

# Fine Mapping Seronegative and Seropositive Rheumatoid Arthritis to Shared and Distinct HLA Alleles by Adjusting for the Effects of Heterogeneity

Buhm Han,<sup>1,2,3</sup> Dorothée Diogo,<sup>1,2,3,4</sup> Steve Eyre,<sup>5,6</sup> Henrik Kallberg,<sup>7</sup> Alexandra Zhernakova,<sup>8,9</sup> John Bowes,<sup>5,6</sup> Leonid Padyukov,<sup>7</sup> Yukinori Okada,<sup>1,2,3,4</sup> Miguel A. González-Gay,<sup>10</sup> Solbritt Rantapää-Dahlqvist,<sup>11</sup> Javier Martin,<sup>12</sup> Tom W.J. Huizinga,<sup>8</sup> Robert M. Plenge,<sup>13</sup> Jane Worthington,<sup>5,6</sup> Peter K. Gregersen,<sup>14</sup> Lars Klareskog,<sup>7</sup> Paul I.W. de Bakker,<sup>1,2,15</sup> and Soumya Raychaudhuri<sup>1,2,3,4,5,\*</sup>

Despite progress in defining human leukocyte antigen (HLA) alleles for anti-citrullinated-protein-autoantibody-positive (ACPA<sup>+</sup>) rheumatoid arthritis (RA), identifying HLA alleles for ACPA-negative (ACPA<sup>-</sup>) RA has been challenging because of clinical heterogeneity within clinical cohorts. We imputed 8,961 classical HLA alleles, amino acids, and SNPs from Immunochip data in a discovery set of 2,406 ACPA<sup>-</sup> RA case and 13,930 control individuals. We developed a statistical approach to identify and adjust for clinical heterogeneity within ACPA<sup>-</sup> RA and observed independent associations for serine and leucine at position 11 in HLA-DRβ1 ( $p = 1.4 \times 10^{-13}$ , odds ratio [OR] = 1.30) and for aspartate at position 9 in HLA-B ( $p = 2.7 \times 10^{-12}$ , OR = 1.39) within the peptide binding grooves. These amino acid positions induced associations at *HLA-DRB1\*03* (encoding serine at 11) and *HLA-B\*08* (encoding aspartate at 9). We validated these findings in an independent set of 427 ACPA<sup>-</sup> case subjects, carefully phenotyped with a highly sensitive ACPA assay, and 1,691 control subjects (HLA-DRβ1 Ser11+Leu11:  $p = 5.8 \times 10^{-4}$ , OR = 1.28; HLA-B Asp9:  $p = 2.6 \times 10^{-3}$ , OR = 1.34). Although both amino acid sites drove risk of ACPA<sup>+</sup> and ACPA<sup>-</sup> disease, the effects of individual residues at HLA-DRβ1 position 11 were distinct ( $p < 2.9 \times 10^{-107}$ ). We also identified an association with ACPA<sup>+</sup> RA at HLA-A position 77 ( $p = 2.7 \times 10^{-8}$ , OR = 0.85) in 7,279 ACPA<sup>+</sup> RA case and 15,870 control subjects. These results contribute to mounting evidence that ACPA<sup>+</sup> and ACPA<sup>-</sup> RA are genetically distinct and potentially have separate autoantigens contributing to pathogenesis. We expect that our approach might have broad applications in analyzing clinical conditions with heterogeneity at both major histocompatibility complex (MHC) and non-MHC regions.

## Introduction

Rheumatoid arthritis (RA [MIM 180300]) has two distinct subtypes—anti-citrullinated-protein-autoantibody-negative (ACPA<sup>-</sup> or seronegative) RA and -positive (ACPA<sup>+</sup> or seropositive) RA—with potentially different genetic risk factors, environmental risk factors, and optimal therapeutic strategies.<sup>1,2</sup> Despite constituting about one-third (~30%) of RA cases,<sup>3</sup> ACPA<sup>-</sup> RA has been relatively understudied in comparison to ACPA<sup>+</sup> RA.<sup>4–7</sup> We and others have demonstrated that the widely established method for identifying ACPA<sup>-</sup> RA subjects on the basis of anticyclic citrullinated peptide (anti-CCP) antibody testing is imperfect in that the absence of antibody is not sufficiently specific to ACPA<sup>-</sup> RA, whereas its presence is specific to ACPA<sup>+</sup> RA.<sup>8–10</sup>

The lack of a specific test for ACPA<sup>-</sup> RA can result in heterogeneity in clinical cohorts, which can confound genetic studies for ACPA<sup>-</sup> disease. For example, ACPA<sup>-</sup> RA subjects might include ACPA<sup>+</sup> RA subjects whose ACPAs have not been detected by conventional anti-CCP testing<sup>8–11</sup> or subjects who have other autoantibody-negative inflammatory arthritic conditions, such as ankylosing spondylitis (AS)<sup>12</sup> or other *HLA-B\*27*-associated conditions. So, although investigators have reported associations between classical HLA alleles and ACPA<sup>-</sup> RA,<sup>13,14</sup> it remains unclear whether these associations are distinct from those alleles driving ACPA<sup>+</sup> disease risk, recently defined by our group.<sup>6</sup> Additionally, the specific amino acid sites and residues driving ACPA<sup>-</sup> RA risk have yet to be defined.

To define HLA alleles driving ACPA<sup>-</sup> RA risk, we first obtained dense SNP genotype data within the major

<sup>1</sup>Division of Genetics, Brigham & Women's Hospital and Harvard Medical School, Boston, MA 02115, USA; <sup>2</sup>Program in Medical and Population Genetics, Broad Institute of Harvard and MIT, Cambridge, MA 02142, USA; <sup>3</sup>Partners HealthCare Center for Personalized Genetic Medicine, Boston, MA 02115, USA; <sup>4</sup>Division of Rheumatology, Immunology, and Allergy, Brigham and Women's Hospital and Harvard Medical School, Boston, MA, 02115, USA; <sup>5</sup>Arthritis Research UK Epidemiology Unit, Musculoskeletal Research Group, University of Manchester, Manchester Academic Health Sciences Centre, Manchester M13 9PT, UK; <sup>6</sup>NIHR Manchester Musculoskeletal Biomedical Research Unit, Central Manchester University Hospitals NHS Foundation Trust, Manchester Academic Health Sciences Centre, Manchester M13 9PT, UK; <sup>7</sup>Rheumatology Unit, Department of Medicine, Karolinska Institutet and Karolinska University Hospital Solna, 171 76 Stockholm, Sweden; <sup>8</sup>Department of Rheumatology, Leiden University Medical Centre, 2300 RC Leiden, the Netherlands; <sup>9</sup>Department of Genetics, University Medical Center Groningen and University of Groningen, 9700 RB Groningen, the Netherlands; <sup>10</sup>Rheumatology Division, Hospital Universitario Marqués de Valdecilla, Instituto de Formación e Investigación Marqués de Valdecilla, 39008 Santander, Spain; <sup>11</sup>Department of Public Health and Clinical Medicine and Department of Rheumatology, Umeå University, 901 85 Umeå, Sweden; <sup>12</sup>Instituto de Parasitología y Biomedicina Lopez-Neyra, Consejo Superior de Investigaciones Científicas, 18100 Armilla, Granada, Spain; <sup>13</sup>Merck Research Laboratories, Merck & Co. Inc., Boston, MA 02115, USA; <sup>14</sup>The Feinstein Institute for Medical Research, North Shore–Long Island Jewish Health System, Manhasset, NY 11030, USA; <sup>15</sup>Departments of Epidemiology and Medical Genetics, University Medical Center Utrecht, 3584 CG Utrecht, the Netherlands

\*Correspondence: [soumya@broadinstitute.org](mailto:soumya@broadinstitute.org)

<http://dx.doi.org/10.1016/j.ajhg.2014.02.013>. ©2014 by The American Society of Human Genetics. All rights reserved.

histocompatibility complex (MHC) region by applying the ImmunoChip custom array<sup>3</sup> to ACPA<sup>-</sup> case and control groups. We then used these data to impute HLA alleles, amino acids, and SNPs with a highly accurate imputation approach.<sup>15</sup> Recognizing that possible clinical heterogeneity within genotyped cohorts might confound associations within the MHC, we developed a statistical approach to correct for the effects of heterogeneity within cohorts; it uses genetic risk scores (GRSs) built from known risk loci for potential confounding diseases as covariates.

We observed that two amino acid positions, HLA-DRβ1 position 11 (in which serine and leucine conferred risk) and HLA-B position 9 (in which aspartate conferred risk), were driving ACPA<sup>-</sup> RA. These two positions are already known to drive ACPA<sup>+</sup> RA as well;<sup>6</sup> however, the specific amino acid residues conferring risk were completely distinct between the two disease subtypes. We also separately tested for associations with ACPA<sup>+</sup> disease. In addition to confirming known associations at positions 11, 71, and 74 in HLA-DRβ1, position 9 in HLA-B, and position 9 in HLA-DPβ1, we identified an additional association at amino acid position 77 within the binding groove of HLA-A. These results contribute to mounting evidence that ACPA<sup>+</sup> and ACPA<sup>-</sup> RA are distinct diseases with certain unique genetic factors.

## Material and Methods

### Samples

#### Case-Control Sample Collections

We used data from six case-control collections (UK, US, Dutch, Spanish, Swedish Umeå, and Swedish Epidemiological Investigation of Rheumatoid Arthritis [EIRA], Table S1, available online).<sup>3</sup> All individuals provided informed consent and were recruited through protocols approved by institutional review boards. Each collection consisted of individuals who were self-described as white and of European descent, and all cases either met the 1987 American College of Rheumatology diagnostic criteria or were diagnosed by board-certified rheumatologists. We previously genotyped all samples with the ImmunoChip custom array, which densely covered the MHC region (7,563 SNPs), in accordance with Illumina protocols.

#### Classifying ACPA<sup>-</sup> RA in Discovery Samples

From these samples, we defined a total of 2,406 ACPA<sup>-</sup> RA case and 13,930 control subjects for discovery from five collections (excluding the Swedish EIRA). To do this, we followed standard clinical practice to identify ACPA<sup>-</sup> RA subjects as those who were not reactive to anti-CCP antibody by using reference cutoff levels defined at local clinical labs. In the UK cohort, we used the commercially available Diastat<sup>TM</sup> ACPA Kit (Axis-Shield Diagnostics Limited). In the US samples, we used a second-generation commercial anti-CCP enzyme immunoassay (Inova Diagnostics).<sup>16</sup> For Spanish samples, we used the Immunoscan ELISA test (Euro Diagnostica). For the Swedish Umeå and Dutch collections, we used the Immunoscan-RA Mark2 ELISA test (Euro Diagnostica).<sup>17</sup> These assays are the standard commercially available assays that are currently being widely used in clinical practice.

#### Clinically Homogeneous ACPA<sup>-</sup> Samples for Replication

To replicate ACPA<sup>-</sup> results, we sought to define an independent replication data set that was as clinically homogeneous as possible. To this end, we used genotype data on 987 case and 1,940 control subjects who were from the Swedish EIRA cohort and who were identified as anti-CCP antibody negative with the Immunoscan-RA Mark2 ELISA test (Euro-Diagnostica). In addition, to stringently ensure clinical homogeneity, we applied a highly sensitive ACPA typing method developed at the Karolinska Institutet<sup>8</sup> to test sera for reactivity to four specific citrullinated peptides ( $\alpha$ -enolase, vimentin, fibrinogen, collagen type II). We considered samples ACPA<sup>-</sup> only if they were negative for all four of these tests. After applying this assay, we removed 106 case individuals who were reactive to the sensitive assay, as well as 381 case individuals to whom we did not apply the assay. We also excluded 73 case and 249 control subjects who were positive for HLA-B\*27. Because HLA-B\*27 is highly sensitive for AS (>90%), excluding HLA-B\*27-positive individuals effectively removed the effect of possible confounding from AS or related spondyloarthropathies. The resulting replication collection consisted of 427 case and 1,691 control subjects.

#### Sample Collections for ACPA<sup>+</sup> RA

For ACPA<sup>+</sup> RA, we used 7,279 anti-CCP-positive individuals from all six cohorts (UK, US, Swedish Umeå, Dutch, Spanish, and Swedish EIRA; Table S1). We used all 15,870 control subjects for ACPA<sup>+</sup> RA analyses.

## Statistical Analyses

### HLA Imputation

We imputed case and control groups together for 8,961 binary markers representing classical HLA alleles, amino acids, and SNPs by using SNP2HLA,<sup>15</sup> which utilizes the Beagle imputation method.<sup>18</sup> The binary markers included every possible grouping of amino acid residues given a multiallelic amino acid position. We used reference data collected by the Type 1 Diabetes Genetics Consortium;<sup>19</sup> these data consisted of genotypes for 5,863 SNPs tagging the MHC and classical alleles for HLA-A, HLA-B, HLA-C, HLA-DRB1, HLA-DQA1, HLA-DQB1, HLA-DPA1, and HLA-DPβ1 at four-digit resolution in 5,225 individuals of European descent.<sup>19</sup>

### Quantifying Imputation Accuracy

To assess accuracy, we took advantage of typed HLA-A, HLA-B, HLA-C, HLA-DQB1, and HLA-DRB1 alleles for 918 individuals in the UK cohort. We calculated imputation accuracy as the proportion of correctly imputed classical alleles:

$$\frac{\sum_i \max(\delta(g_{i,1} = x_{i,1}) + \delta(g_{i,2} = x_{i,2}), \delta(g_{i,1} = x_{i,2}) + \delta(g_{i,2} = x_{i,1}))}{2n}$$

where  $g_{i,1}$  and  $g_{i,2}$  are genotyped alleles of individual  $i$  and  $x_{i,1}$  and  $x_{i,2}$  are imputed alleles. For each gene, we used individuals successfully typed for four-digit alleles. The  $\delta$  function is 1 if the genotyped allele is the imputed allele and 0 otherwise. The term  $n$  is the number of samples.

### Statistical Framework for Association Testing

We tested associations at all 8,961 binary markers by using probabilistic genotypic dosages that take uncertainty in imputation into account. We used logistic regression under the assumption that each marker conferred a fixed log additive effect across each case-control collection. To account for population stratification, we included ten principal components (PCs) as covariates for each collection. We calculated PCs by using EIGENSOFT v.4.2<sup>20</sup> with HapMap Phase 2 samples as reference populations on a

subset of SNPs (minor allele frequency > 0.05) filtered for minimizing intermarker linkage disequilibrium (LD).<sup>3</sup> This resulted in the following logistic regression model:

$$\log(\text{odds}_i) = \theta + \beta_a g_{a,i} + \sum_{j \in \text{collections}} \delta_{i,j} \left( \gamma_j + \sum_{k=1 \dots 10} \pi_{j,k} p_{i,k} \right), \quad (\text{Equation 1})$$

where  $a$  indicates the marker being tested,  $g_{a,i}$  is the dosage of  $a$  in individual  $i$ , and  $\beta_a$  is the additive effect of  $a$ . In the collection-specific term,  $\delta_{i,j}$  is an indicator variable that is 1 only if individual  $i$  is in collection  $j$ . The  $\gamma_j$  parameter is the collection-specific effect due to the differences in case-control proportions; it is set to 0 for one arbitrarily selected reference collection. The  $\pi_{j,k}$  parameter is the effect of the  $k^{\text{th}}$  PC, and  $p_{i,k}$  is the  $k^{\text{th}}$  PC value for individual  $i$ .

#### Adjusting for Clinical Heterogeneity in ACPA<sup>-</sup> Discovery

In the discovery analysis for ACPA<sup>-</sup> disease, we adjusted for possible clinical heterogeneity within the collections. Our approach was to extend Equation 1 to include GRSs of potentially confounding diseases as covariates:

$$\log(\text{odds}_i) = \theta + \beta_a g_{a,i} + \sum_{j \in \text{Collections}} \delta_{i,j} \left( \gamma_j + \sum_{k=1 \dots 10} \pi_{j,k} p_{i,k} + \sum_{h=1 \dots H} \alpha_{j,h} s_{i,h} \right), \quad (\text{Equation 2})$$

where  $h$  indicates a confounding disease we want to adjust for and  $H$  is the total number of confounding diseases.  $s_{i,h}$  is the GRS of individual  $i$  for disease  $h$  and is defined as the sum of risk-allele dosages weighted by effect sizes:

$$s_{i,h} = \sum_l \beta_{l,h} g_{l,i}, \quad (\text{Equation 3})$$

where  $l$  iterates over known risk alleles for  $h$ ,  $\beta_{l,h}$  is the effect size of  $l$  for  $h$ , and  $g_{l,i}$  is the dosage of  $l$  in individual  $i$ .  $\alpha_{j,h}$  is the effect of  $s_{i,h}$ , which approximates the sample proportion of confounding disease in the collection. For a detailed description of the method, see Appendix A.

For our analysis, we adjusted for both ACPA<sup>+</sup> RA and AS. For the ACPA<sup>+</sup> RA GRS,  $l$  iterated over 47 independent SNPs associated with ACPA<sup>+</sup> RA (Table S2),<sup>3</sup> all four-digit *HLA-DRB1* alleles, *HLA-B Asp9*, *HLA-DPβ1 Phe9*, and *HLA-A Asn77*. We estimated  $\beta_l$  from our ACPA<sup>+</sup> RA case-control data set presented in this paper. To estimate  $\beta_l$  for all four-digit *HLA-DRB1* alleles in a multivariate model, we included in the logistic regression all four-digit alleles with allele frequency > 0.1%, except for the reference allele we chose (*HLA-DRB1\*15:01*). To avoid reusing the same controls both to estimate  $\beta_l$  and to map ACPA<sup>-</sup> RA, which could result in bias as a result of overfitting, we estimated  $\beta_l$  for each collection by using the other five collections. Similarly, for the AS GRS,  $l$  iterated over *HLA-B\*27* and 19 AS-associated SNPs that passed our quality control (QC) (Table S2).<sup>12</sup> We used reported effect sizes  $\beta_l$  in Cortes et al.<sup>12</sup>

#### Two-Step Approach for Adjusting for Heterogeneity

Using GRSs as covariates in regression might be overly conservative and could remove true associations if the causal loci are shared between the disease of interest and the confounding disease. To account for the shared genetic structure between the two RA subtypes, we employed an alternative two-step approach: (1) we estimated the confounding proportions  $\alpha_{j,h}$  in Equation 2 by using GRSs based on nonshared loci first, which gave us an unbiased estimate of  $\alpha_{j,h}$ , and then (2) we used this  $\alpha_{j,h}$  as a fixed value in the regression framework presented above. Because we did not definitively know which loci were shared, we used a heuristic to

choose nonshared loci by using 38 non-MHC SNPs not associated with ACPA<sup>-</sup> RA at a nominal significance threshold ( $p > 0.01$ )<sup>3</sup> (Table S2).

#### Genomic-Control Inflation Factor

We assessed the genomic-control inflation factor,  $\lambda_{\text{GC}}$ , by testing associations at “reading-writing-ability SNPs” included on the ImmunoChip platform. Out of 1,469 SNPs, we used 1,250 that passed QC in all six collections. We obtained chi-square statistics at these SNPs by using logistic regression as described above to assess  $\lambda_{\text{GC}}$ .

#### Forward Conditional Search

Once we identified an associated marker, we forward searched further associations by including the identified marker as a covariate in the logistic regression.

#### Exhaustive Search

To find the best pair of associations in *HLA-DRB1* and *HLA-B* for ACPA<sup>-</sup> disease, we examined every possible combination of 495 binary markers within *HLA-DRB1* and 774 binary markers within *HLA-B* (383,130 tests). We extend the single-marker model in Equation 2 to the following two-marker model:

$$\log(\text{odds}_i) = \theta + \beta_a g_{a,i} + \beta_b g_{b,i} + \sum_{j \in \text{collections}} \delta_{i,j} \left( \gamma_j + \sum_{k=1 \dots 10} \pi_{j,k} p_{i,k} + \sum_{h=1 \dots H} \alpha_{j,h} s_{i,h} \right), \quad (\text{Equation 4})$$

where  $a$  and  $b$  are the pair of binary markers being tested. We calculated the log-likelihood difference ( $\Delta\text{LL}$ ) in model fit due to this pair and assessed significance by comparing the deviance ( $-2 \times \Delta\text{LL}$ ) to a chi-square distribution with 2 degrees of freedom.

#### Joint Analysis of Discovery and Replication Data

In order to jointly analyze five discovery collections and a replication cohort for ACPA<sup>-</sup> disease, we combined them into one logistic regression framework, including GRSs as covariates for five discovery cohorts to adjust for heterogeneity.

#### Forward Search outside of HLA-DRB1 for ACPA<sup>+</sup> RA

Because *HLA-DRB1* has a very strong effect in ACPA<sup>+</sup> disease, to examine the associations beyond *HLA-DRB1*, we conditioned on the *HLA-DRB1* effects by including binary variables as covariates corresponding to all four-digit *HLA-DRB1* alleles, excluding one allele as a reference (*HLA-DRB1\*15:01*). If we forward searched by conditioning on an amino acid position with  $m$  residues, such as position 9 of *HLA-B*, we included binary variables corresponding to the  $m - 1$  residues, excluding the most frequent one.

#### Testing for Discordant Effect Sizes

Given a multiallelic amino acid position with  $m$  residues, we wanted to test whether the effect sizes of  $m$  residues were concordant between two different conditions (e.g., ACPA<sup>-</sup> versus ACPA<sup>+</sup>). To this end, we calculated multivariate odds ratios (ORs) of residues by including in the logistic regression  $m - 1$  binary markers corresponding to  $m - 1$  residues, excluding one residue as the reference. Let  $a_1, \dots, a_{m-1}$  and  $b_1, \dots, b_{m-1}$  be the multivariate log ORs in two different conditions. Let  $v_1, \dots, v_{m-1}$  and  $u_1, \dots, u_{m-1}$  be their variances. To test discordance of effect sizes between two conditions, we used the statistic

$$\sum_{i=1 \dots m} \frac{(a_i - b_i)^2}{v_i + u_i}, \quad (\text{Equation 5})$$

which is chi-square distributed with  $m - 1$  degrees of freedom under the null.

To test the accuracy of our approach to adjust for clinical heterogeneity in fine mapping, we simulated an ACPA<sup>-</sup> RA case-control study confounded by ACPA<sup>+</sup> RA. We simulated a large study (50,000 case and 50,000 control subjects) to assess the asymptotic results. We first simulated control subjects by sampling with replacement from the UK control subjects. Then we assumed that specific amino acid positions were conferring risk to ACPA<sup>-</sup> RA with predefined ORs, and we sampled ACPA<sup>-</sup> RA subjects from the UK control subjects on the basis of the ORs. Finally, we replaced 26.3% of the case group with individuals randomly sampled from the UK ACPA<sup>+</sup> RA case group. We performed an association test with and without adjusting for heterogeneity to examine whether we could fine map the risk-conferring amino acid positions correctly. To adjust for heterogeneity, we used GRSs built from the effect sizes estimated from the other five cohorts, excluding the UK cohort.

## Results

### ACPA<sup>-</sup> RA Discovery Collection and HLA Imputation

To define HLA alleles driving ACPA<sup>-</sup> RA risk, we analyzed a discovery data set of 2,406 ACPA<sup>-</sup> RA case and 13,930 control subjects (from the UK, the US, Spain, Sweden, and the Netherlands, see Table S1) genotyped on the Immunochip custom array with 7,563 SNPs across the MHC region.<sup>3</sup> This platform represents greater SNP density than most standard genome-wide-association-study arrays and offers the potential for higher HLA imputation accuracy. Indeed, applying SNP2HLA,<sup>15</sup> we observed an overall imputation accuracy of 96.9% for four-digit HLA alleles in a subset of UK control subjects separately typed for HLA alleles (Table S3). We classified RA samples as ACPA<sup>-</sup> on the basis of anti-CCP antibody amounts according to standard clinical practice (see Material and Methods). After adjusting for ten PCs, we observed little evidence of population stratification ( $\lambda_{GC} = 0.98$ , see Material and Methods).

### Correcting for Clinical Heterogeneity in ACPA<sup>-</sup> RA Collections

We considered that other syndromes clinically indistinguishable from ACPA<sup>-</sup> RA might be embedded within ACPA<sup>-</sup> RA and thus confound associations. Indeed, in an analysis unadjusted for clinical heterogeneity, we observed that as we defined ACPA<sup>-</sup> samples by increasing the level of stringency of the anti-CCP cutoff, the frequency of HLA-DR $\beta$ 1 Val11 (the strongest risk factor for ACPA<sup>+</sup> disease) decreased in our ACPA<sup>-</sup> cohort ( $p = 6.9 \times 10^{-5}$ ), suggesting confounding from ACPA<sup>+</sup> RA (Figure S1). We also noticed significant association at *HLA-B\*27* ( $p = 2.8 \times 10^{-9}$ ), a well-known risk factor for AS,<sup>12,21,22</sup> but not at *HLA-C\*06:02* ( $p > 0.001$ ), a risk factor for psoriatic arthritis.<sup>23–25</sup> However, as in most clinical settings, the phenotypic information that would be essential for identifying and excluding the specific individuals with conditions other than ACPA<sup>-</sup> RA was not available.

To correct for the effects of heterogeneous samples within our ACPA<sup>-</sup> cohort, we applied a statistical approach to adjust

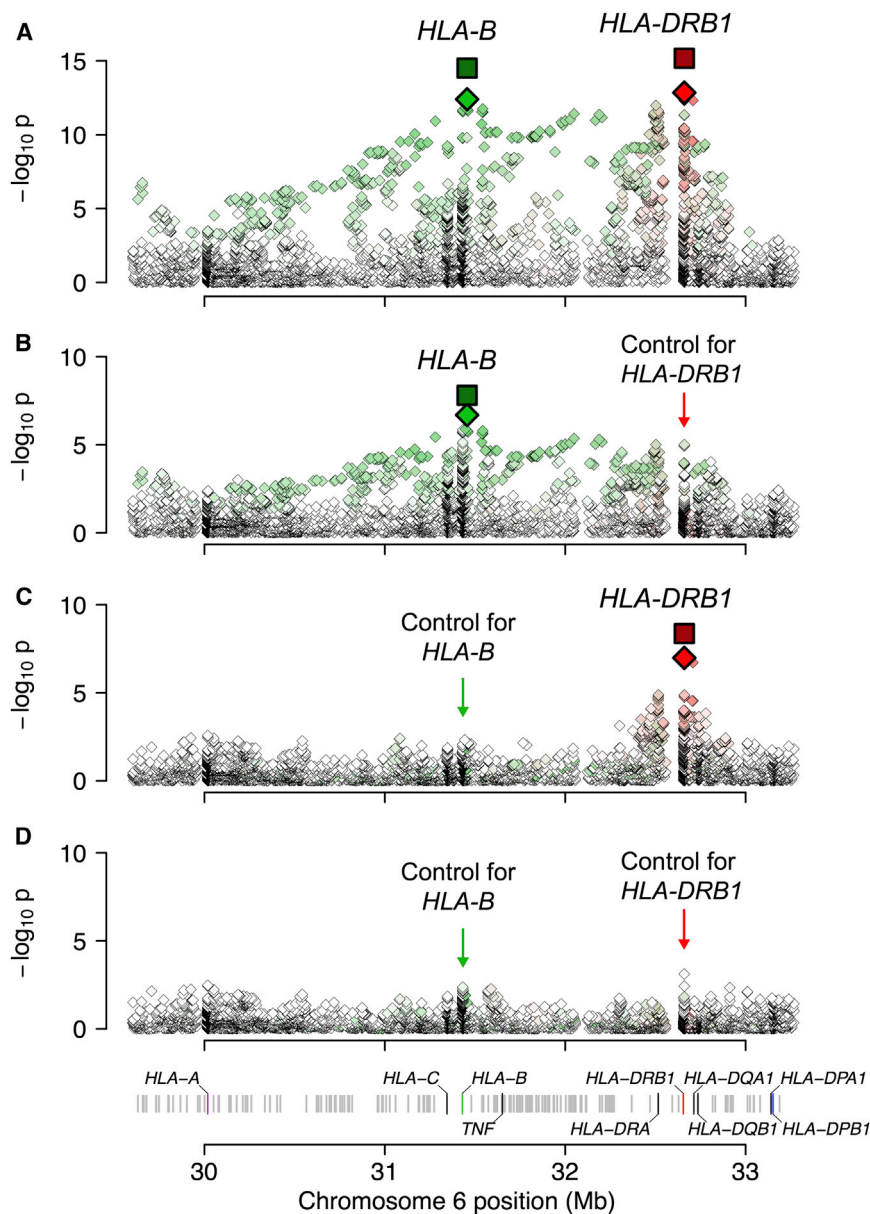
for confounding diseases (ACPA<sup>+</sup> RA and AS, Material and Methods). We constructed GRSs representing the log OR for an individual for the confounding disease on the basis of the known-risk-allele dosages weighted by effect sizes.<sup>26–28</sup> Then, adjusting association statistics in a logistic regression model for GRSs could successfully control for the effects of confounding diseases (see Appendix A).

### ACPA<sup>-</sup> RA Is Associated with Ser11 and Leu11 in HLA-DR $\beta$ 1 and Asp9 in HLA-B

After correcting for clinical heterogeneity as described above, we tested for allelic associations in ACPA<sup>-</sup> RA. Taking into account multiple hypothesis testing, we considered  $p < 5.6 \times 10^{-6}$  (0.05/8,961 binary MHC-marker association tests) to be significant. After testing all amino acids and classical and SNP alleles, we observed that the strongest association was at amino acid residues at position 11 in HLA-DR $\beta$ 1 (presence of Ser or Leu, OR = 1.30,  $p = 1.4 \times 10^{-13}$ ), encoded by *HLA-DRB1* (see Figure 1A, Table 1, and Figure S2). This allele exceeded the significance of all other SNPs and classical alleles that we tested. The variation of amino acid residues at this position was attributable to a triallelic SNP (rs9269955, G/C/A) and a quadallelic SNP (rs17878703) at the first and second base positions of the codon, respectively. The association at position 11 was statistically indistinguishable ( $p > 0.09$ ) from the association at position 13 (presence of Ser, Gly, or Phe, OR = 1.29,  $p = 4.7 \times 10^{-13}$ ). The most strongly associated classical allele was *HLA-DRB1\*03* ( $p = 6.7 \times 10^{-10}$ ).<sup>13,14</sup> After conditioning on *HLA-DRB1\*03*, we observed that Ser11+Leu11 remained highly significant ( $p = 2.4 \times 10^{-8}$ ), suggesting that *HLA-DRB1\*03* does not fully explain *HLA-DRB1* associations. We also observed a separate, strong association 23 kb away from *HLA-B* at SNP rs9266669 (OR = 1.38,  $p = 4.0 \times 10^{-13}$ ; Figure 1A). This SNP was statistically indistinguishable ( $p > 0.01$ ) from the presence of Asp9 in HLA-B (OR = 1.39,  $p = 2.7 \times 10^{-12}$ ); these two alleles were in tight LD ( $r^2 = 0.8$ ). HLA-B Asp9 was almost perfectly correlated with *HLA-B\*08* in our data set ( $r^2 = 0.997$ ). The *HLA-B\*08* classical allele, Asp9, and SNP rs9266669 thus could not be distinguished on the basis of genetics alone. Both of these amino acid sites mapped to the binding grooves of their respective HLA receptors (Figure 2).

The *HLA-DRB1* and *HLA-B* associations were independent of each other and explained most of the MHC association with ACPA<sup>-</sup> RA. After conditioning on Ser11+Leu11 effects in HLA-DR $\beta$ 1, we observed that rs9266669 in *HLA-B* (or Asp9 in HLA-B) remained the most significant association ( $p = 2.0 \times 10^{-7}$ , OR = 1.27; Figure 1B). Similarly, we observed that after conditioning on Asp9 in HLA-B, Ser11+Leu11 in HLA-DR $\beta$ 1 remained the most significant association ( $p = 1.0 \times 10^{-7}$ , OR = 1.22; Figure 1C). When we conditioned on both Ser11+Leu11 in HLA-DR $\beta$ 1 and Asp9 in HLA-B, no further significant association was found ( $p > 0.0007$ ; Figure 1D).

Because the so-called 8.1 ancestral haplotype<sup>29</sup> harbors both HLA-DR $\beta$ 1 Ser11 and HLA-B Asp9, we considered



**Figure 1. Association Results within the MHC to ACPA<sup>-</sup> RA**

(A) We observed the most significant association at position 11 of HLA-DRβ1 (encoded by *HLA-DRB1*), where Ser and Leu conferred risk (red diamond). We also observed an independent association at SNP rs9266669, which was statistically indistinguishable from HLA-B Asp9 (green diamond). The dark-red and dark-green squares denote the statistical significance of the two positions in a joint analysis including both discovery and replication data.

(B) Conditioning on HLA-DRβ1 Ser11+Leu11, we found that the association at rs9266669 remained the most significant.

(C) Conditioning on HLA-B Asp9, we found that the association at HLA-DRβ1 Ser11+Leu11 remained the most significant.

(D) Conditioning on both HLA-DRβ1 Ser11+Leu11 and HLA-B Asp9, we did not observe any more statistically significant association within MHC ( $p > 0.0007$ ).

shared loci between two subtypes of RA. To address this concern, we developed a two-step alternative approach that estimates the confounding proportion (proportion of misdiagnosed ACPA<sup>+</sup> RA samples within ACPA<sup>-</sup> RA cohorts) by using a GRS calculated on the basis of an approximated set of nonshared loci (i.e., known loci associated with ACPA<sup>+</sup> RA but with  $p > 0.01$  association in ACPA<sup>-</sup> RA) and then regresses out only this amount from the model (see [Material and Methods](#)). The confounding proportion estimates by this approach were comparable to the estimates by the previous approach with the full GRS

the possibility that these associations were driven by that haplotype alone and not the individual amino acid sites. Given that our imputation provided phased haplotypes spanning the whole MHC region, we inferred the ancestral haplotype dosage for each individual. Then, using a trivariate logistic regression model including dosages for the 8.1 ancestral haplotype, HLA-DRβ1 Ser11+Leu11, and HLA-B Asp9, we observed that association at the ancestral haplotype was not significant ( $p = 0.21$ ). In contrast, the other two HLA amino acid variables retained statistical significance even after adjustment for the effect of the 8.1 ancestral haplotype ( $p = 1.6 \times 10^{-7}$  at HLA-DRβ1 Ser11+Leu11 and  $p = 3.4 \times 10^{-3}$  at HLA-B Asp9). These results suggest that the association was driven primarily by the amino acid sites and not by the effect of the 8.1 haplotype alone.

We further considered that our approach to correcting for heterogeneity might be conservative and might remove

(mean proportion across cohorts was 26.3% with the full GRS and 28.3% with the nonshared-loci GRS; see [Figure S3](#)). Consistent with the previous approach, this two-step approach produced the most significant associations at rs9266669 ( $p = 1.8 \times 10^{-13}$ , OR = 1.38 at HLA-B Asp9) and HLA-DRβ1 Ser11+Leu11 ( $p = 2.3 \times 10^{-13}$ , OR = 1.27). Again, these two associations were independent ( $p = 5.4 \times 10^{-8}$ ).

### Replicating HLA Associations in a Clinically Homogeneous ACPA<sup>-</sup> Collection

We wanted to validate these findings in an independent cohort without significant clinical heterogeneity. To this end, we assessed association in an independent data set of 427 phenotypically homogeneous ACPA<sup>-</sup> individuals and 1,691 control subjects (Swedish EIRA). According to a state-of-the-art commercially unavailable assay,<sup>8</sup> these

**Table 1. Effect Estimates for Amino Acids Associated with Risk of ACPA<sup>-</sup> and ACPA<sup>+</sup> RA**

RA Subtypes	HLA Protein	Amino Acid Position	Amino Acid Residue	OR after Adjustment for Known Associated Positions (95% CI)			Frequency in Control Group	Frequency in Case Group	Classical Alleles
				Discovery	Replication	Joint			
ACPA <sup>-</sup>	HLA-DRβ1	11	Ser+Leu	1.22 (1.14–1.32)	1.22 (1.04–1.43)	1.22 (1.14–1.31)	0.514	0.548	<i>HLA-*01, HLA-*03, HLA-*08, HLA-*11, HLA-*12, HLA-*13, HLA-*14</i>
	HLA-B	9	Asp	1.27 (1.15–1.40)	1.23 (0.99–1.52)	1.26 (1.15–1.38)	0.131	0.161	<i>HLA-*08</i>
ACPA <sup>+</sup>	HLA-A	77	Asn	0.85 (0.81–0.90)			0.343	0.279	<i>HLA-*01, HLA-*23, HLA-*24, HLA-*26, HLA-*29, HLA-*30, HLA-*36, HLA-*80</i>

For each amino acid identified in this study, we show the OR and 95% confidence interval (95% CI), unadjusted frequencies in the case and control groups, and corresponding classical HLA alleles. All ORs were conditioned on known associated positions; for ACPA<sup>-</sup> RA, we estimated ORs of HLA-DRβ1 Ser11+Leu11 and HLA-B Asp9 by conditioning on each other. For ACPA<sup>+</sup> RA, we estimated the OR of HLA-A Asn77 by conditioning on all alleles at *HLA-DRB1*, amino acids at HLA-B position 9, and amino acids at HLA-DPβ1 position 9. See Table S7 for the complete table, including previously identified positions.

ACPA<sup>-</sup> individuals were negative for not only anti-CCP antibody but also antibodies for four specific citrullinated peptide antigens. We also excluded *HLA-B\*27*-positive individuals (>90% sensitive for AS) from case and control groups. We tested for association without any adjustment for heterogeneity. We confirmed associations both at HLA-DRβ1 Ser11+Leu11 ( $p = 5.8 \times 10^{-4}$ , OR = 1.28) and at HLA-B Asp9 ( $p = 2.6 \times 10^{-3}$ , OR = 1.34) with comparable effect sizes (Table 1). These associations were again independent of each other. Conditioning on HLA-DRβ1 Ser11+Leu11, we observed an independent effect at HLA-B Asp9 ( $p = 0.03$ , OR = 1.23). Conversely, conditioning on HLA-B Asp9, we observed an independent effect at HLA-DRβ1 Ser11+Leu11 ( $p = 0.007$ , OR = 1.22).

In a joint analysis of the discovery and replication cohorts, we observed increased significance at both HLA-DRβ1 and HLA-B positions ( $p = 6.7 \times 10^{-16}$  and OR = 1.30 for HLA-DRβ1 Ser11+Leu11;  $p = 5.3 \times 10^{-14}$  and OR = 1.38 for HLA-B Asp9; Figure 1A and Table S4) and that their effects were independent ( $p < 2 \times 10^{-8}$ ; Figures 1B and 1C and Table S4). Conditioning on both of these effects, we observed no other independent association throughout the MHC ( $p > 0.0002$ ).

#### Exhaustive Search Confirms Associations with Ser11 and Leu11 in HLA-DRβ1 and Asp9 in HLA-B

Because the conditional forward search might miss the best explanations, we exhaustively tested every possible pair of binary markers in *HLA-DRB1* and *HLA-B* in a joint analysis. Out of 383,130 pairs we tested, HLA-DRβ1 Ser11+Leu11 and HLA-B Asp9 in HLA-B (or equivalently *HLA-B\*08* and *HLA-B\*0801*) constituted the most significant pair ( $p = 1.1 \times 10^{-20}$ ; Table S5), confirming that our model provides the most parsimonious explanation of the data.

#### Associations Are Independent of Rheumatoid Factor Status

We examined whether the associations we identified were independent of rheumatoid factor (RF) status. We obtained

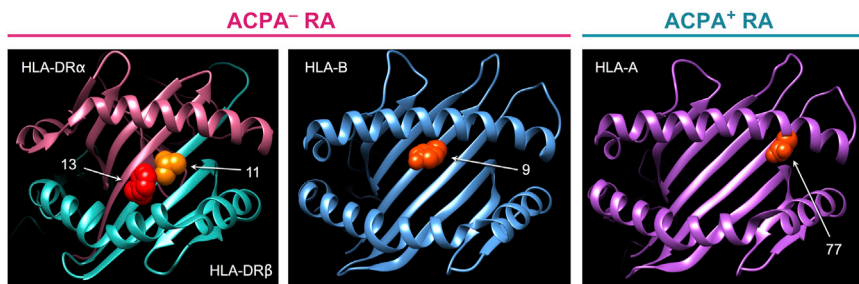
RF data for 1,016 affected individuals in the UK cohort; 470 individuals (46%) were RF<sup>+</sup>, and 546 individuals (54%) were RF<sup>-</sup>. We stratified the samples into two groups on the basis of RF status. The associations were consistent between the two groups in that they showed the same direction of effects at both HLA-DRβ1 Ser11+Leu11 and HLA-B Asp9 (Table S6). We observed that effect sizes tended to be greater in the RF<sup>+</sup> subjects than in the RF<sup>-</sup> subjects at both loci ( $p = 0.02$ ). A thorough investigation of this phenomenon will require larger sample sizes.

#### Asn77 at HLA-A Is Associated with ACPA<sup>+</sup> RA

We also mapped associations within the MHC to ACPA<sup>+</sup> RA in 7,279 ACPA<sup>+</sup> RA subjects and 15,870 control subjects (see Table S1 and Material and Methods). We observed little evidence of stratification after adjusting for ten PCs ( $\lambda_{GC} = 1.07$ ). We confirmed previously published associations in HLA-DRβ1 at amino acid positions 11 ( $p < 10^{-692}$ ), 71 ( $p < 10^{-37}$ ), and 74 ( $p < 10^{-23}$ ) (Table S7). Conditioning on *HLA-DRB1* alleles, we confirmed associations at Asp9 in HLA-B ( $p < 10^{-36}$ , OR = 1.93) and Phe9 in HLA-DPβ1 ( $p < 10^{-19}$ , OR = 1.31)<sup>6</sup> (Figure S4). Conditioning on all of these previously known associated positions (the *HLA-DRB1* alleles, position 9 in HLA-B, and position 9 in HLA-DPβ1), we observed an independent association with ACPA<sup>+</sup> RA with the presence of Asn77 in HLA-A ( $p = 2.7 \times 10^{-8}$ , OR = 0.85; Figure S4D and Table 1). Similar to the other amino acid sites associated with RA,<sup>6</sup> position 77 in HLA-A was also located in the binding groove (Figure 2 and Figure S5). We previously observed that Ser77 in HLA-A confers protection in HIV controllers.<sup>31</sup> After conditioning on this sixth position, we observed no convincing associations ( $p > 4 \times 10^{-6}$ ).

#### Discussion

In this study, we observed that associations with ACPA<sup>-</sup> RA within the MHC were driven by *HLA-DRB1* and *HLA-B*. In



**Figure 2. 3D Models of Amino Acid Positions Identified in This Study**

Key amino acid positions are highlighted as spheres. We used Protein Data Bank entries 3pdo (HLA-DR), 2bvp (HLA-B), and 1x7q (HLA-A) with UCSF Chimera to prepare the figure.<sup>30</sup> See Figure S5 for all known associated positions.

addition, we identified the specific residues and specific amino acid sites that parsimoniously explained these associations. These positions mapped to the peptide binding grooves of these receptors, pointing to an important role for antigen recognition. The success of this study was contingent on our ability to distinguish the effects from other conditions contributing to heterogeneity within the case individuals.

Intriguingly, the positions that drove ACPA<sup>-</sup> risk were the same positions that drove most risk for ACPA<sup>+</sup> RA as well (Table S8). The risk of Asp9 in HLA-B in ACPA<sup>-</sup> RA was shared with ACPA<sup>+</sup> disease but had a more modest effect size (OR = 1.38 in ACPA<sup>-</sup> versus OR = 1.93 in ACPA<sup>+</sup>). This allele, also associated with myasthenia gravis,<sup>32</sup> might affect nonspecific immune reactivity.

In contrast, at position 11 of HLA-DRβ1, different residues drove risk of the two diseases (discordance  $p < 2.9 \times 10^{-107}$ ; Figure 3). For example, Ser11 conferred risk of ACPA<sup>-</sup> disease (OR = 1.31) but was protective against ACPA<sup>+</sup> disease (OR = 0.39). On the other hand, Gly11 and Pro11 showed protective effects for both subsets. We speculate that citrullinated antigens that drive ACPA<sup>+</sup> RA risk might be biochemically distinct from the antigens driving ACPA<sup>-</sup> RA risk, for example, carbamylated antigens.<sup>33</sup> The different set of risk and protective residues for the two disease subsets might be related to differential binding affinity and reactivity to these autoantigens.

In a multicohort study where allele frequencies can differ between cohorts, it is crucial to account for population stratification. For example, the frequency of ancestral 8.1 haplotype differed from 5% to 17% depending on cohorts (Table S9). As described in the Material and Methods, we took two approaches to account for population structure: (1) we stratified the data by country of origin, and (2) we used ten PCs to aggressively adjust for any residual population effects. The effectiveness of this standard approach is reflected in the relatively modest inflation factors for the study ( $\lambda_{1,000} = 1.00$  for ACPA<sup>-</sup> RA and  $\lambda_{1,000} = 1.01$  for ACPA<sup>+</sup> RA).

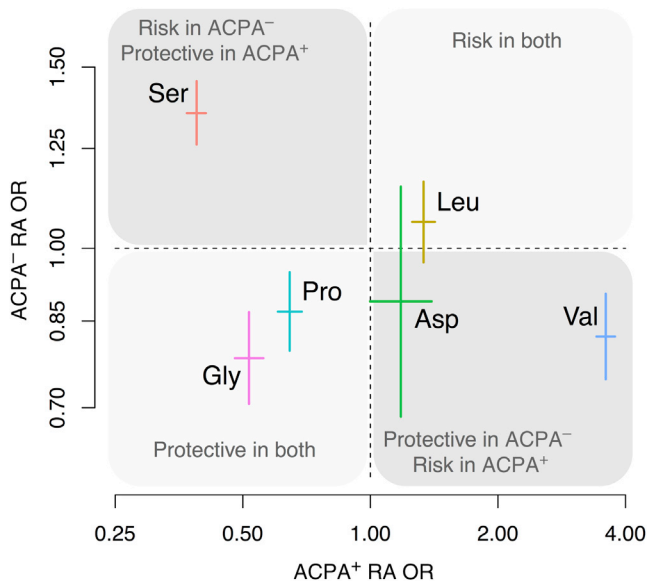
In this study, we addressed the issue of heterogeneity within cohorts. Like for population stratification, if the heterogeneity is present and we fail to adequately adjust for it, spurious associations can occur. For example, without adjusting for heterogeneity, the top ACPA<sup>-</sup> RA association appeared to be at Leu67 in HLA-DRβ1 ( $p = 2.9 \times 10^{-28}$ ). Despite its remarkable significance in our het-

erogeneous discovery sample, Leu67 failed to replicate when we examined it in our homogenous replication data set ( $p = 0.26$ ). In contrast, after adjusting for heterogeneity in our discovery data set, we observed the strongest effect at position 11 of HLA-DRβ1 (Table 1); not only did this effect replicate in our homogenous replication data set, but the effect sizes of each amino acid residue at that site were also highly concordant between discovery and replication sets (discordance  $p > 0.4$  after adjustment; Figure S6).

To further demonstrate the potential for accounting for heterogeneity in fine mapping, we performed simulations. We simulated a study under the assumption that HLA-DRβ1 Ser11+Leu11 (OR = 1.30) and HLA-B Asp9 (OR = 1.39) confer risk, which is the model that we found in this study, and included ACPA<sup>+</sup> RA subjects in 26.3% of affected individuals (Material and Methods). Without adjustment for heterogeneity, the top association was deceptively at HLA-DRβ1 Leu67 ( $p < 10^{-331}$ ), which was exactly what we observed in discovery cohorts without adjusting for heterogeneity. Using our statistical approach to adjust for heterogeneity, we were able to map the correct positions we simulated; the top associations were HLA-DRβ1 Ser11+Leu11 ( $p = 1.3 \times 10^{-189}$ ), and conditioned on this, rs2853986 ( $p = 7.2 \times 10^{-59}$ ), which was statistically indistinguishable ( $p > 0.05$ ) from HLA-B Asp9. We also showed that adjusting for heterogeneity not only removed spurious associations but also provided accurate estimation of the proportion of confounding samples under the null model (Figure S7).

We note that we adjusted for possible confounding from AS by correcting for AS GRSs in discovery cohorts and removing HLA-B\*27-positive individuals in the replication cohort. This approach effectively adjusted for putative HLA-B\*27 associations with ACPA<sup>-</sup> RA if there were any. Currently, it is difficult to distinguish true HLA-B\*27 associations from confounding from AS. We expect that we will be able to accurately distinguish these two situations as we identify a greater number of non-MHC AS risk loci in the future.

The concern of clinical heterogeneity extends beyond RA to a wide range of diseases where clinical classification might be uncertain because of imperfect diagnostic tests, for example, (1) subclassification of inflammatory bowel disease (MIM 266600) into Crohn disease or ulcerative colitis or (2) distinguishing early bipolar disease (MIM



**Figure 3. Distinct Effect Sizes of Amino Acid Residues at HLA-DRβ1 Position 11 for ACPA<sup>-</sup> and ACPA<sup>+</sup> RA**

For each residue, we show the univariate OR (OR with respect to the other residues as a reference) and the 95% confidence interval. Effect sizes were distinct between the two disease subsets ( $p < 2.9 \times 10^{-107}$ ).

125480) from major depressive disorder (MIM 608516). We expect that our statistical approach might have application to genetic studies of these conditions as well. The applicability of our approach is contingent on adequate power to detect confounding genetic effects; such power is only possible when sufficient numbers of genetic loci for confounding diseases are known. We also expect that our approach might have utility in better characterizing non-HLA loci of the conditions with clinical heterogeneity.

Our results have important implications for the clinical practice of ACPA<sup>-</sup> RA. Investigators have long speculated that individuals diagnosed with ACPA<sup>-</sup> RA might have other inflammatory arthritic conditions, such as AS, that mimic RA and have atypical clinical presentations. Our analysis supports this; we estimated here that each ACPA<sup>-</sup> RA cohort contained 4%–11% of the affected individuals who most likely had AS and 15%–37% of affected individuals who most likely had ACPA<sup>+</sup> RA (Table S10 and Figure S3). We note the possibility that other conditions that we did not account for, such as Sjögren syndrome (MIM 270150),<sup>34</sup> might have been included within the ACPA<sup>-</sup> RA samples. These subjects were identified through research protocols, and in clinical practice, these diagnostic uncertainties can be even more pronounced. Clinical misclassifications can be particularly concerning in this setting given that optimal pharmacological treatment and long-term prognosis for these different arthritic conditions vary. Our data not only underscore the need for more accurate clinical tests than the conventional anti-CCP antibody testing but also illuminate the potential

role of genetic data in helping categorize individuals with ACPA<sup>-</sup> inflammatory arthritis.

## Appendix A

### Asymptotic Mean of Effect-Size Estimate in the Presence of Confounding

We first consider linear regression for quantitative traits. We assume a single locus, which we will extend to multiple loci later. Suppose that two groups of samples are mixed in a cohort. Let  $x_1$  and  $x_2$  be the genotype vectors of the two groups at the locus and  $y_1$  and  $y_2$  be the phenotype vectors. Let  $\beta_1$  and  $\beta_2$  be the effect sizes, such that the true model is  $y_1 = x_1\beta_1 + \varepsilon_1$  and  $y_2 = x_2\beta_2 + \varepsilon_2$ , where  $\varepsilon_1$  and  $\varepsilon_2$  are error terms. Without loss of generality, assume that  $x_1$ ,  $x_2$ ,  $y_1$ , and  $y_2$  have zero mean. Because of sample mixture, what we observe are  $x = (x_1^T | x_2^T)^T$  and  $y = (y_1^T | y_2^T)^T$ . The standard linear regression formula gives us the least-squares estimate of effect size:

$$\begin{aligned} \hat{\beta} &= (x^T x)^{-1} x^T y \\ &= (x_1^T x_1 + x_2^T x_2)^{-1} (x_1^T | x_2^T) \left( (x_1 \beta_1 + \varepsilon_1)^T | (x_2 \beta_2 + \varepsilon_2)^T \right)^T \\ &= (x_1^T x_1 + x_2^T x_2)^{-1} \left( (x_1^T x_1 \beta_1 + x_1^T \varepsilon_1) + (x_2^T x_2 \beta_2 + x_2^T \varepsilon_2) \right) \\ &= (x_1^T x_1 + x_2^T x_2)^{-1} \left( (x_1^T x_1) \left( \beta_1 + (x_1^T x_1)^{-1} x_1^T \varepsilon_1 \right) \right. \\ &\quad \left. + (x_2^T x_2) \left( \beta_2 + (x_2^T x_2)^{-1} x_2^T \varepsilon_2 \right) \right) \end{aligned}$$

Given that  $E[(x_1^T x_1)^{-1} x_1^T \varepsilon_1] = 0$  and  $E[(x_2^T x_2)^{-1} x_2^T \varepsilon_2] = 0$ ,

$$E[\hat{\beta}] = (x_1^T x_1 + x_2^T x_2)^{-1} (x_1^T x_1 \beta_1 + x_2^T x_2 \beta_2)$$

If we assume that the minor allele frequency of the variant is the same for the two groups and the genotypes follow Hardy-Weinberg equilibrium,  $(x_1^T x_1) / (x_2^T x_2) \approx N_1 / N_2$ , where  $N_1$  and  $N_2$  are the sample sizes of the two groups. Thus, the effect-size estimate asymptotically converges to an average effect size weighted by the sample sizes of two groups.

This result has the following implication. Suppose that  $\beta_1$  is the true effect size of interest and  $\beta_2$  is the effect size for confounding samples. Consider the null model ( $\beta_1 = 0$ ). What we observe will be  $E[\hat{\beta}] = \alpha \beta_2$ , where  $\alpha$  is the confounding proportion. Thus, we will have spurious association ( $E[\hat{\beta}] \neq 0$ ). Suppose that we build GRSs with respect to confounding disease as  $s = x \beta_2$ . If we regress out  $s$  as a covariate, it will remove spurious association. Moreover, the regression coefficient of  $s$  will be an unbiased estimator of  $\alpha$ .

Under the alternative model ( $\beta_1 \neq 0$ ), using risk score as a covariate might be conservative and remove true association. If we know  $\alpha$  a priori, one approach is fixing the coefficient of  $s$  to the constant  $\alpha$ . That is, we subtract  $s\alpha = x\beta_2\alpha$  from  $y$ . This approach will retain true association. The effect-size estimate can still be conservative, given that what we would want to subtract is actually  $x(\beta_2 - \beta_1)\alpha$ , which is unknown.



## Logistic Regression

Similar results extend to logistic regression. For simplicity, we assume the null model (true OR is 1). Suppose that  $\alpha\%$  of the case group is confounded by a disease whose OR is  $\gamma \neq 1$ . Let  $p$  be the control minor allele frequency. Then, the asymptotic mean of the observed log OR  $\hat{\beta}$  will be

$$E[\hat{\beta}] = \pi = \log \frac{(\alpha p_A + (1 - \alpha)p)(1 - p)}{(\alpha(1 - p_A) + (1 - \alpha)(1 - p))p}$$

where  $p_A = \gamma p / ((\gamma - 1)p + 1)$  is the case minor allele frequency of the confounding disease. Thus, we will have spurious association ( $E[\hat{\beta}] \neq 0$ ).

If  $\gamma$  is small, we can establish an approximate relationship,  $\pi \approx \alpha \log(\gamma)$ , which we show by simulations (Figure S8). Thus, using risk score  $s = \log(\gamma)x$  as a covariate, we can not only remove spurious association but also approximate  $\alpha$  from the regression coefficient of  $s$ .

## Generalization to Multiple Loci

We can generalize our approach to multiple loci. Suppose that we know  $m$  independent loci associated with the confounding disease. Let  $\beta_1, \dots, \beta_m$  be their effect sizes. We build GRSs for each individual locus,

$$s_i = x_i \beta_i \quad i \in \{1, \dots, m\},$$

where  $x_i$  is the genotype vector at locus  $i$ . In order to estimate the confounding proportion  $\alpha$ , we look at all loci together by including all  $s_i$  in the regression:

$$y = \alpha s_1 + \alpha s_2 + \dots + \alpha s_m + \varepsilon.$$

Application to logistic regression is also straightforward. Because  $\alpha$  is invariant across loci, this is equivalent to the model using a combined GRS,  $y = \alpha S + \varepsilon$ , where  $S = \sum s_i = \sum x_i \beta_i$ , which results in the approach presented in the [Material and Methods](#). The advantage of a combined GRS over multiple loci is that it can be less conservative under the alternative model. For example, if we test locus  $i$  and include  $s_i$  as a covariate, it will remove true association. However, if we include  $S$  as a covariate, the information from other loci ( $s_1, s_2, \dots, s_{i-1}, s_{i+1}, \dots, s_m$ ) will help in finding correct  $\alpha$  and preventing overly regressing out  $s_i$ . Another possible way to more strictly prevent overly regressing out GRS can be estimating  $\alpha$  with nonoverlapping loci first, as presented in the [Material and Methods](#).

## Supplemental Data

Supplemental Data include eight figures and ten tables and can be found with this article online at <http://www.cell.com/ajhg>.

## Acknowledgments

This work was supported by funds from the National Institutes of Health (K08AR055688, 1R01AR062886-01, 1R01AR063759-01A1, and 5U01GM092691-04), the Arthritis Foundation, and the Doris Duke Foundation and in part through the Be the Cure For Rheumatoid Arthritis grant funded by the Innovative Medicine Initia-

tive program from the European Union. This research used data provided by the Type 1 Diabetes Genetics Consortium (a collaborative clinical study sponsored by the National Institute of Diabetes and Digestive and Kidney Diseases, National Institute of Allergy and Infectious Diseases, National Human Genome Research Institute, National Institute of Child Health and Human Development, and Juvenile Diabetes Research Foundation International). A.Z. was supported by a grant from the Dutch Reuma-fonds (11-1-101) and the Rosalind Franklin Fellowship from the University of Groningen (the Netherlands). These data also included data generously provided by the Rheumatoid Arthritis International Consortium. P.I.W.d.B. is the recipient of a Vidi award from the Netherlands Organization for Scientific Research (project 016.126.354). This work was partially supported by the Red de Investigación en Inflamación y Enfermedades Reumáticas (RD12/0009) of the Redes Temáticas de Investigación Cooperativa en Salud from the Instituto de Salud Carlos III Health Ministry (Spain).

Received: December 16, 2013

Accepted: February 24, 2014

Published: March 20, 2014

## Web Resources

The URLs for data presented herein are as follows:

Online Mendelian Inheritance in Man (OMIM), <http://www.omim.org/>

Protein Data Bank (PDB), <http://www.rcsb.org/pdb/home/home.do>

## References

1. Daha, N.A., and Toes, R.E.M. (2011). Rheumatoid arthritis: Are ACPA-positive and ACPA-negative RA the same disease? *Nat. Rev. Rheumatol.* **7**, 202–203.
2. van der Helm-van Mil, A.H., and Huizinga, T.W. (2008). Advances in the genetics of rheumatoid arthritis point to subclassification into distinct disease subsets. *Arthritis Res. Ther.* **10**, 205.
3. Eyre, S., Bowes, J., Diogo, D., Lee, A., Barton, A., Martin, P., Zhernakova, A., Stahl, E., Viatte, S., McAllister, K., et al.; Biologics in Rheumatoid Arthritis Genetics and Genomics Study Syndicate; Wellcome Trust Case Control Consortium (2012). High-density genetic mapping identifies new susceptibility loci for rheumatoid arthritis. *Nat. Genet.* **44**, 1336–1340.
4. Ding, B., Padyukov, L., Lundström, E., Seielstad, M., Plenge, R.M., Oksenberg, J.R., Gregersen, P.K., Alfredsson, L., and Klareskog, L. (2009). Different patterns of associations with anti-citrullinated protein antibody-positive and anti-citrullinated protein antibody-negative rheumatoid arthritis in the extended major histocompatibility complex region. *Arthritis Rheum.* **60**, 30–38.
5. Stahl, E.A., Raychaudhuri, S., Remmers, E.F., Xie, G., Eyre, S., Thomson, B.P., Li, Y., Kurzeeman, F.A.S., Zhernakova, A., Hinks, A., et al.; BIRAC Consortium; YEAR Consortium (2010). Genome-wide association study meta-analysis identifies seven new rheumatoid arthritis risk loci. *Nat. Genet.* **42**, 508–514.

6. Raychaudhuri, S., Sandor, C., Stahl, E.A., Freudenberg, J., Lee, H.-S., Jia, X., Alfredsson, L., Padyukov, L., Klareskog, L., Worthington, J., et al. (2012). Five amino acids in three HLA proteins explain most of the association between MHC and seropositive rheumatoid arthritis. *Nat. Genet.* *44*, 291–296.
7. Raychaudhuri, S., Remmers, E.F., Lee, A.T., Hackett, R., Guiducci, C., Burt, N.P., Gianniny, L., Korman, B.D., Padyukov, L., Kurreeman, F.A.S., et al. (2008). Common variants at CD40 and other loci confer risk of rheumatoid arthritis. *Nat. Genet.* *40*, 1216–1223.
8. Lundberg, K., Bengtsson, C., Kharlamova, N., Reed, E., Jiang, X., Källberg, H., Pollak-Dorocic, I., Israelsson, L., Kessel, C., Padyukov, L., et al. (2013). Genetic and environmental determinants for disease risk in subsets of rheumatoid arthritis defined by the anticitrullinated protein/peptide antibody fine specificity profile. *Ann. Rheum. Dis.* *72*, 652–658.
9. Wiik, A.S., van Venrooij, W.J., and Pruijn, G.J.M. (2010). All you wanted to know about anti-CCP but were afraid to ask. *Autoimmun. Rev.* *10*, 90–93.
10. van der Linden, M.P.M., van der Woude, D., Ioan-Facsinay, A., Levarht, E.W.N., Stoeken-Rijsbergen, G., Huizinga, T.W.J., Toes, R.E.M., and van der Helm-van Mil, A.H.M. (2009). Value of anti-modified citrullinated vimentin and third-generation anti-cyclic citrullinated peptide compared with second-generation anti-cyclic citrullinated peptide and rheumatoid factor in predicting disease outcome in undifferentiated arthritis and rheumatoid arthritis. *Arthritis Rheum.* *60*, 2232–2241.
11. Viatte, S., Plant, D., and Raychaudhuri, S. (2013). Genetics and epigenetics of rheumatoid arthritis. *Nat. Rev. Rheumatol.* *9*, 141–153.
12. Cortes, A., Hadler, J., Pointon, J.P., Robinson, P.C., Karaderi, T., Leo, P., Cremin, K., Pryce, K., Harris, J., Lee, S., et al.; International Genetics of Ankylosing Spondylitis Consortium (IGAS); Australo-Anglo-American Spondyloarthritis Consortium (TASC); Groupe Française d'Etude Génétique des Spondylarthrites (GFECS); Nord-Trøndelag Health Study (HUNT); Spondyloarthritis Research Consortium of Canada (SPARCC); Wellcome Trust Case Control Consortium 2 (WTCCC2) (2013). Identification of multiple risk variants for ankylosing spondylitis through high-density genotyping of immune-related loci. *Nat. Genet.* *45*, 730–738.
13. Verpoort, K.N., van Gaalen, F.A., van der Helm-van Mil, A.H.M., Schreuder, G.M.T., Breedveld, F.C., Huizinga, T.W.J., de Vries, R.R.P., and Toes, R.E.M. (2005). Association of HLA-DR3 with anti-cyclic citrullinated peptide antibody-negative rheumatoid arthritis. *Arthritis Rheum.* *52*, 3058–3062.
14. Irigoyen, P., Lee, A.T., Wener, M.H., Li, W., Kern, M., Batliwalla, F., Lum, R.F., Massarotti, E., Weisman, M., Bombardier, C., et al. (2005). Regulation of anti-cyclic citrullinated peptide antibodies in rheumatoid arthritis: contrasting effects of HLA-DR3 and the shared epitope alleles. *Arthritis Rheum.* *52*, 3813–3818.
15. Jia, X., Han, B., Onengut-Gumuscu, S., Chen, W.-M., Concannon, P.J., Rich, S.S., Raychaudhuri, S., and de Bakker, P.I.W. (2013). Imputing amino acid polymorphisms in human leukocyte antigens. *PLoS ONE* *8*, e64683.
16. Lee, H.-S., Irigoyen, P., Kern, M., Lee, A., Batliwalla, F., Khalili, H., Wolfe, F., Lum, R.F., Massarotti, E., Weisman, M., et al. (2007). Interaction between smoking, the shared epitope, and anti-cyclic citrullinated peptide: a mixed picture in three large North American rheumatoid arthritis cohorts. *Arthritis Rheum.* *56*, 1745–1753.
17. Klareskog, L., Stolt, P., Lundberg, K., Källberg, H., Bengtsson, C., Grunewald, J., Rönnelid, J., Harris, H.E., Ulfgren, A.-K., Rantapää-Dahlqvist, S., et al. (2006). A new model for an etiology of rheumatoid arthritis: smoking may trigger HLA-DR (shared epitope)-restricted immune reactions to autoantigens modified by citrullination. *Arthritis Rheum.* *54*, 38–46.
18. Browning, B.L., and Browning, S.R. (2009). A unified approach to genotype imputation and haplotype-phase inference for large data sets of trios and unrelated individuals. *Am. J. Hum. Genet.* *84*, 210–223.
19. Brown, W.M., Pierce, J., Hilner, J.E., Perdue, L.H., Lohman, K., Li, L., Venkatesh, R.B., Hunt, S., Mychaleckyj, J.C., and Deloukas, P.; Type 1 Diabetes Genetics Consortium (2009). Overview of the MHC fine mapping data. *Diabetes Obes. Metab.* *11 (Suppl 1)*, 2–7.
20. Price, A.L., Patterson, N.J., Plenge, R.M., Weinblatt, M.E., Shadick, N.A., and Reich, D. (2006). Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet.* *38*, 904–909.
21. Brown, M.A., Pile, K.D., Kennedy, L.G., Calin, A., Darke, C., Bell, J., Wordsworth, B.P., and Cornélis, F. (1996). HLA class I associations of ankylosing spondylitis in the white population in the United Kingdom. *Ann. Rheum. Dis.* *55*, 268–270.
22. Reveille, J.D., Sims, A.M., Danoy, P., Evans, D.M., Leo, P., Pointon, J.J., Jin, R., Zhou, X., Bradbury, L.A., Appleton, L.H., et al.; Australo-Anglo-American Spondyloarthritis Consortium (TASC) (2010). Genome-wide association study of ankylosing spondylitis identifies non-MHC susceptibility loci. *Nat. Genet.* *42*, 123–127.
23. Tiilikainen, A., Lassus, A., Karvonen, J., Vartiainen, P., and Julin, M. (1980). Psoriasis and HLA-Cw6. *Br. J. Dermatol.* *102*, 179–184.
24. Nair, R.P., Stuart, P.E., Nistor, I., Hiremagalore, R., Chia, N.V.C., Jenisch, S., Weichenthal, M., Abecasis, G.R., Lim, H.W., Christophers, E., et al. (2006). Sequence and haplotype analysis supports HLA-C as the psoriasis susceptibility 1 gene. *Am. J. Hum. Genet.* *78*, 827–851.
25. Ho, P.Y.P.C., Barton, A., Worthington, J., Thomson, W., Silman, A.J., and Bruce, I.N. (2007). HLA-Cw6 and HLA-DRB1\*07 together are associated with less severe joint disease in psoriatic arthritis. *Ann. Rheum. Dis.* *66*, 807–811.
26. Karlson, E.W., Chibnik, L.B., Kraft, P., Cui, J., Keenan, B.T., Ding, B., Raychaudhuri, S., Klareskog, L., Alfredsson, L., and Plenge, R.M. (2010). Cumulative association of 22 genetic variants with seropositive rheumatoid arthritis risk. *Ann. Rheum. Dis.* *69*, 1077–1085.
27. Morrison, A.C., Bare, L.A., Chambless, L.E., Ellis, S.G., Malloy, M., Kane, J.P., Pankow, J.S., Devlin, J.J., Willerson, J.T., and Boerwinkle, E. (2007). Prediction of coronary heart disease risk using a genetic risk score: the Atherosclerosis Risk in Communities Study. *Am. J. Epidemiol.* *166*, 28–35.
28. Meigs, J.B., Shrader, P., Sullivan, L.M., McAteer, J.B., Fox, C.S., Dupuis, J., Manning, A.K., Florez, J.C., Wilson, P.W.F., D'Agostino, R.B., Sr., and Cupples, L.A. (2008). Genotype score in addition to common risk factors for prediction of type 2 diabetes. *N. Engl. J. Med.* *359*, 2208–2219.
29. Price, P., Witt, C., Allcock, R., Sayer, D., Garlepp, M., Kok, C.C., French, M., Mallal, S., and Christiansen, F. (1999). The genetic basis for the association of the 8.1 ancestral haplotype (A1, B8, DR3) with multiple immunopathological diseases. *Immunol. Rev.* *167*, 257–274.

30. Pettersen, E.F., Goddard, T.D., Huang, C.C., Couch, G.S., Greenblatt, D.M., Meng, E.C., and Ferrin, T.E. (2004). UCSF Chimera—a visualization system for exploratory research and analysis. *J. Comput. Chem.* *25*, 1605–1612.
31. Pereyra, F., Jia, X., McLaren, P.J., Telenti, A., de Bakker, P.I., Walker, B.D., Ripke, S., Brumme, C.J., Pulit, S.L., Carrington, M., et al.; International HIV Controllers Study (2010). The major genetic determinants of HIV-1 control affect HLA class I peptide presentation. *Science* *330*, 1551–1557.
32. Gregersen, P.K., Kosoy, R., Lee, A.T., Lamb, J., Sussman, J., McKee, D., Simpfendorfer, K.R., Pirskanen-Matell, R., Piehl, F., Pan-Hammarstrom, Q., et al. (2012). Risk for myasthenia gravis maps to a (151) Pro→Ala change in TNIP1 and to human leukocyte antigen-B\*08. *Ann. Neurol.* *72*, 927–935.
33. Shi, J., Knevel, R., Suwannalai, P., van der Linden, M.P., Jansen, G.M.C., van Veelen, P.A., Levarht, N.E.W., van der Helm-van Mil, A.H.M., Cerami, A., Huizinga, T.W.J., et al. (2011). Autoantibodies recognizing carbamylated proteins are present in sera of patients with rheumatoid arthritis and predict joint damage. *Proc. Natl. Acad. Sci. USA* *108*, 17372–17377.
34. Boire, G., Ménard, H.A., Gendron, M., Lussier, A., and Myhal, D. (1993). Rheumatoid arthritis: anti-Ro antibodies define a non-HLA-DR4 associated clinicoserological cluster. *J. Rheumatol.* *20*, 1654–1660.