Prism 7.00: 1-way ANOVA with multiple comparisons and correction
261 for multiple testing) between haplotypes across the regions of peak epitope
262 binding where 20-mer epitopes shared common 9-mer core epitopes.

263

264 **RESULTS**

265

266 **Genome-wide Association Study**

267 We conducted a GWAS for rheumatic heart disease (RHD) in 1263 individuals
268 comprising 398 RHD case and 865 control participants from communities in the
269 Northern Territory of Australia. From direct genotyping on the Illumina
270 HumanCoreExome array, we achieved 4.46 million high quality imputed variants
271 (92.33% of variants imputed to high accuracy, $r^2 > 0.80$) with moderate to high
272 imputation accuracy genome-wide (Figure S1A). Genetic population structure
273 was clearly evident from principal components analysis, largely capturing the
274 geographic distribution of the remote Aboriginal Australian communities (data
275 not shown). The use of a linear mixed model framework with genetic relatedness
276 matrix (FastLMM) to perform a genome-wide association study for RHD
277 effectively controlled this stratification, as evidenced by a quantile-quantile plot

278 of the p-values from the genome-wide scan ($\lambda = 1.021$; Figure S1B). A single
279 major signal was detected within the class II region of the HLA gene family on
280 chromosome 6 which peaked at the imputed variant rs9272622 (32607986bp,
281 $P=1.86 \times 10^{-7}$, OR=0.897 for protective allele C) within intron 1 of *HLA-DQA1*
282 (Figure 2).
283
284 **Fine-mapping the HLA Class II Region**
285 Regional plots of the class II region showed that the top SNP rs9272622 tagged a
286 linkage disequilibrium block ($r^2 > 0.8$) across the *HLA-DQA1* to *HLA-DQB1* region
287 (Figure 3). There were no residual signals across the HLA Class II region after
288 conditioning on the index variant rs9272622 (Figure S2). To further understand
289 the potential functional variants across the HLA Class II region, we typed and
290 imputed traditional 4-digit HLA alleles, converted alleles to amino acid calls, and
291 applied a multiple stepwise regression analysis. The top 4-digit HLA alleles for
292 risk and protection were *HLA-DQB1*0601* ($P=4.06 \times 10^{-4}$, OR=1.07) and *HLA-*
293 *DQA1*0301* ($P=2.71 \times 10^{-4}$, OR=0.92), respectively. The top 4-digit HLA-DRB1
294 association was HLA-DRB1*0803 ($P=0.005$, OR=1.06), and no significant
295 associations were observed for classical alleles across the SNP poor region
296 (Figure 3) of HLA-DRB3/DRB4/DRB5. The strongest amino acid associations
297 (Figure 4A) were at positions AA_DQA1_16_32713236 ($P=2.08 \times 10^{-6}$, OR=0.91)
298 and AA_DQA1_69_32717257_L ($P=2.08 \times 10^{-6}$, OR=0.91) in exons 1 and 2 of DQA1,
299 respectively, which were in 100% linkage disequilibrium with each other, and at
300 AA_DQB1_38_32740723 ($P=2.17 \times 10^{-6}$, OR=0.91) in exon 2 of DQB1. As when
301 conditioning on the top SNP (Figure 4B), there was no residual signal across the
302 HLA region when conditioning on either the top DQA1 AA variant (Figure 4C) or

both the top SNP and the top DQA1 AA variant (Figure 4D), suggesting that associations across the HLA-DQA1 to HLA-DQB1 region are all due to linkage disequilibrium with top variants at *HLA-DQA1*.

HLA-DQ haplotype risk

HLA-DQA1 and *HLA-DQB1* genes encode alpha and beta chains, respectively, forming DQ alpha/beta heterodimers that together bind antigenic epitopes to present to CD4+ T cells. For antigen presentation via HLA-DQ class II molecules, variation at both the alpha and beta chains contribute to epitope binding to the peptide groove encoded by exons 2 of both alpha and beta chains. Variants at both genes may thus contribute together to determine risk versus protection from RHD. We therefore looked for associations between RHD and *HLA-DQA1_HLA-DQB1* haplotypes. Haplotype analysis in PLINK identified *HLA-DQA1*0101_DQB1*0503* (OR 1.44, 95%CI 1.09-1.90, $P=9.56 \times 10^{-3}$) and *HLA-DQA1*0103_DQB1*0601* (OR 1.27, 95%CI 1.07-1.52, $P=7.15 \times 10^{-3}$) as risk haplotypes, with *HLA-DQA1*0301_DQB1*0402* (OR 0.30, 95%CI 0.14-0.65, $P=2.36 \times 10^{-3}$) as the protective haplotype for RHD in the study population (Figure 5). These haplotypes were taken forward in *in silico* functional analyses.

Mapping Group A Streptococcus epitopes to risk versus protective HLA DQ haplotypes

There are two important ways in which association between HLA-DQ haplotypes could impact on disease susceptibility and control programs: (i) in the pathogenesis of disease, particularly in relation to an autoimmune mechanism for RHD through GAS epitopes that cross-react with self; and (ii) in the ability of

high risk individuals to respond to proposed vaccine antigens. To address the first, we initially assessed the binding affinities of epitopes across the M-proteins M5 and M6 from rheumatogenic GAS strains to the alpha/beta heterodimers specific to the observed risk versus protective HLA-DQA-HLA-DQB haplotypes. Figure 6 shows the epitope binding affinities mapped for these haplotypes across the full-length M5 and M6 proteins, together with annotation indicating the positions along each protein where experimentally validated cross-reactive epitopes have been identified (Table S2). Several epitope peaks that correspond to key cross-reactive epitopes are shown to bind with higher affinity to the two risk haplotypes compared to the protective haplotype (Figure 6), notably in the B repeat regions previously shown to contain key cross-reactive T cell epitopes with human cardiac myosin (e.g. Cunningham et al., 1997 [10]; see also Table S2). The peak differences in binding affinities for 20-mer epitopes in these regions of previously experimentally validated cross reactivity for the M5 (see arrows, Figure 6A) and M6 (see arrows, Figure 6B) proteins were highly significant ($P < 0.0001$) between risk and protective haplotypes (Figure 7). No differences in epitope binding to risk versus protective haplotypes were observed when we mapped epitopes across GAS M proteins (e.g. E pattern M4 and M49 types [35]) from non-RHD GAS strains (Figure S3). Nor did we observe regions of differential epitope binding affinities across other GAS proteins (HSP70, STRP1; Figure S3) reported in the literature to contain epitopes cross reactive with human proteins implicated in RHD pathogenesis (Table S2).

350

Also annotated in Figure 6 are the C-terminal regions of the M5/M6 proteins that contain peptides incorporated into the two candidate vaccines currently in

353 advanced stages of development that include antigens from this M protein
354 region, J8-DT [36] (vertical blue strip) and StreptinCor [37] (vertical apricot
355 strip). Whilst the risk haplotype HLA-DQA1*0103_DQB1*0601 binds to epitopes
356 across this region with higher affinity, all three haplotypes show similar patterns
357 of epitope binding across this region. None show the low level of binding affinity
358 such as that observed for the protective haplotype for cross-reactive epitopes
359 across the B-repeat region. These results suggest that individuals genetically at
360 risk of developing RHD have the potential to make HLA-DQ-driven CD4+ T cell
361 responses to these vaccines.

362

363 **DISCUSSION**

364

365 The results of an unbiased genome-wide evaluation of genetic determinants for
366 RHD in Aboriginal Australians living in northern Australia provide evidence for a
367 prominent association in the class II gene region of HLA, consistent with prior
368 data from more limited genetic studies. Strong linkage disequilibrium across
369 HLA, together with variable selection of candidate HLA genes, likely contributes
370 to the inconsistency in the HLA genes/alleles associated with risk versus
371 protective from RHD in prior studies [6, 19, 20] even though experimental
372 studies support HLA-DQ restriction of T cell clones involved in T cell mimicry in
373 RHD [11]. In contrast, our study benefitted from dense fine mapping across HLA,
374 allowing us to identify specific risk (HLA-DQA1*0101_DQB1*0503; HLA-
375 DQA1*0103_DQB1*0601) *versus* protective (HLA-DQA1*0301_DQB1*0402)
376 haplotypes across the genes encoding alpha and beta chains of HLA-DQ. While
377 our conditional analysis suggested only a single HLA signal, we cannot discount

the possibility that other genes may contribute to genetic susceptibility to RHD in this population. It is of specific interest, however, that our study did not find evidence for replication for variants at the IGH locus recently shown to be significantly associated with RHD in a GWAS of New Caledonian and Fijian populations [21]. Differences in study design and phenotype classification may have contributed, as could genetic heterogeneity between indigenous populations which is known to occur for autoimmune and infectious diseases [38]. It is reassuring, nevertheless, that both GWAS have found evidence consistent with autoimmune genetic architecture. Ultimately, meta-analyses of greater statistical power will be required to investigate population-specific differences and detect additional RHD loci.

Our identification of risk versus protective haplotypes across HLA-DQA/DQB provided an opportunity to revisit the molecular mimicry hypothesis in relation to RHD pathogenesis. Dimers created from alpha and beta chains of HLA class II molecules present epitopes processed from foreign proteins to CD4 T cells, the preferred outcome of which would be to provide an immune response that will protect against infection. In the context of autoimmune disease, self-epitopes are presented and recognised as non-self, leading to detrimental immune pathology. The molecular mimicry hypothesis proposes that GAS contains proteins with AA sequences that mimic (or are cross-reactive with) human proteins, thus leading the immune system to recognise them as auto-antigens that drive immune pathology rather than (or in addition to) immunity against GAS itself.[1, 6] In the case of HLA-DQ, variation in exons 2 of both alpha and beta chains encoded by DQA and DQB, respectively, contribute to variation in shape and structure of the

403 epitope binding pocket [39]. This means that the specific alpha/beta dimers
 404 encoded by DQA/DQB genes carried on the same haplotype will create binding
 405 pockets that have different characteristics in terms of ability to bind and present
 406 epitopes to CD4+ T cells. Using the current gold standard NetMHCIIpan 3.1[34,
 407 40] predictive algorithm to map specific epitopes across GAS proteins allowed us
 408 to identify significant differences in the ability of dimers created from risk *versus*
 409 protective haplotypes to bind cross-reactive epitopes. In particular, cross-
 410 reactive epitopes from cardiac myosin, one of the key cardiac proteins thought to
 411 contribute to the molecular mimicry hypothesis in RHD [1, 6, 10], were predicted
 412 to bind to dimers created from risk haplotypes but have no predicted binding to
 413 dimers created from the protective haplotype. Thus we identify a potential
 414 molecular mechanism to account for immune pathogenesis causing RHD in this
 415 population. Although we carried out our epitope mapping studies on just two
 416 M5 and M6 GAS strains most studied for the presence of human cross-reactive
 417 epitopes, our results are relevant to all GAS strains carrying cross-reactive N-
 418 terminal or B repeats. Relevance to our study population is consistent with
 419 global-scale transmission of GAS strains in this remote Aboriginal population
 420 [15]. Of interest too is the observation that, whilst rare cases of dimers created
 421 by *trans* association of alpha/beta chains encoded on opposite strands of the
 422 chromosome have been observed to contribute to susceptibility to type 1
 423 diabetes, the predominant observation is that dimers are formed by alpha/beta
 424 chains encoded in *cis* [39]. This likely contributes to our ability to identify risk
 425 *versus* protective haplotypes across the HLA DQA1-DQB1 region, since strong
 426 linkage disequilibrium will keep particular combinations of DQA/DQB genes
 427 together in *cis*.

428

429 More broadly, this study represents a rare example of a genome-wide
430 association study in a remote Indigenous population, yet one which shows that
431 such studies can be successfully undertaken and uncover insights which have the
432 potential to inform pathogenesis and vaccination strategy.

433

434 In conclusion, we here present results of the first GWAS undertaken for RHD
435 in an Aboriginal Australian population. We report strong evidence for a role for
436 HLA DQ/DB Class II molecules, and we link this to significant differences in
437 affinity of binding of cross-reactive epitopes from GAS M proteins to antigen
438 presenting heterodimers formed by risk versus protective DQ-DB haplotypes.
439 Further functional analysis of T cell responses to cross-reactive T cell epitopes,
440 as carried out in previous studies [11], could now be targeted at these specific
441 DQ-DB heterodimers. Overall, our results provide new data on mechanisms that
442 may contribute to risk of RHD caused by GAS strains.

443

444 Main text: (3,621 words)

445

446 **Author contributions.** L-AG, HD'A, and SYCT contribute equally to the work.
 447 JMB, MII, and JRC contributed equally to supervision of the work. L-AG managed
 448 the data and carried out the genetic statistical and bioinformatic analyses, and
 449 prepared the first draft of the manuscript. HD'A and MMcK project managed in
 450 Darwin, including management of ethical, legal and social aspects of the study.
 451 MMcK carried out the field work and sample collection. DB and NB made
 452 significant contributions to governance and helped design the community engagement
 453 arms of the project. SYCT, BR and AS provided the major clinical inputs for
 454 diagnosis and review of patient records. GS prepared the DNAs including quality
 455 control, liaised with providers for both chip genotyping and sequence-based HLA
 456 typing. JMB and MII supervised the GWAS and HLA fine mapping analysis. JMB
 457 devised, supervised and interpreted the *in silico* analyses, and undertook major
 458 revisions of manuscript. JRC was the lead investigator on the project. All authors
 459 reviewed and approved the final manuscript.

460 **Acknowledgements.** We would like to acknowledge all Chief Investigators of
 461 this study, the project team including the community based researchers, the
 462 communities, agencies and all the participants for their invaluable contribution
 463 to this project. We also acknowledge the contribution of Paul I.W. de Bakker,
 464 now Vertex Pharmaceuticals, to the design and initial analysis of the study, and
 465 thank Kara Imbrogno and Grace Chua for assistance with preparation of some of
 466 the DNAs used for this study.

467 **Financial support.** National Health and Medical Research Council (NHMRC)
 468 Grant APP1023462, NHMRC and National Heart Foundation of Australia Career
 469 Development Fellow (#1061435), NHMRC Career Development Fellow
 470 (#1065736).

471 ***Potential conflicts of interest.*** The authors disclose no conflicts of interest. The
472 study sponsor had no role in study design, data collection, data analysis, data
473 interpretation, or writing of the report. The corresponding author had full
474 access to all study data and had final responsibility for the decision to submit for
475 publication.

476

477

REFERENCES

479

- 480 1. Carapetis JR, Beaton A, Cunningham MW, et al. Acute rheumatic fever and
481 rheumatic heart disease. *Nat Rev Dis Primers* **2016**; 2:15084.
- 482 2. Roberts KV, Maguire GP, Brown A, et al. Rheumatic heart disease in Indigenous
483 children in northern Australia: differences in prevalence and the challenges of
484 screening. *Med J Aust* **2015**; 203:221 e1-7.
- 485 3. Zuhlke LJ, Steer AC. Estimates of the global burden of rheumatic heart disease.
486 *Glob Heart* **2013**; 8:189-95.
- 487 4. Carapetis JR, Currie BJ. Mortality due to acute rheumatic fever and rheumatic
488 heart disease in the Northern Territory: a preventable cause of death in
489 aboriginal people. *Aust N Z J Public Health* **1999**; 23:159-63.
- 490 5. Guilherme L, Kalil J. Rheumatic fever and rheumatic heart disease: cellular
491 mechanisms leading autoimmune reactivity and disease. *J Clin Immunol* **2010**;
492 30:17-23.
- 493 6. Martin WJ, Steer AC, Smeesters PR, et al. Post-infectious group A streptococcal
494 autoimmune syndromes and the heart. *Autoimmun Rev* **2015**; 14:710-25.
- 495 7. Tandon R, Sharma M, Chandrashekhar Y, Kotb M, Yacoub MH, Narula J.
496 Revisiting the pathogenesis of rheumatic fever and carditis. *Nat Rev Cardiol*
497 **2013**; 10:171-7.
- 498 8. Cunningham MW. Rheumatic fever, autoimmunity, and molecular mimicry: the
499 streptococcal connection. *Int Rev Immunol* **2014**; 33:314-29.
- 500 9. Guilherme L, Kohler KF, Pommerantseff P, Spina G, Kalil J. Rheumatic Heart
501 Disease: Key Points on Valve Lesions Development. *J Clin Exp Cardiol* **2013**;
502 S:3:006.
- 503 10. Cunningham MW, Antone SM, Smart M, Liu R, Kosanke S. Molecular analysis
504 of human cardiac myosin-cross-reactive B- and T-cell epitopes of the group A
505 streptococcal M5 protein. *Infect Immun* **1997**; 65:3913-23.
- 506 11. Ellis NM, Li Y, Hildebrand W, Fischetti VA, Cunningham MW. T cell mimicry
507 and epitope specificity of cross-reactive T cell clones from rheumatic heart
508 disease. *J Immunol* **2005**; 175:5448-56.

509 12. Fae KC, da Silva DD, Oshiro SE, et al. Mimicry in recognition of cardiac myosin
510 peptides by heart-intralesional T cell clones from rheumatic heart disease.
511 JImmunol **2006**; 176:5662-70.

512 13. Kirvan CA, Galvin JE, Hilt S, Kosanke S, Cunningham MW. Identification of
513 streptococcal m-protein cardiopathogenic epitopes in experimental autoimmune
514 valvulitis. J Cardiovasc Transl Res **2014**; 7:172-81.

515 14. Bright PD, Mayosi BM, Martin WJ. An immunological perspective on
516 rheumatic heart disease pathogenesis: more questions than answers. Heart
517 **2016**; 102:1527-32.

518 15. Towers RJ, Carapetis JR, Currie BJ, et al. Extensive diversity of Streptococcus
519 pyogenes in a remote human population reflects global-scale transmission rather
520 than localised diversification. PLoS ONE **2013**; 8:e73851.

521 16. Williamson DA, Smeesters PR, Steer AC, et al. Comparative M-protein analysis
522 of Streptococcus pyogenes from pharyngitis and skin infections in New Zealand:
523 Implications for vaccine development. BMC Infect Dis **2016**; 16:561.

524 17. Carapetis JR, Currie BJ, Mathews JD. Cumulative incidence of rheumatic fever
525 in an endemic region: a guide to the susceptibility of the population? Epidemiol
526 Infect **2000**; 124:239-44.

527 18. Engel ME, Stander R, Vogel J, Adeyemo AA, Mayosi BM. Genetic susceptibility
528 to acute rheumatic fever: a systematic review and meta-analysis of twin studies.
529 PLoS ONE **2011**; 6:e25326.

530 19. Anastasiou-Nana MI, Anderson JL, Carlquist JF, Nanas JN. HLA-DR typing and
531 lymphocyte subset evaluation in rheumatic heart disease: a search for immune
532 response factors. Am Heart J **1986**; 112:992-7.

533 20. Hafez M, Chakravarti A, el-Shennawy F, el-Morsi Z, el-Sallab SH, Al-Tonbary Y.
534 HLA antigens and acute rheumatic fever: evidence for a recessive susceptibility
535 gene linked to HLA. Genet Epidemiol **1985**; 2:273-82.

536 21. Parks T, Mirabel MM, Kado J, et al. Association between a common
537 immunoglobulin heavy chain allele and rheumatic heart disease risk in Oceania.
538 Nat Commun **2017**; 8:14946.

539 22. Bessarab D, Read C, D'Antoine H, et al. A proposed framework for
540 engagement, informed consent, and governance for conducting genetic research
541 with Indigenous communities. J Med Ethics **2017**; submitted.

542 23. Gewitz MH, Baltimore RS, Tani LY, et al. Revision of the Jones Criteria for the
543 diagnosis of acute rheumatic fever in the era of Doppler echocardiography: a
544 scientific statement from the American Heart Association. *Circulation* **2015**;
545 131:1806-18.

546 24. Remenyi B, Wilson N, Steer A, et al. World Heart Federation criteria for
547 echocardiographic diagnosis of rheumatic heart disease--an evidence-based
548 guideline. *Nat Rev Cardiol* **2012**; 9:297-309.

549 25. Abraham G, Inouye M. Fast principal component analysis of large-scale
550 genome-wide data. *PLoS ONE* **2014**; 9:e93766.

551 26. Genomes Project C, Auton A, Brooks LD, et al. A global reference for human
552 genetic variation. *Nature* **2015**; 526:68-74.

553 27. Delaneau O, Marchini J, Zagury JF. A linear complexity phasing method for
554 thousands of genomes. *Nat Methods* **2011**; 9:179-81.

555 28. Howie BN, Donnelly P, Marchini J. A flexible and accurate genotype
556 imputation method for the next generation of genome-wide association studies.
557 *PLoS Genet* **2009**; 5:e1000529.

558 29. Lippert C, Listgarten J, Liu Y, Kadie CM, Davidson RI, Heckerman D. FaST
559 linear mixed models for genome-wide association studies. *Nat Methods* **2011**;
560 8:833-5.

561 30. Pe'er I, Yelensky R, Altshuler D, Daly MJ. Estimation of the multiple testing
562 burden for genomewide association studies of nearly all common variants. *Genet*
563 *Epidemiol* **2008**; 32:381-5.

564 31. Raychaudhuri S, Sandor C, Stahl EA, et al. Five amino acids in three HLA
565 proteins explain most of the association between MHC and seropositive
566 rheumatoid arthritis. *Nat Genet* **2012**; 44:291-6.

567 32. Zheng X, Shen J, Cox C, et al. HIBAG--HLA genotype imputation with attribute
568 bagging. *Pharmacogenomics J* **2014**; 14:192-200.

569 33. Purcell S, Neale B, Todd-Brown K, et al. PLINK: a tool set for whole-genome
570 association and population-based linkage analyses. *Am J Hum Genet* **2007**;
571 81:559-75.

572 34. Andreatta M, Karosiene E, Rasmussen M, Stryhn A, Buus S, Nielsen M.
573 Accurate pan-specific prediction of peptide-MHC class II binding affinity with
574 improved binding core identification. *Immunogenetics* **2015**; 67:641-50.

575 35. Smeesters PR, McMillan DJ, Sriprakash KS. The streptococcal M protein: a
576 highly versatile molecule. Trends Microbiol **2010**; 18:275-82.

577 36. Batzloff MR, Hayman WA, Davies MR, et al. Protection against group A
578 streptococcus by immunization with J8-diphtheria toxoid: contribution of J8- and
579 diphtheria toxoid-specific antibodies to protection. J Infect Dis **2003**; 187:1598-
580 608.

581 37. Guerino MT, Postol E, Demarchi LM, et al. HLA class II transgenic mice
582 develop a safe and long lasting immune response against StreptInCor, an anti-
583 group A streptococcus vaccine candidate. Vaccine **2011**; 29:8250-6.

584 38. Ramos PS, Shedlock AM, Langefeld CD. Genetics of autoimmune diseases:
585 insights from population genetics. J Hum Genet **2015**; 60:657-64.

586 39. Tollefsen S, Hotta K, Chen X, et al. Structural and functional studies of trans-
587 encoded HLA-DQ2.3 (DQA1*03:01/DQB1*02:01) protein molecule. J Biol Chem
588 **2012**; 287:13611-9.

589 40. Zhang L, Udaka K, Mamitsuka H, Zhu S. Toward more accurate pan-specific
590 MHC-peptide binding prediction: a review of current methods and tools. Brief
591 Bioinform **2012**; 13:350-64.

592

593

594 **Figure legends** 595

596 **Figure 1.** Locations of study populations. Locations are given by latitude and
597 longitude for 19 Aboriginal communities in the Northern Territory of Australia that
598 participated in the study. Each dot indicates a single community, with wedges
599 indicating the proportion of case (filled in wedge) compared to control (open
600 wedges) samples for each population.

601
602 **Figure 2.** Manhattan plot of GWAS results for the 4.46M high quality 1000G imputed
603 SNP variants. Data are for analysis in FastLMM looking for association between SNPs
604 and RHD. The top SNP rs9272622 occurred within the HLA region on Chromosome
605 6p21, as shown.

606
607 **Figure 3.** LocusZoom plot of SNP associations with RHD across the Class II region of
608 the HLA complex. The $-\log_{10}P$ values (left Y-axis) are shown in the upper section of
609 the plot. Dots representing individual SNPs are coloured (see key) based on their
610 linkage disequilibrium r^2 with the top SNP rs9272622. The right Y-axis is for
611 recombination rate (blue line), based on HapMap data. The bottom section of the
612 plot shows the positions of genes across the region. For clarity, 5 genes were
613 removed upstream of 32.8Mb (PSMB8-9, HLA-DOA, LOC100507463, LOC100294145).

614
615 **Figure 4.** Plots of association between RHD and imputed classical 4-digit and amino
616 acid (AA) HLA alleles. Results are for association analyses in FastLMM: (A) without
617 conditioning; (B) after conditioning on the top SNP rs9272622; (C) after conditioning

on the top AA variant at DQA1 AA position 16; and (D) after conditioning on both of these variants.

620

Figure 5. Forest plot showing associations between RHD and phased HLA DQ_DB haplotypes. The plot show odds ratios (OR) and 95% confidence intervals for two risk (OR>1) and one protective (OR<1) haplotypes. Information to the right of the plot shows values for the OR, the haplotype frequency (HF), and the p-value for the haplotype association.

626

Figure 6. Plots showing binding affinities for predicted epitopes of GAS M proteins recognised by HLA DQ-DB heterodimers. Epitope binding predictions were performed in NetMHCIIpan3.1. The legend between parts (A) and (B) of the figure applies to both parts. The y-axis shows the relative binding affinity (expressed as $1 - \log_{50,000}$ of the nM binding affinity) for heterodimers formed from risk (red, brown) and protective (blue) DQ_DB haplotypes (see legend); the x-axis indicates the amino acid sequence locations for mature proteins, also equivalent to the start position of overlapping 20mers (1-mer sliding window) in (A) the GAS M5 sequence (Accession number CAM31002.1) and (B) the GAS M6 sequence (Accession number AAA26920.1). Horizontal dotted lines show different nM binding affinities. Negative binding affinity is indicated at >10,000 nM (i.e. below the red dotted line). Vertical arrows indicate the N-terminal or B-repeat cross-reactive epitopes used to compare binding affinities in Figure 7. The linear positions of known cross-reactive epitopes with human cardiac myosin and/or human heart valve tissue are shown in green; black lines indicate the regions of known experimentally determined human T cell

642 epitopes (see Table S2). The apricot and pale blue vertical bars indicate the positions
643 of C-repeat region peptides incorporated into the StreptinCor and J8-DT vaccines,
644 respectively.

645

646 **Figure 7.** Mean binding affinities for GAS M protein epitopes cross-reactive with
647 human cardiac myosin. (A) The y-axis (as for Figure 6) shows mean plus SD for
648 predicted M5 and M6 GAS protein epitopes (NT and B repeat regions; as annotated
649 with arrows in Figure 6) recognised by risk (red and brown bars) or protective (blue
650 bar) DQ-DB heterodimers formed from DQA1_DQB1 haplotypes, as labelled. ****
651 indicates $P < 0.0001$. (B) Shows the 20-mer epitope at the peak of the differences for
652 binding affinity of risk versus protective haplotypes, together with the predicted 9-
653 mer cores for each haplotype.

654

655

656 Online data supplements

657

658 **Figure S1.** (A) Imputation accuracy measured as average dosage R^2 for 235,942
659 type 2 variants filtered for an information metric >0.4 and genotype probability
660 call > 0.90 . Data shown separately by chromosome and for different minor allele
661 frequency (MAF) ranges, as indicated in the key. (B) Quantile-quantile plot of
662 GWAS p-values.

663

664 **Figure S2.** Manhattan plots of GWAS results for the 4.46M high quality 1000G
665 imputed SNP variants after conditioning on the top SNP rs9272622. Data are for
666 analysis in FastLMM looking for association between SNPs and RHD.

667

668 **Figure S3.** Plots showing binding affinities for predicted epitopes of GAS M
669 proteins recognised by HLA DQ-DB heterodimers. Epitope binding predictions
670 were performed in NetMHCIIpan3.1. The y-axis shows the relative binding
671 affinity (expressed as $1-\log_{50,000}$ of the nM binding affinity) for heterodimers
672 formed from risk (red, brown) and protective (blue) DQ_DB haplotypes (see
673 key); the x-axis indicates the amino acid sequence locations for mature proteins,
674 also equivalent to the start position of overlapping 20mers (1-mer sliding
675 window) in non-rheumatogenic GAS M4 (Accession number CAA33269) and
676 M49 (Accession number AAA26868.1) sequences, GAS HSP70 (Accession
677 number AAB39223.1) and GAS STRP1 (Accession number AAA26987.1)
678 sequences. Horizontal lines indicate 500nM (upper) and 1000nM (lower)
679 binding affinities.

680

681 **Table S1.** Basic demographic details (by gender, age at collection) for the 396
682 cases and 867 controls that passed all QC and were used in the GWAS analysis.

683

684 **Table S2.** Summary of experimentally confirmed published epitopes for (A) GAS
685 M proteins for which there is evidence of cross-reaction with human heart-
686 related proteins. (B) GAS M proteins for which there is evidence of non-cross-
687 reactive T- and B-epitopes, and (C) other GAS proteins with evidence of cross-
688 reactive epitopes.

689

690

Figure 1

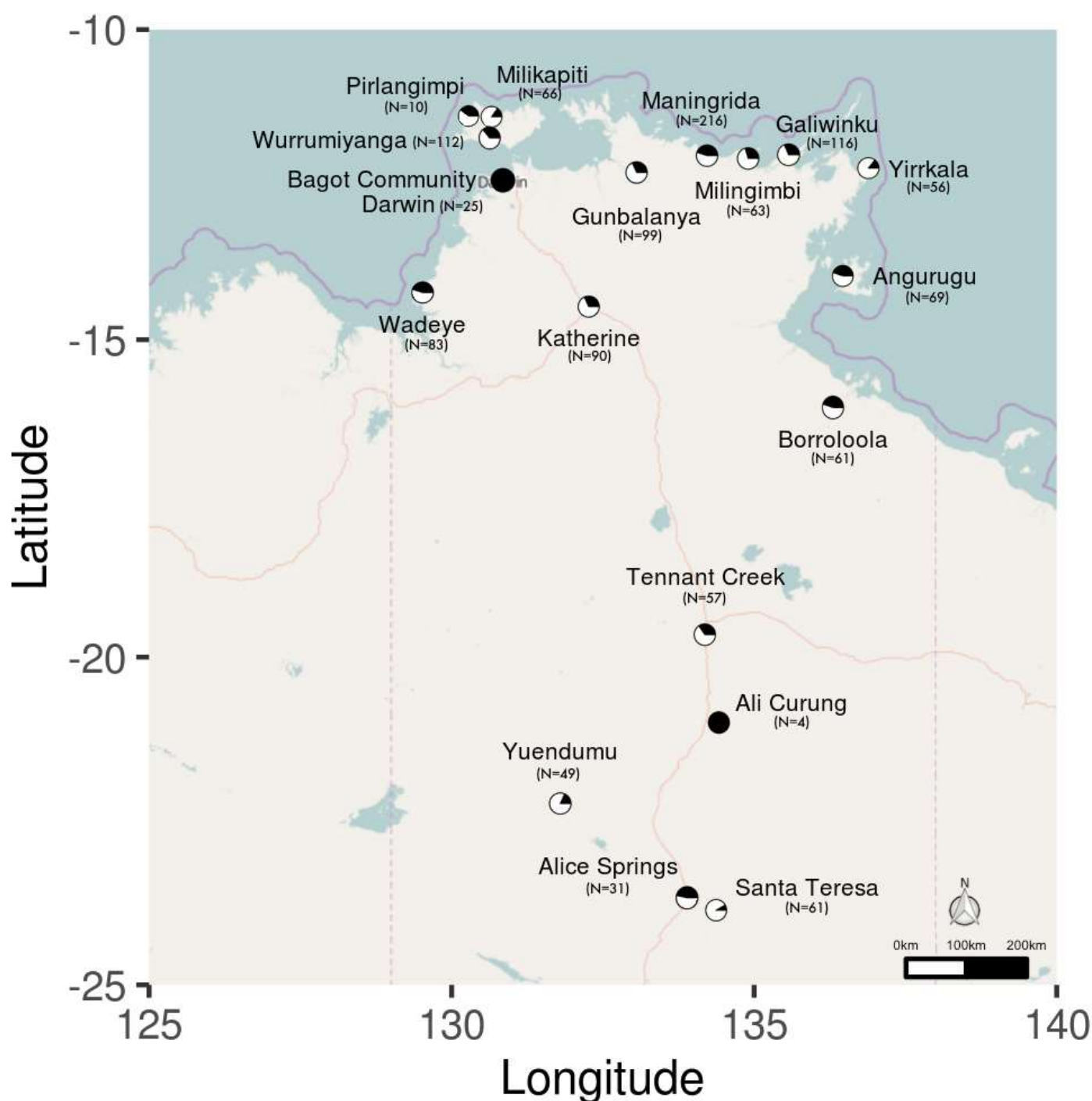


Figure 2

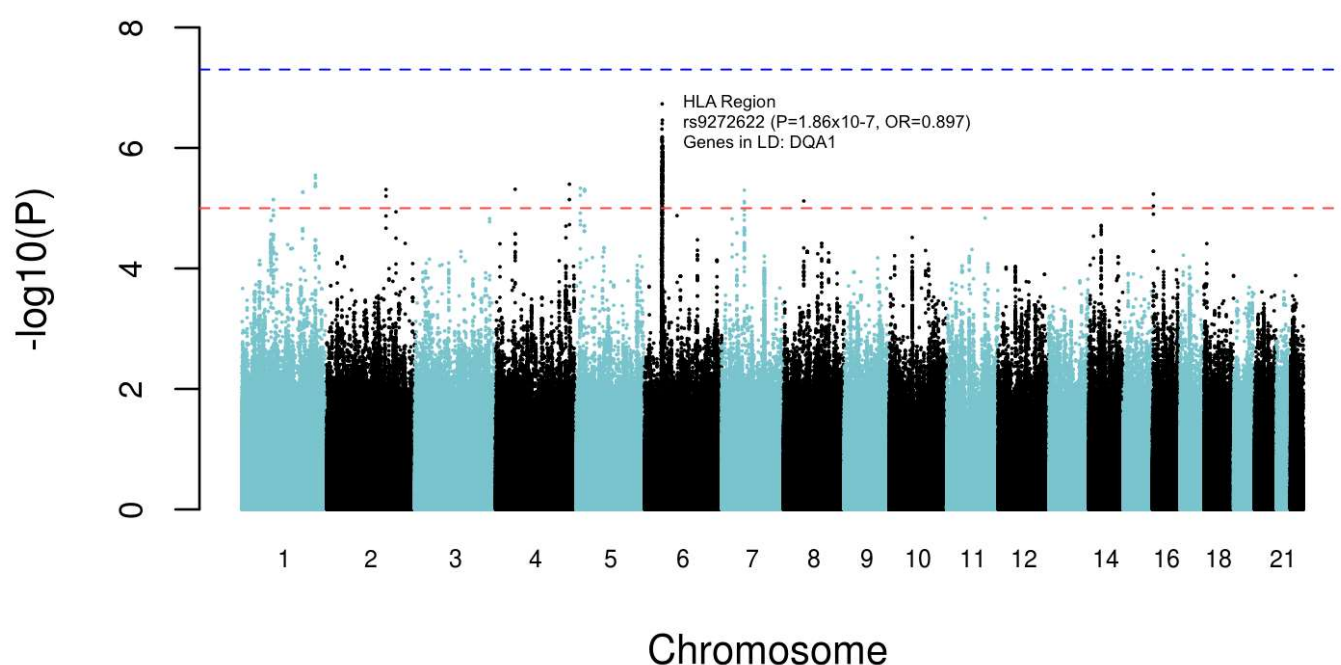


Figure 2. Manhattan plot of GWAS results for the 4.46M high quality 1000G imputed SNP variants. Data are for analysis in FastLMM looking for association between SNPs and RHD. The top SNP rs9272622 occurred within the HLA region on Chromosome 6p21, as shown.

Figure 3

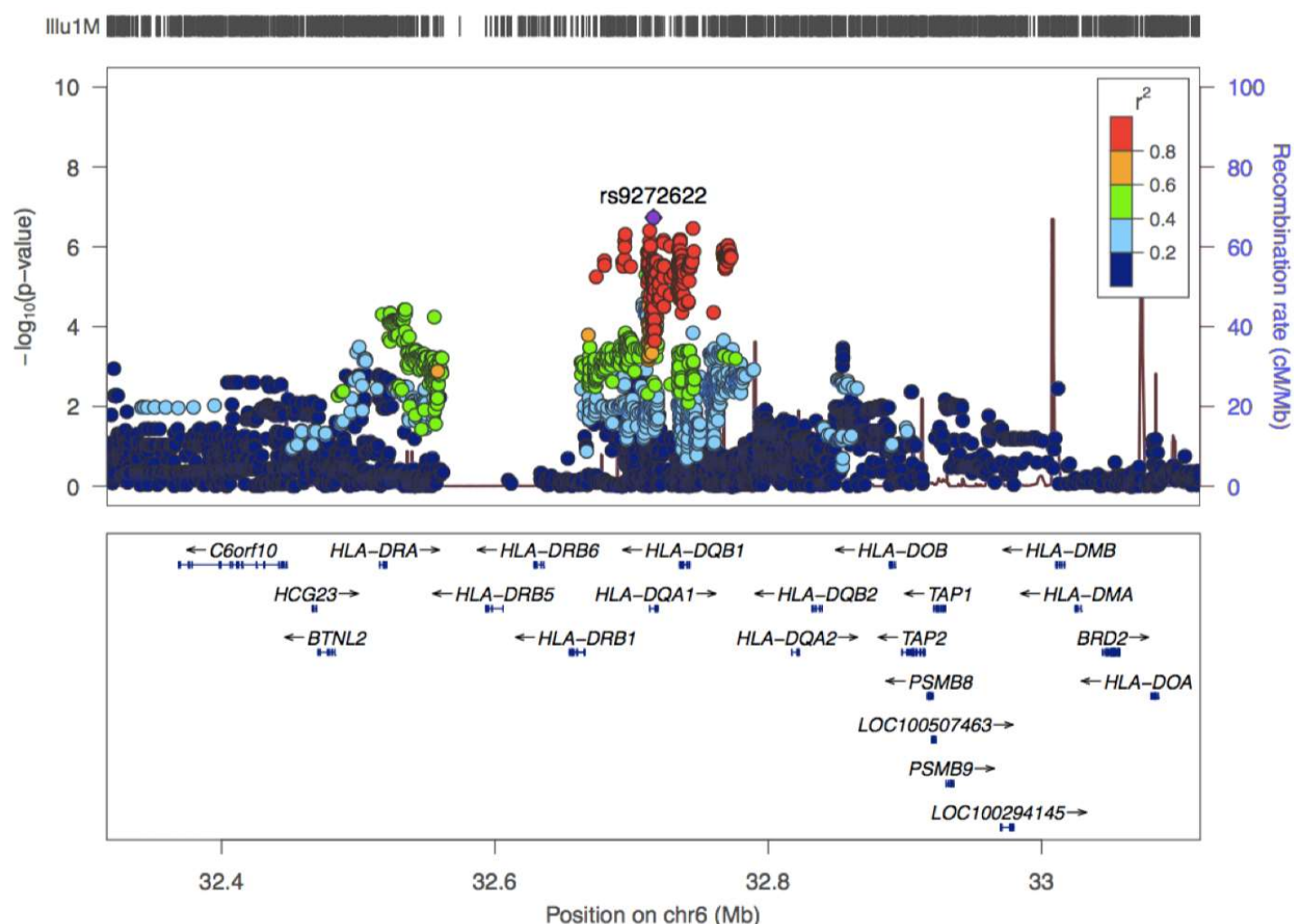


Figure 3. LocusZoom plot of SNP associations with RHD across the Class II region of the HLA complex. The $-\log_{10}P$ values (left Y-axis) are shown in the upper section of the plot. Dots representing individual SNPs are coloured (see key) based on their linkage disequilibrium r^2 with the top SNP rs9272622. The right Y-axis is for recombination rate (blue line), based on HapMap data. The bottom section of the plot shows the positions of genes across the region. For clarity, 5 genes were removed upstream of 32.8Mb (PSMB8-9, HLA-DOA, LOC100507463, LOC100294145).

Figure 4

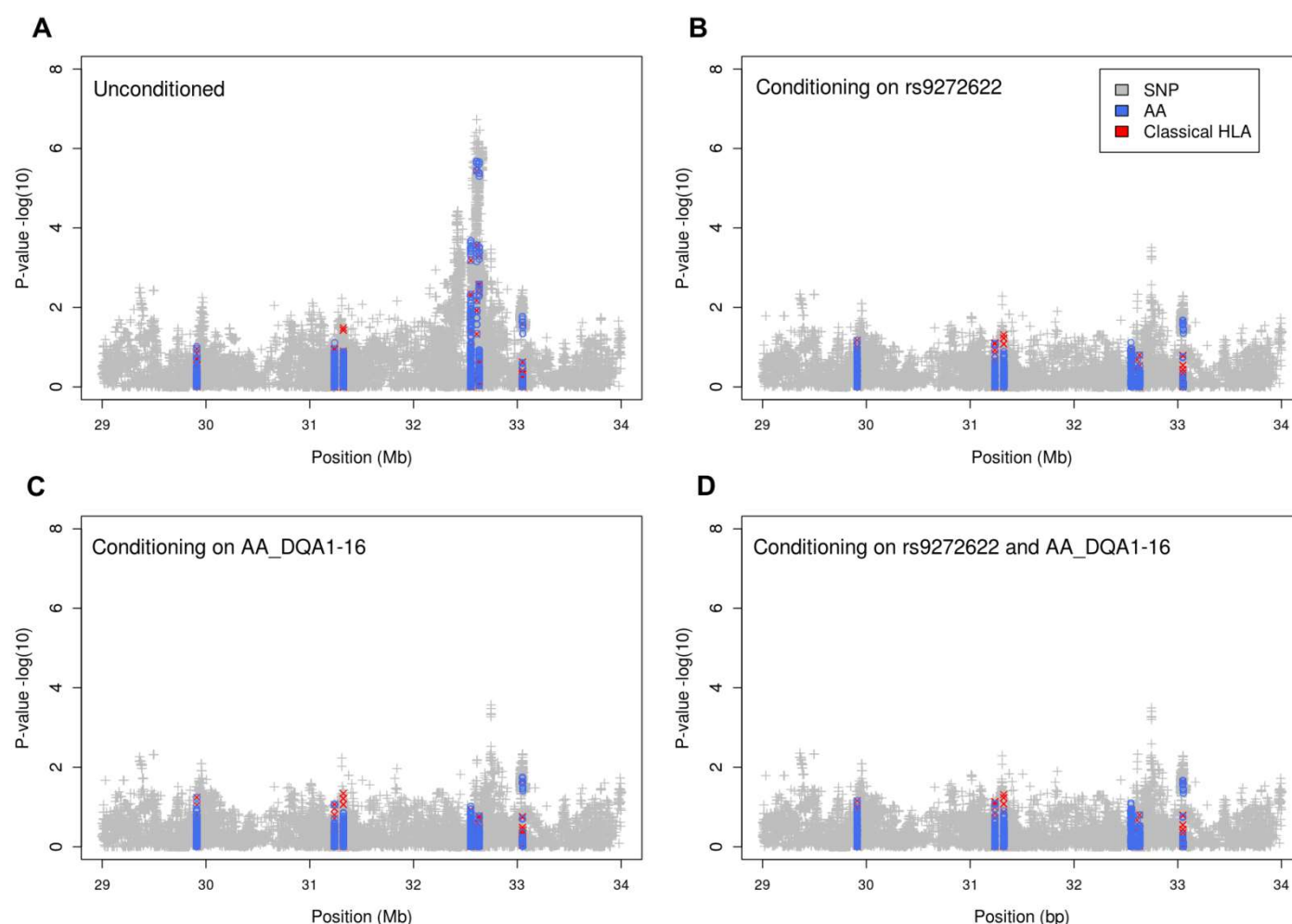


Figure 4. Plots of association between RHD and imputed classical 4-digit and amino acid (AA) HLA alleles. Results are for association analyses in FastLMM: (A) without conditioning; (B) after conditioning on the top SNP rs9272622; (C) after conditioning on the top AA variant at DQA1 AA position 16; and (D) after conditioning on both of these variants.

Figure 5

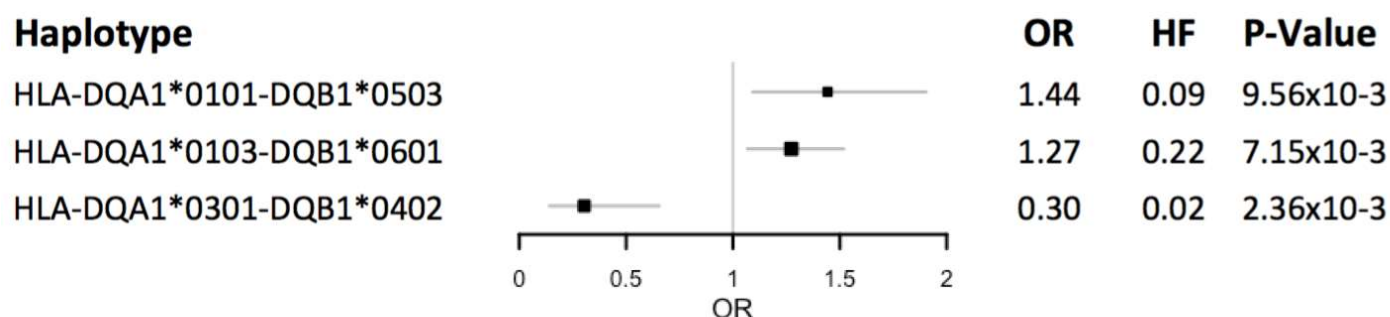


Figure 5. Forest plot showing associations between RHD and phased HLA DQ_DB haplotypes. The plot show odds ratios (OR) and 95% confidence intervals for two risk (OR>1) and one protective (OR<1) haplotypes. Information to the right of the plot shows values for the OR, the haplotype frequency (HF), and the p-value for the haplotype association.

Figure 6

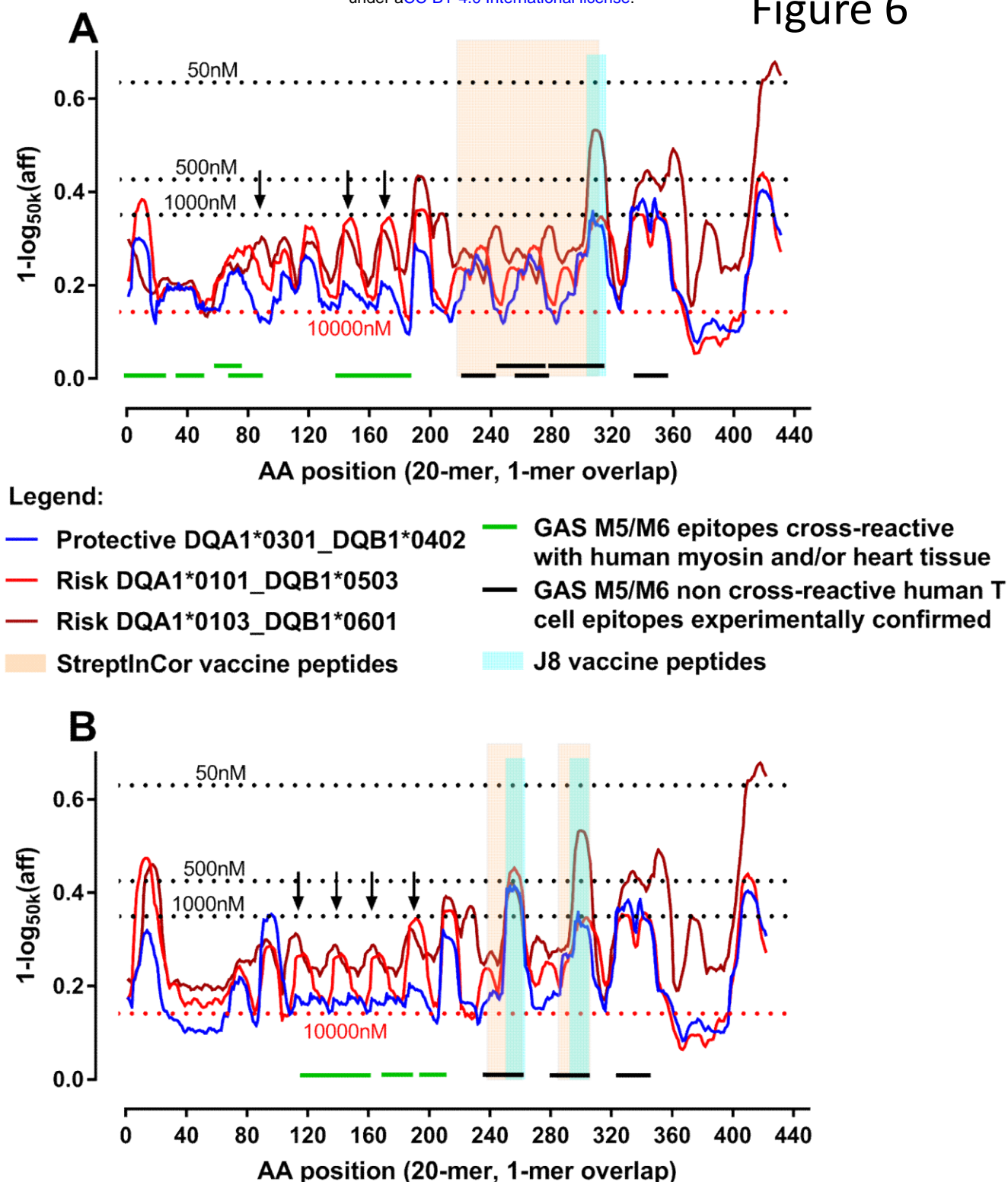


Figure 6. Plots showing binding affinities for predicted epitopes of GAS M proteins recognised by HLA DQ-DB heterodimers. Epitope binding predictions were performed in NetMHCIIpan3.1. The y-axis shows the relative binding affinity (expressed as $1-\log_{50,000}$ of the nM binding affinity) for heterodimers formed from risk (red, brown) and protective (blue) DQ-DB haplotypes (see key); the x-axis indicates the amino acid sequence locations for mature proteins, also equivalent to the start position of overlapping 20mers (1-mer sliding window) in (A) the GAS M5 sequence and (B) the GAS M6 sequence. Horizontal dotted lines show different nM binding affinities. Negative binding affinity is indicated at >10,000 nM (i.e. below the red dotted line). Vertical arrows indicate the N-terminal or B-repeat cross-reactive epitopes used to compare binding affinities in Figure 7. The linear positions of known cross-reactive epitopes with human cardiac myosin and/or human heart valve tissue are shown in green; black lines indicate the regions of known experimentally determined human T cell epitopes (see Table S2). The apricot and pale blue vertical bars indicate the positions of C-repeat region peptides incorporated into the J8-DT and StreptInCor vaccines, respectively.

Figure 7

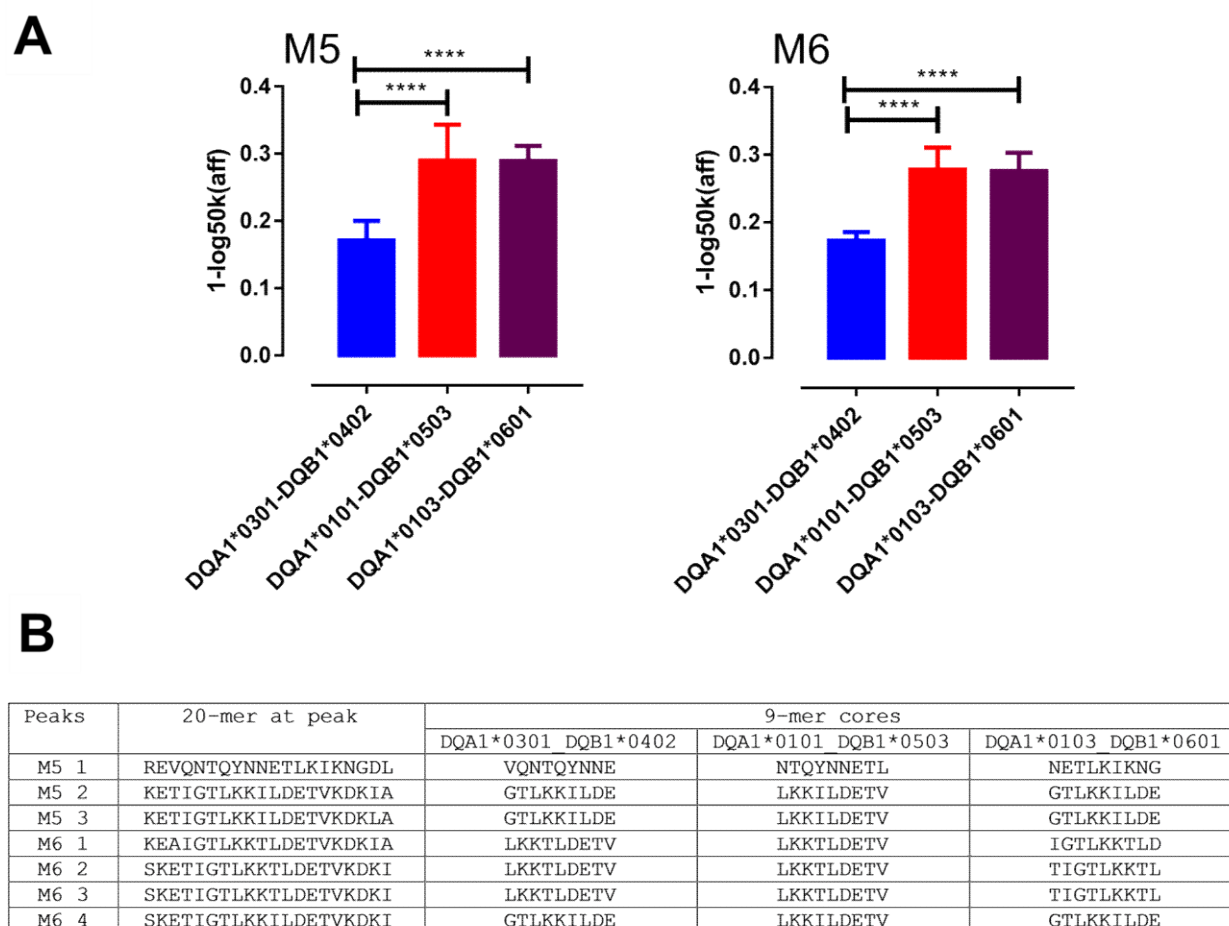


Figure 7. Mean binding affinities for GAS M protein epitopes cross-reactive with human cardiac myosin. (A) The y-axis (as for Figure 6) shows mean plus SD for predicted M5 and M6 GAS protein epitopes (NT and B repeat regions; as annotated with arrows in Figure 6) recognised by risk (red and brown bars) or protective (blue bar) DQ-DB heterodimers formed from DQA1_DQB1 haplotypes, as labelled. **** indicated $P < 0.0001$. (B) Shows the 20-mer epitope at the peak of the differences for binding affinity of risk versus protective haplotypes, together with the predicted 9-mer cores for each haplotype.

Figure S1

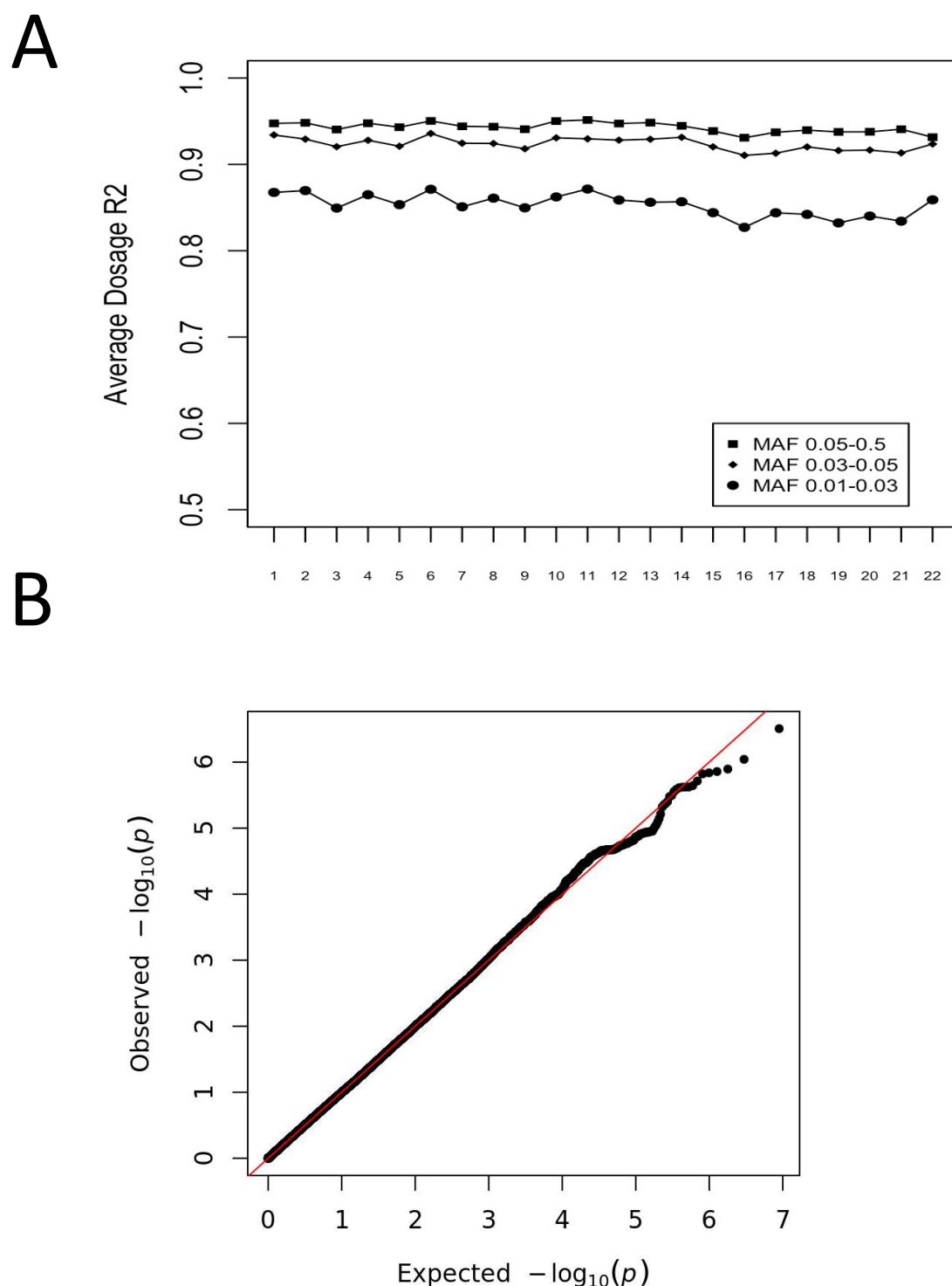


Figure S1. (A) Imputation accuracy measured as average dosage R^2 for 235,942 type 2 variants filtered for an information metric >0.4 and genotype probability call > 0.90 . Data shown separately by chromosome and for different minor allele frequency (MAF) ranges, as indicated in the key. (B) Quantile-quantile plot of GWAS p-values.

Figure S2

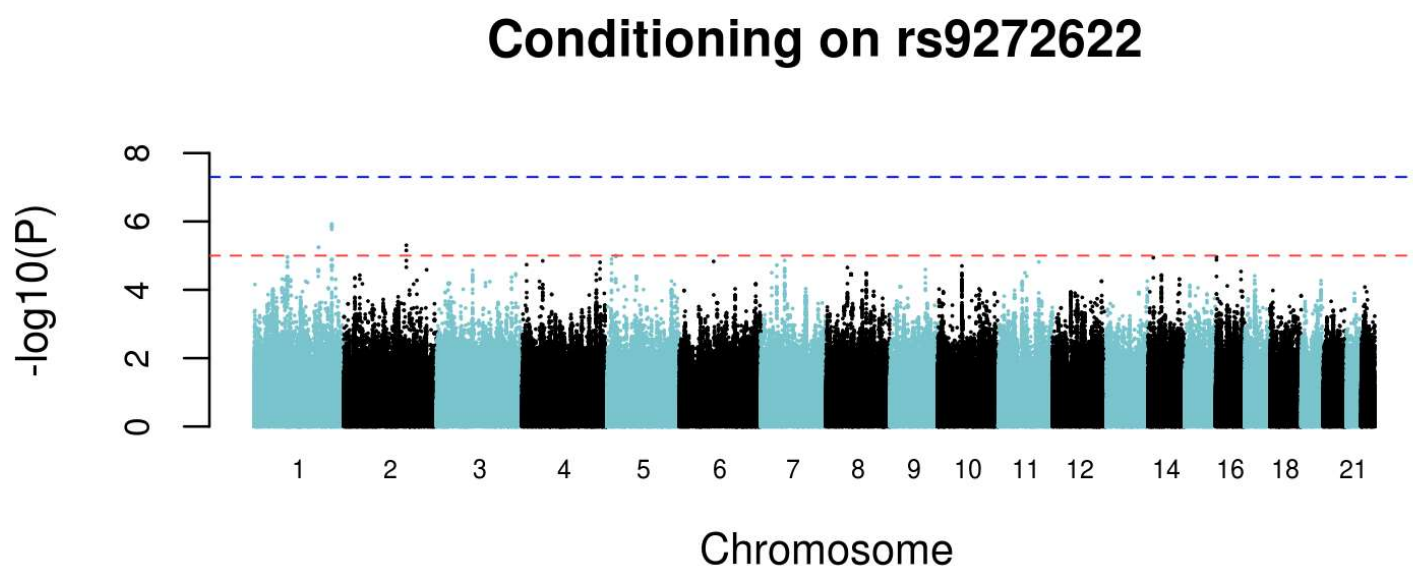


Figure S2. Manhattan plots of GWAS results for the 4.46M high quality 1000G imputed SNP variants after conditioning on the top SNP rs9272622. Data are for analysis in FastLMM looking for association between SNPs and RHD.

Figure S3

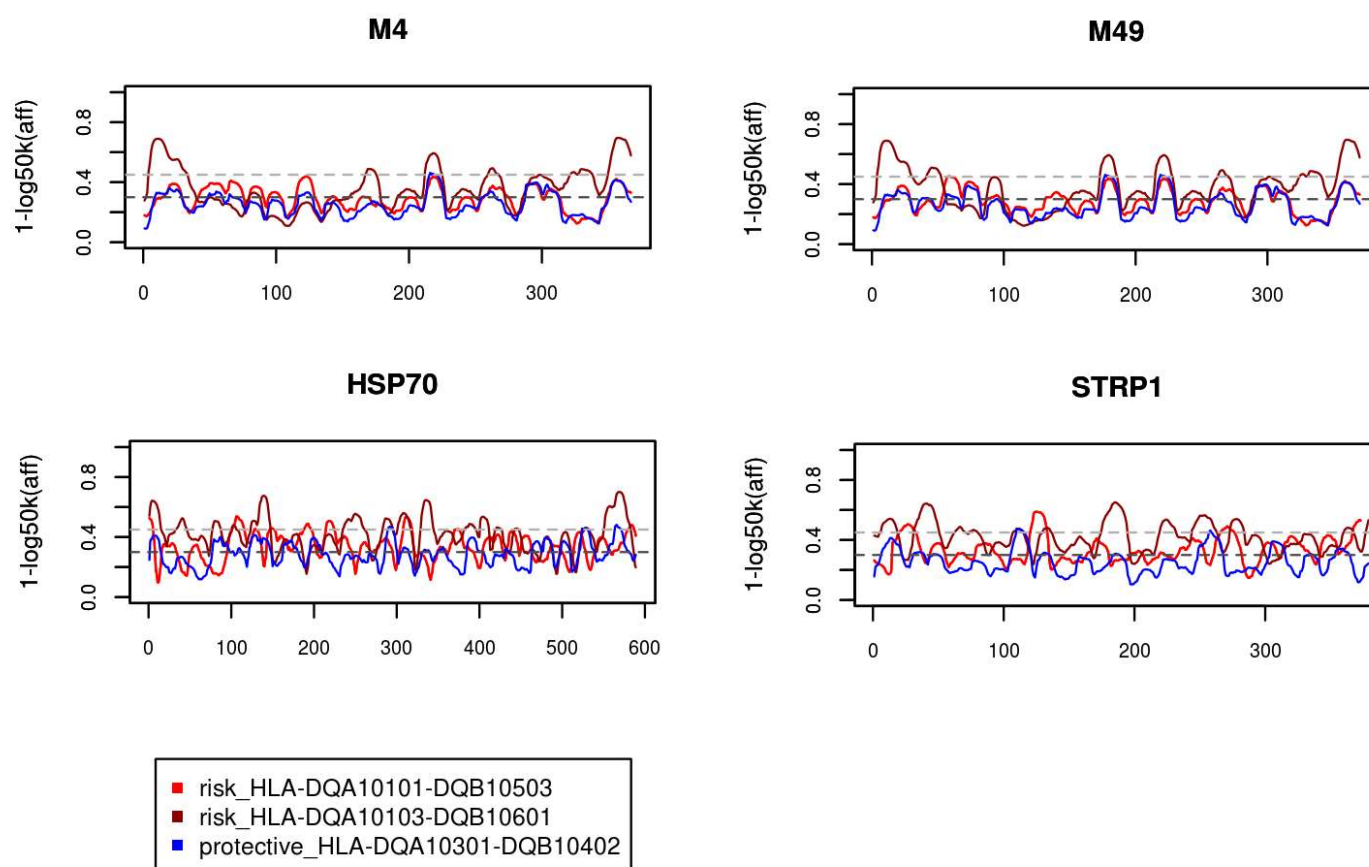


Figure S3. Plots showing binding affinities for predicted epitopes of GAS M proteins recognised by HLA DQ-DB heterodimers. Epitope binding predictions were performed in NetMHCIIpan3.1. The y-axis shows the relative binding affinity (expressed as $1-\log_{50,000}$ of the nM binding affinity) for heterodimers formed from risk (red, brown) and protective (blue) DQ_DB haplotypes (see key); the x-axis indicates the amino acid sequence locations for mature proteins, also equivalent to the start position of overlapping 20mers (1-mer sliding window) in non-rheumatogenic GAS M4 (Accession number CAA33269) and M49 (Accession number AAA26868.1) sequences, GAS HSP70 (Accession number AAB39223.1) and GAS STRP1 (Accession number AAA26987.1) sequences. Horizontal lines indicate 500nM (upper) and 1000nM (lower) binding affinities.

Table S1. Basic demographic details (by gender, age at collection) for the 396 cases and 867 controls that passed all QC and were used in the GWAS analysis.

	RHD Cases	Healthy Controls
Sample size		
Males	130	353
Females	266	514
Total	396	867
Age at collection		
Males		
Mean±SD	37.41±12.80	38.21±12.84
Range	18-74	18.2-83.8
Females		
Mean±SD	37.4±12.94	39.97±13.38
Range	18.3-82.9	18-73.8
Total		
Mean±SD	37.40±12.88	39.25±13.18
Range	18-82.9	18-83.8

Table S2. Summary of experimentally confirmed published epitopes for (A) GAS M proteins for which there is evidence of cross-reaction with human heart-related proteins. (B) GAS M proteins for which there is evidence of non-cross-reactive T- and B-epitopes, and (C) other GAS proteins with evidence of cross-reactive epitopes.

Function/peptide name*	Position **	Peptide Studied	AA Location M5	AA Location M6	Reference
(A) GAS M5 and M6 epitopes cross-reactive with human myosin and/or heart valve tissue					
Valve cross-reactive	NT	AVTRGTISDPQRAKEALDKYELENH	1-25	-	1-4
Myosin cross-reactive NT4	NT	GLKTENEGLKTENEGLKTE	33-51	-	5,6
Myosin cross-reactive NT5	NT	KKEHEAENDKLKQQRDTL	59-76	-	5
Valve cross-reactive	NT	DKLKQQRDTLSTQKET LKQQRDTLSTQKETLEREVQN STQKETLEREVQN	67-82 69-89 77-89	- - -	3,4,7
Myosin cross-reactive	NT	QRDTLSTQKETLEREVQN	72-89	-	1-3,5
Myosin cross-reactive B.6 B repeat (B1-B2)	B	TVKDKIAKEQENKETIGTLK VKDKIAKEQENKETIGTL	136-155 161-180	130-149 155-174 180-199	8 5
Myosin cross-reactive B repeat (B1-B2)	B	ETIGTLKKILDETVK	149-163 174-188	143-157 168-182 193-207	4
Myosin cross-reactive B repeat (B2, B3A)	B	TIGTLKKILDETVKDKIA IGTLKKILDETVKDKLAK	150-167 176-193	144-162 169-88; 194-211	9 5
(B) Vaccine candidates and related epitopes (not cross-reactive)					
StreptInCor Vaccine	C	KGLRRDLASREAKKQLEAEQQKLEEQNKISEASRKGLRRDLASREAKKQVEKA KGLRRDLASREAKKQLEAEQQKLEEQNKISEASRKGLRRDLASREAKKQVEKA	223-242; 244-273; 258-277; 279-312	237-260 285-302	10-12
P145 minimal T epitope	C	RDLASREAKKQ	227-238; 262-273; 298-308	246-257 288-299	13

Human responses India; C-repeat J14	C	KQAEDKVK ASREAKKQ VEKALEQLEDKRVK	231-238; 301-314	249-261 292-305	¹⁴
T cell (C1-A)	C	NKISEASRKGLRRDLASRE	250-269 285-304	234-253 281-295	⁸
C-term p145	C	LRRDLASREAKKQVEKALE	224-237 295-314	243-260 286-305	^{13,15-18}
P145 minimal B epitope	C	ASREAKKQVEKALE	231-238 301-314	249-260 293-305	¹³
J8 vaccine peptide	C	QAEDKVKQ SREAKKQVEK ALKQLEDKVQ	302-313	251-260 293-304	¹⁹
C-term T and B cell	C	KLTEKEKAELQAKLEAEAKA	335-354	325-345	¹⁸
(C) Other GAS proteins with cross-reactive epitopes					
STRP1 streptopain x vimentin		KKKLGVRLLSLA	3-15		²⁰
HSP70 x vimentin		AYFNDAQRQATKDA	118-131		²⁰

Note 1: Bold and grey highlights indicate the region of the peptide epitope(s) that matches the AA location for the M5 (Accession Number CAM31002.1) and M6 (Accession Number AAA26920.1) protein sequences used in our epitopes mapping studies. AA locations are for the mature protein sequence (i.e. after removal of the signal peptide). * Peptide name from the relevant publication; **NT=N terminal region; B= B repeat region; C = C repeat region.

Note 2: Many of these epitopes have been worked on in murine and rat²¹ models of disease, as recently reviewed,²² which are not all referenced in this table.

References

1. Guilherme, L. *et al.* T cell response in rheumatic fever: crossreactivity between streptococcal M protein peptides and heart tissue proteins. *Curr Protein Pept Sci* **8**, 39-44 (2007).
2. Guilherme, L. & Kalil, J. Rheumatic fever: from innate to acquired immune response. *Ann N Y Acad Sci* **1107**, 426-33 (2007).
3. Guilherme, L., Ramasawmy, R. & Kalil, J. Rheumatic fever and rheumatic heart disease: genetics and pathogenesis. *Scand J Immunol* **66**, 199-207 (2007).
4. Fae, K.C. *et al.* Mimicry in recognition of cardiac myosin peptides by heart-intralesional T cell clones from rheumatic heart disease. *J Immunol* **176**, 5662-70 (2006).
5. Cunningham, M.W., Antone, S.M., Smart, M., Liu, R. & Kosanke, S. Molecular analysis of human cardiac myosin-cross-reactive B- and T-cell epitopes of the group A streptococcal M5 protein. *Infect Immun* **65**, 3913-23 (1997).

6. Cunningham, M.W. *et al.* Human and murine antibodies cross-reactive with streptococcal M protein and myosin recognize the sequence GLN-LYS-SER-LYS-GLN in M protein. *J Immunol* **143**, 2677-83 (1989).
7. Guilherme, L. *et al.* Human heart-infiltrating T-cell clones from rheumatic heart disease patients recognize both streptococcal and cardiac proteins. *Circulation* **92**, 415-20 (1995).
8. Gorton, D., Govan, B., Olive, C. & Ketheesan, N. B- and T-cell responses in group A streptococcus M-protein- or Peptide-induced experimental carditis. *Infect Immun* **77**, 2177-83 (2009).
9. Ellis, N.M., Li, Y., Hildebrand, W., Fischetti, V.A. & Cunningham, M.W. T cell mimicry and epitope specificity of cross-reactive T cell clones from rheumatic heart disease. *J Immunol* **175**, 5448-56 (2005).
10. Postol, E. *et al.* StreptInCor: a candidate vaccine epitope against *S. pyogenes* infections induces protection in outbred mice. *PLoS One* **8**, e60969 (2013).
11. Guerino, M.T. *et al.* HLA class II transgenic mice develop a safe and long lasting immune response against StreptInCor, an anti-group A streptococcus vaccine candidate. *Vaccine* **29**, 8250-6 (2011).
12. Guilherme, L. *et al.* Towards a vaccine against rheumatic fever. *Clin Dev Immunol* **13**, 125-32 (2006).
13. Hayman, W.A. *et al.* Mapping the minimal murine T cell and B cell epitopes within a peptide vaccine candidate from the conserved region of the M protein of group A streptococcus. *Int Immunol* **9**, 1723-33 (1997).
14. Gupta, V.K. *et al.* Immune response against M protein-conserved region peptides from prevalent group A Streptococcus in a North Indian population. *J Microbiol Immunol Infect* **49**, 352-8 (2016).
15. Brandt, E.R. *et al.* Opsonic human antibodies from an endemic population specific for a conserved epitope on the M protein of group A streptococci. *Immunology* **89**, 331-7 (1996).
16. Relf, W.A. *et al.* Mapping a conserved conformational epitope from the M protein of group A streptococci. *Pept Res* **9**, 12-20 (1996).
17. Brandt, E.R. *et al.* Functional analysis of IgA antibodies specific for a conserved epitope within the M protein of group A streptococci from Australian Aboriginal endemic communities. *Int Immunol* **11**, 569-76 (1999).
18. Pruksakorn, S., Galbraith, A., Houghten, R.A. & Good, M.F. Conserved T and B cell epitopes on the M protein of group A streptococci. Induction of bactericidal antibodies. *J Immunol* **149**, 2729-35 (1992).
19. Batzloff, M.R. *et al.* Protection against group A streptococcus by immunization with J8-diphtheria toxoid: contribution of J8- and diphtheria toxoid-specific antibodies to protection. *J Infect Dis* **187**, 1598-608 (2003).
20. Delunardo, F. *et al.* Streptococcal-vimentin cross-reactive antibodies induce microvascular cardiac endothelial proinflammatory phenotype in rheumatic heart disease. *Clin Exp Immunol* **173**, 419-29 (2013).
21. Galvin, J.E. *et al.* Induction of myocarditis and valvulitis in lewis rats by different epitopes of cardiac myosin and its implications in rheumatic carditis. *Am J Pathol* **160**, 297-306 (2002).
22. Kirvan, C.A., Galvin, J.E., Hilt, S., Kosanke, S. & Cunningham, M.W. Identification of streptococcal m-protein cardiopathogenic epitopes in experimental autoimmune valvulitis. *J Cardiovasc Transl Res* **7**, 172-81 (2014).