

Epigenome-wide association data implicate DNA methylation as an intermediary of genetic risk in rheumatoid arthritis

Yun Liu^{1,2,12}, Martin J Aryee^{1,3,12}, Leonid Padyukov^{4,5,12}, M Daniele Fallin^{1,6,7,12}, Espen Hesselberg^{4,5}, Arni Runarsson^{1,2}, Lovisa Reinius⁸, Nathalie Acevedo⁹, Margaret Taub^{1,6}, Marcus Ronninger^{4,5}, Klementy Shchetynsky^{4,5}, Annika Scheynius⁹, Juha Kere⁸, Lars Alfredsson¹⁰, Lars Klareskog^{4,5}, Tomas J Ekström^{5,11} & Andrew P Feinberg^{1,2,6}

Epigenetic mechanisms integrate genetic and environmental causes of disease, but comprehensive genome-wide analyses of epigenetic modifications have not yet demonstrated robust association with common diseases. Using Illumina HumanMethylation450 arrays on 354 anti-citrullinated protein antibody-associated rheumatoid arthritis cases and 337 controls, we identified two clusters within the major histocompatibility complex (MHC) region whose differential methylation potentially mediates genetic risk for rheumatoid arthritis. To reduce confounding factors that have hampered previous epigenome-wide studies, we corrected for cellular heterogeneity by estimating and adjusting for cell-type proportions in our blood-derived DNA samples and used mediation analysis to filter out associations likely to be a consequence of disease. Four CpGs also showed an association between genotype and variance of methylation. The associations for both clusters replicated at least one CpG ($P < 0.01$), with the rest showing suggestive association, in monocyte cell fractions in an independent cohort of 12 cases and 12 controls. Thus, DNA methylation is a potential mediator of genetic risk.

Epigenetic mechanisms can cause durable changes of gene expression that are heritable during cell division by covalent modifications of DNA bases and potentially other chromatin alterations. They might influence disease development in a manner complementary to direct mutations of the DNA sequence. The role of epigenetic modifications in cancer etiology and progression is well established¹, and a number of small surveys of DNA methylation in common disease have been carried out^{2–5}. We and others have suggested that genetic and epigenetic modifications could interact biologically^{6–8}, and that methylation analysis might uncover heritable genetic variants contributing to disease that are invisible to conventional genome-wide association studies (GWAS). A comprehensive genome-wide methylation analysis has not yet demonstrated robust association of specific methylation alterations with a common disease, however. Several limitations to such studies may explain why not, including (i) the cellular heterogeneity of the sample material, and (ii) the potential for methylation changes that are a consequence of a disease rather than part of its etiology. Here, we apply a series of *ad hoc* filtering steps, which address these issues, to identify CpG methylation that probably mediates genetic risk for rheumatoid arthritis, from genome-wide

epigenetic and genetic data. This process may serve as a guidepost for epigenetic epidemiological studies generally.

Rheumatoid arthritis is a complex and heterogeneous disease, where onset as well as disease course is dependent on interactions between different genetic and environmental or life-style factors^{9,10}. Several meta-analyses of GWAS have identified close to 40 genetic variants that confer risk for the anti-citrullinated protein antibody (ACPA)-associated subtype of rheumatoid arthritis^{11–14}. However, the fact that these discoveries can only explain <20% of disease variance suggests that other factors are likely involved in the disease¹³.

Two additional factors make rheumatoid arthritis an ideal test case for analyzing the relationships between genes, methylation and disease pathogenesis. In rheumatoid arthritis, one of the main classes of cells involved in the disease, leukocytes, is readily available for DNA methylation analysis and disease state can be reproducibly determined by the presence of antibodies to citrullinated protein antigens.

In our present study, 354 rheumatoid arthritis patients (cases) with citrullinated protein antibodies¹⁵ and 337 healthy individuals (controls) were selected from the epidemiological investigation of rheumatoid arthritis (EIRA)^{16,17}, a Swedish population-based

¹Center for Epigenetics, Johns Hopkins University School of Medicine, Baltimore, Maryland, USA. ²Department of Medicine, Johns Hopkins University School of Medicine, Baltimore, Maryland, USA. ³Department of Oncology, Johns Hopkins University School of Medicine, Baltimore, Maryland, USA. ⁴Rheumatology Unit, Department of Medicine, Karolinska Institute, Stockholm, Sweden. ⁵Center for Molecular Medicine, Karolinska Institute, Stockholm, Sweden. ⁶Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, Baltimore, Maryland, USA. ⁷Department of Epidemiology, Johns Hopkins Bloomberg School of Public Health, Baltimore, Maryland, USA. ⁸Department of Biosciences and Nutrition, Karolinska Institute, Stockholm, Sweden. ⁹Department of Medicine Solna, Karolinska Institute, Stockholm, Sweden. ¹⁰Institute of Environmental Medicine, Karolinska Institute, Stockholm, Sweden. ¹¹Department of Clinical Neuroscience, Karolinska Institute, Stockholm, Sweden. ¹²These authors contributed equally to this work. Correspondence should be addressed to L.K. (lars.klareskog@ki.se), T.J.E. (tomas.ekstrom@ki.se) or A.P.F. (afeinberg@jhu.edu).

Received 9 November 2012; accepted 18 December 2012; published online 20 January 2013; doi:10.1038/nbt.2487

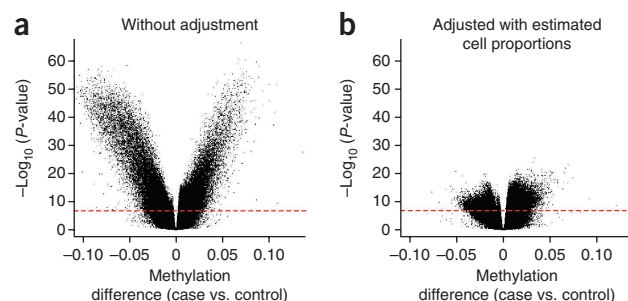
Figure 1 Differential cell counts in identifying rheumatoid arthritis–associated differentially methylated positions (DMPs). (a,b) Volcano plot of $-\log_{10}$ (P -value) against beta, representing the methylation difference between ACPA–rheumatoid arthritis cases and controls, without (a) or with (b) adjusting for differential cell counts, in addition to controlling for age, sex and smoking status. The dashed red lines represent the threshold used for statistical cutoff (Bonferroni-adjusted $P = 0.05$).

case-control study. Cases were recruited at the first visit to a rheumatology clinic before initiation of treatment with disease-modifying small-molecule or biological agents. At this first visit, blood samples were collected for DNA analysis and serology^{16,18}. Control subjects were selected from the same study to match rheumatoid arthritis cases in terms of age, gender, smoking status and residential area at the time of diagnosis (**Supplementary Table 1**). An additional advantage of these samples is that genome-wide single-nucleotide polymorphism (SNP) data were available on the same individuals, enabling us to determine the relationship between genotype, epigenotype and phenotype. On these 691 samples, we first performed genome-wide DNA methylation analysis using the Illumina 450K methylation array, to identify epigenetic differences associated with rheumatoid arthritis. After excluding two samples with poor quality and 187,468 probes containing SNPs, which might affect the measurement of DNA methylation, the final data set used for downstream analysis comprised 354 cases and 335 controls for 298,109 CpG positions (Online Methods).

RESULTS

Correcting for cellular heterogeneity

Our first challenge is that the DNA samples available for methylation analysis are generally derived from heterogeneous cell populations. For example, the DNA samples most readily available from large numbers of individuals are from whole blood, which consists of many distinct populations in varying proportions. It has been shown that these functionally distinct populations have unique DNA methylation signatures¹⁹, thus cell heterogeneity may act as a potential confounder when investigating DNA methylation differences between cases and controls, if cell distribution itself differs by disease status. To address



this, we attempted to adjust for cell proportion using linear regression models in our epigenetic association analysis. To obtain sample-specific estimates of the proportion of the major cell types in blood, we applied a statistical algorithm²⁰ that uses reference information on cell-specific methylation signatures to estimate cell proportions from genome-scale methylation data. The estimated cell type distribution from this algorithm was consistent with our experimental results from flow cytometry and did show that patterns in rheumatoid arthritis cases were distinct from those of controls (**Supplementary Table 2**), suggesting that it is critical to adjust for cell type distribution in the downstream analyses. For example, **Figure 1** shows epigenome-wide association results before and after adjustment using estimated cell proportions, showing a notable reduction in association signals after adjustment.

Establishing epigenetic mediation of genetic risk

Our second challenge is that many methylation differences are probably a consequence of rheumatoid arthritis. To filter these out and reveal biology related to the cause, we applied methods from the causal inference literature^{21–24}. In this approach, which employs a series of conditional correlation analyses, one considers the possible directed relationships between a causal factor, a potential mediator and an outcome (**Fig. 2a**).

As we were particularly interested in epigenetic marks that may mediate the genetic risk for rheumatoid arthritis, we applied this method with genotype as a causal factor (G, **Fig. 2a**), DNA methylation as a potential mediator (M, **Fig. 2a**) and rheumatoid arthritis as the outcome (Y, **Fig. 2a**). We developed a three-step filtering process followed by the application of the causal inference test (CIT)²⁴ to identify the differentially methylated positions (DMPs) associated with rheumatoid arthritis that are most likely to be acting as mediators of genetic risk rather than being consequences of rheumatoid arthritis. These three filtering steps are (i) establish the relationship between the potential mediator (M) and outcome (Y); (ii) from these, establish the relationship between the primary cause (G) and the mediator (M); and (iii) from these, establish the relationship between the cause (G) and outcome (Y) (**Fig. 2b**). We then applied the CIT to establish that methylation (M) is

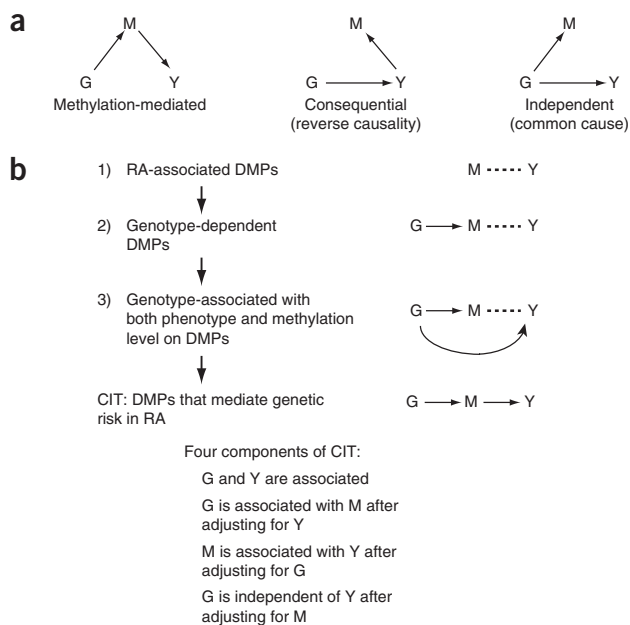


Figure 2 Identification of epigenetically mediated genetic risk factors for rheumatoid arthritis disease. (a) Possible relationships between a causal factor (G), a possible mediator (M) and an outcome (Y). Left, the methylation-mediated relationship, in which genotype (G) acts on phenotype (Y) through methylation (M). Middle, the consequential methylation model, in which DNA methylation (M) changes are the consequence of phenotype (Y). Right, the independent model, in which the genotype (G) acts on DNA methylation (M) and phenotype (Y) independently. (b) Summary workflow for identifying epigenetically mediated genetic risk factors for rheumatoid arthritis. The diagrams on the right represent the relationships between genotype (G), DNA methylation (M) and rheumatoid arthritis phenotype (Y). Dashed lines, the association relationship; arrows, the causal relationship. RA, rheumatoid arthritis.

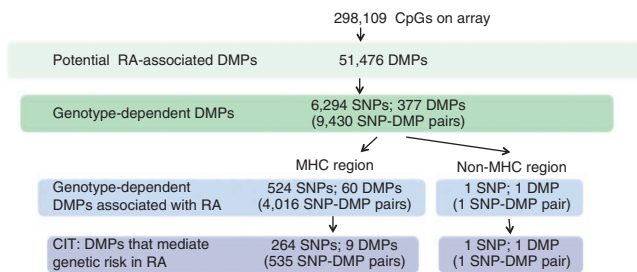


Figure 3 Summary workflow and results for identifying epigenetically mediated genetic risk factors for rheumatoid arthritis.

the mediator between the cause (G) and outcome (Y). Our approach results in a candidate set of mediators, but it should be noted that, as in all epidemiological studies, it is impossible to conclusively prove causal relationships on the basis of observational data alone.

Briefly, the CIT requires the following criteria²⁴. G and Y are associated; G is associated with M after adjusting for Y; M is associated with Y after adjusting for G; and G is independent of Y after adjusting for M. If methylation is a consequence of Y (Fig. 2a, middle panel) or independently controlled by G (Fig. 2a, right panel), rather than a mediator in the path from G to Y (Fig. 2a, left panel), the estimated effect of G on Y should not be affected by conditioning on M. However, if methylation is indeed a mediator, this conditioning should drastically reduce the observed effect of G on Y (Fig. 2a, left panel).

In our first filtering step, we performed epigenome-wide association analysis using adjustment for estimated cell proportions as well as age, sex and smoking status in the context of a linear model. We took all

DMPs, putatively associated with rheumatoid arthritis, that achieved a $P < 0.05$ after Bonferroni correction, into the next step (51,476; step 1, Fig. 3 and Supplementary Table 3). We then performed genome-wide SNP association analysis for each of these DMPs (step 2, Fig. 3) to identify the subset where methylation appears to be under genetic control. We fit an allelic dosage model for each DMP and each of 1,196,263 SNPs (300,987 genotyped SNPs and 895,276 SNPs imputed from the HapMap 3 panel)¹¹ and identified 9,430 SNP-DMP pairs (Supplementary Table 4) with genome-wide significance (Bonferroni-adjusted $P < 0.05$). These SNP-DMP pairs comprised 6,294 unique SNPs and 377 unique DMPs (step 2, Fig. 3). More than half of SNP-DMP pairs are spread over a 5-Mb region covering the MHC cluster, which is known to harbor several loci associated with rheumatoid arthritis risk¹⁴. Given that the MHC cluster was heavily over-represented, we decided to analyze the MHC region and non-MHC region separately.

In-depth analysis of the MHC region

It has been shown that the major genetic risk loci for seropositive rheumatoid arthritis are located within the MHC region, which accounts for >10% of the phenotypic variance¹⁴. Given that the MHC cluster was heavily over-represented in the results from the genome-wide SNP-DMP scan, we decided to explore this region in detail using increased-density genotype data, which were imputed based on a large reference panel¹⁴. We again fit an allelic dosage model for each rheumatoid arthritis-associated DMP and each of 5,009 imputed SNPs within the MHC region and identified 7,242 significantly associated SNP-DMP pairs (Bonferroni-adjusted $P < 0.05$) in the MHC (Supplementary Table 5). These SNP-DMP pairs comprised 1,952 unique SNPs and 76 unique DMPs. We then carried

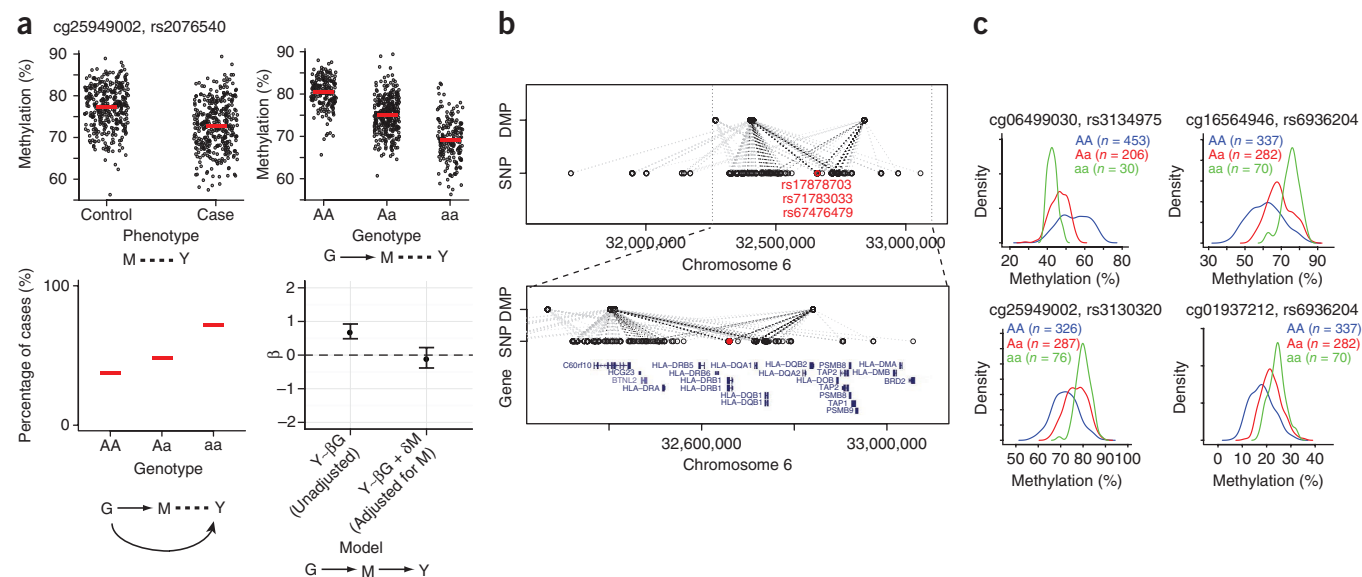


Figure 4 Genotype-dependent candidate DMPs that mediate genetic risk within the MHC region. (a) Top panels: association between DNA methylation level at a DMP that mediates genetic risk in rheumatoid arthritis and phenotype (top left) or genotype (top right). Red horizontal bars mark average DNA methylation levels. Bottom left panel: association between genotype and rheumatoid arthritis phenotype. Red horizontal bars mark percentage of cases for each genotype. Bottom right panel: coefficient (β) represents the dependence of rheumatoid arthritis phenotype (Y) on genotype (G), with or without adjusting for DNA methylation (M). The error bars represent the 95% confidence interval for the coefficient, β . In the case of the methylation-mediated model, the absolute value of the observed G:Y relationship strength reduces toward zero when adjusting for methylation (M). (b) Association between candidate genetic risk-mediating DMPs and genotype within a 1.5 Mb section of the MHC region. Each dashed line represents a potential mediation relationship between a SNP and a DMP as determined by the CIT. The color of the line indicates significance of the P -value for the statistical mediation test. Genotype data were imputed based on a large reference panel. Three previously identified rheumatoid arthritis-associated genetic variants are illustrated in red. Bottom panel: zoomed-in images of top panel with gene annotations. (c) Examples of rheumatoid arthritis-associated DMP in which genotype is associated with the change in both mean and variance of DNA methylation. The methylation density plot is color coded by genotype. The number of individuals in each genotype group is shown on the corner.

Table 1 DMPs that mediate genetic risk in rheumatoid arthritis

	Rheumatoid arthritis-associated CpGs (DMPs)			SNPs associated with CpG methylation					
	DMP	Beta ^a	<i>P</i> -value (meth vs. pheno)	Gene name	SNP ^b	<i>P</i> -value (geno vs. meth)	Adjusted <i>P</i> -value (geno vs. pheno)	<i>P</i> -value (CIT)	
MHC	cg21325723	-0.028	1.49E-09	<i>C6orf10</i>	rs2395163	<2E-16	1.00E-04	2.96E-19	
	cg16609995	0.015	2.88E-09	<i>PBX2</i>	DRB1_AA104_E2_32659926_AE	1.33E-15	3.00E-04	5.95E-07	
	cg06499030	0.038	4.01E-09	<i>HLA-DQB2</i>	CHR6_POS32657567	<2E-16	1.00E-04	6.62E-15	
	cg16564946	-0.046	9.74E-09	<i>C6orf10</i>	rs9267954	<2E-16	1.00E-04	7.30E-08	
	cg25949002	-0.026	2.57E-08	<i>C6orf10</i>	rs2076540	<2E-16	1.00E-04	1.24E-10	
	cg14704780	-0.066	2.87E-08	<i>C6orf10</i>	rs3916765	<2E-16	1.00E-04	8.72E-10	
	cg19555708	0.012	5.11E-08	<i>GPSM3</i>	DRB1_AA104_E2_32659926_AE	8.88E-16	3.00E-04	5.95E-07	
	cg01937212	-0.021	6.00E-08	<i>C6orf10</i>	rs477005	<2E-16	1.00E-04	5.39E-10	
	cg19321684	0.023	8.26E-08	<i>GPSM3</i>	DRB1_AA104_E2_32659926_AE	<2E-16	3.00E-04	5.95E-07	
	Non-MHC	cg00462104	-0.023	8.84E-08	<i>GSTA2</i>	rs3996993	<2E-16	0.038	0.0016

Geno, genotype; meth, methylation; pheno, rheumatoid arthritis phenotype.

^aAdjusted methylation difference between cases and controls.

^bFor each DMP, only the SNP with smallest CIT *P*-value is shown here. For the full SNP list, see **Supplementary Table 7**.

out prefiltering step 3 and the CIT, outlined above, to distinguish between methylation that is a result of disease and methylation that may be part of the causal path to disease.

We tested the association between each of 1,952 SNPs and rheumatoid arthritis using an allelic dosage model. As all of these SNPs are located within the MHC region and many are correlated, we permuted adjusted *P*-values using the step-down maxT multiple testing procedure to control the family wise type I error rate²⁵. Of the 1,952 SNPs, we identified 524 that were significantly associated with rheumatoid arthritis at an adjusted *P* < 0.05. These 524 SNPs form 4,016 SNP-DMP pairs with 60 unique DMPs (step 3, **Fig. 3** and **Supplementary Table 6**). By performing the CIT test²⁴, we found that for 535 of the 4,016 SNP-DMP pairs, the SNP effect on rheumatoid arthritis was reduced after adjusting for methylation (Bonferroni-adjusted CIT *P* < 0.05; Online Methods), suggesting mediation (**Fig. 4a**). These 535 MHC SNP-DMP pairs comprised 264 unique SNPs and 9 unique DMPs (CIT in **Fig. 3**, **Table 1** and **Supplementary Table 7**) and represent potential methylation-mediated relationships between SNPs and rheumatoid arthritis disease risk.

Genetic control of mean and variance of methylation

It has been demonstrated that variations in genes coding for five different amino acid positions in the binding grooves of HLA-DR, HLA-DP and HLA-B account for most hitherto described associations between genetic differences in the MHC region and ACPA-rheumatoid arthritis¹⁴. Three of these five genetic variants are located within the *HLA-DRB1* gene and show significant (Bonferroni-adjusted *P* < 0.05) evidence of association with DNA methylation loci that appear to at least partially mediate the genetic risk effect (**Fig. 4b**). Some of these rheumatoid arthritis-associated DMPs are located >100 kb from the associated genetic risk variants, possibly as a consequence of the relative sparseness in CpG coverage on the Illumina 450K methylation array.

We proposed recently that genetic variants might regulate phenotypic variability in addition to mean phenotype and that this connection between genotype and phenotypic plasticity would be mediated epigenetically²⁶. Such a mechanism would provide a non-Lamarckian basis for an epigenetic role in natural selection because the variants themselves would be transmitted genetically, but they would also allow increased phenotypic plasticity in response to a varying environment²⁶. DNA methylation represents a promising candidate for mediating such plasticity, as methylation levels are measured as

proportions where variance changes are intrinsically related to mean shifts. Notably, five out of nine DMPs identified here also showed a significant association between genotype and variance of methylation, as suggested by the model (**Fig. 4c** and **Supplementary Table 8**).

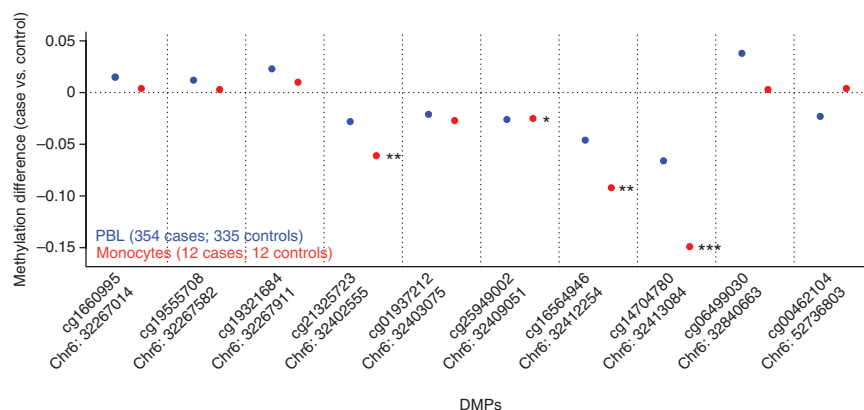
Methylation-mediated genetic risks in non-MHC regions

We further analyzed the 4,619 SNP-DMP pairs (including 4,540 SNPs and 343 DMPs) outside of the MHC region that were identified in the genome-wide scan. Of the 4,540 SNPs, one SNP is significantly associated with rheumatoid arthritis phenotype (step-down maxT adjusted *P* < 0.05) (step 3, **Fig. 3**). Using the CIT²⁴ as described before, we concluded that the effect of genotype on rheumatoid arthritis risk appears to be mediated by a DNA methylation change in the promoter region of a gene called *GSTA2* (**Supplementary Fig. 1** and **Table 1**). *GSTA2* belongs to the glutathione *S*-transferase supergene family, which is important in the detoxification of electrophilic compounds, including environmental toxins (such as tobacco smoke)^{27–29}. Polymorphisms within the μ (*GSTM1*), θ (*GSTT1*) and π (*GSTP1*) classes of GST previously have been identified and are associated with rheumatoid arthritis risk and severity^{30–33}, although *GSTA2* has not previously been implicated in rheumatoid arthritis through GWAS¹³. This underscores the usefulness of our approach. In fact, we did not identify any non-MHC genetic variants associated with rheumatoid arthritis by fitting a standard allelic dosage model at each SNP within our own samples (**Supplementary Fig. 2**).

Replication of methylation differences in monocytes

Having identified ten DMPs (nine within the MHC region and one outside) whose methylation level putatively mediates genetic risk in rheumatoid arthritis (**Table 1**), we attempted to replicate these methylation differences in fresh flow cytometry-sorted cell populations from untreated rheumatoid arthritis cases and controls. We separated peripheral blood lymphocytes (PBLs) from 12 case-control pairs into separate cell fractions. For the monocyte cell fraction, nine out of ten DMPs showed methylation changes in the same direction as that seen in the large-scale PBL analysis. Of these, three of the CpG sites were significant at *P* < 0.05, one at *P* = 0.063 and a fifth at 0.11, even with this small sample number (**Supplementary Table 9**). The three with greatest significance also showed larger beta values than seen in PBL (**Fig. 5**), suggesting that monocytes are more proximal to the pathogenic cell type. Given that monocytes represented less than 10% of the PBL fraction, this may explain the smaller effect size seen using total PBLs.

Figure 5 Replication data for ten candidate DMPs that mediate genetic risk in rheumatoid arthritis from sorted monocytes. DNA methylation data on CD14⁺ monocytes from a replication set of 12 ACPA–rheumatoid arthritis cases and 12 controls were measured using the Illumina 450K methylation array. For PBLs, all 10 DMPs $P < 10^{-7}$. For sorted monocytes: * $P < 0.1$; ** $P < 0.01$; *** $P < 0.001$.



DISCUSSION

In summary, we have applied an approach that corrects for the confounding influence of cell heterogeneity and filters out signals likely due to the disease itself. Using a strategy of three filtering steps followed by the application of mediation analysis using the CIT algorithm, we performed genome-scale methylation and SNP analysis and identified ten putative DMPs that mediate genetic risk for rheumatoid arthritis, nine in the MHC cluster, and one outside on the same chromosome (6p12.1).

Our approach for adjusting for cell heterogeneity should be applicable for many tissue sources, if cell-specific methylation signatures for the particular mixture in question are available. Even samples from primary affected tissues tend to consist of a mixture of many cell types making an adjustment for cell type proportions a prerequisite for epigenetic association analysis, somewhat analogous to the correction for population stratification using empirically estimated ancestry proportions in GWAS studies^{34,35}. This adjustment for cell proportions does not address the question of whether the chosen tissue is the appropriate surrogate tissue for the disease in question, but simply handles the heterogeneity issue regardless of surrogate or primary status. In this report, we have assumed that blood is the primary tissue for an inflammatory disease.

Although we show that our cell type adjustment is a notable improvement over unadjusted analyses and reduces confounding by cell type bias, there may be residual confounding not fully accommodated in the specific proportion estimation and linear adjustment we pursued. Further methodological work to improve this estimation and modeling approach is important. Other sources of confounding in array analyses must also be considered. In addition to age, sex and other demographic confounders, batch effects such as date and laboratory should be evaluated. For example, although we did not anticipate strong batch effects for 450K methylation data to date, we examined our methylation data by means of principal components analysis (PCA) and did observe a relationship between date of assay and first component of PCA. Although ideally one would design assay runs to have an equivalent spread of phenotypes across dates and/or laboratories, this is often not practical. In our study, there was an imbalance between the number of cases and controls run per date, and thus batch effects in these 450K data could potentially confound associations³⁶. To address these issues, we re-analyzed the results for our top ten CpGs using a procedure that simultaneously corrects for batch and cell type composition (Online Methods and **Supplementary Fig. 3**). Although the statistical significance of all CpGs is reduced by this adjustment, it is notable that the five CpGs in the two regions that were replicated in flow-sorted monocytes retained the strongest effect size (**Supplementary Fig. 4** and **Supplementary Table 10**). We would recommend in future studies addressing batch effect issues by principal components analysis or surrogate variable analysis (SVA)³⁷, which is designed to identify and estimate the sources of heterogeneity that are not captured by variables included in the model, as a first step. This may in fact improve cell proportion estimation for subsequent

adjustment. We also noted that although our approach to cell proportion adjustment, which is a convenient tool for blood-derived samples, is a considerable improvement over no adjustment, residual confounding due to cell type may remain. This can be dealt with by replication in cell-sorted samples as we have shown, or through advanced estimation and statistical adjustment methods—a call for additional methodological work.

Another issue that complicates epigenetic studies over purely genetic analysis is that the primary tissue may harbor DNA methylation changes that are a consequence of the disease, rather than a marker of causal mechanism. To address this, we applied mediation approaches already used in the gene expression and epidemiology literature, but not previously applied to epigenetic studies. We emphasize that our findings, as in all epidemiological studies, are hypotheses that will ultimately require verification in independent and/or mechanistic studies. In particular, there exist conditions, such as the presence of unmeasured confounders, where it may be impossible to distinguish causal from consequential methylation events based on observational data alone³⁸. Although much will need to be worked out over time, just as it was in the development of GWAS, we feel that our approach directly addresses the fundamental question of epigenetic epidemiology, that is, how one can link genetics to epigenetics to phenotype. Similarly, mediation analysis can be applied to the other component of epigenetic epidemiology—the role of the environment—if one assumes that the environmental factors are causal in the disease.

It is notable that our top ten CpGs represent signals across five genomic regions, and the five CpGs that replicate most robustly in monocytes cluster in two regions (**Supplementary Fig. 4**). This supports use of region-based statistical approaches such as “bump hunting”³⁹ epigenetic association analyses and further suggests that denser coverage than the 450K array will be better in identifying methylation differences moving forward either by a new array design or capture bisulfite sequencing. It is notable that monocyte subfractions showed effect sizes comparable to unfractionated PBLs, with statistical significance for three of the DMPs and marginal significance for another two, with only 12 case-control pairs, supporting a role for monocytes in rheumatoid arthritis pathogenesis, something that is also suggested from many previous cell biologic studies in rheumatoid arthritis⁴⁰. In addition, MHC class II gene expression in macrophages, which are derived from monocytes, show a strong relationship to rheumatoid arthritis progression⁴¹.

A byproduct of the analysis presented here was the identification of suggestive evidence for vSNPs for epigenetic modification, that is, SNPs regulating DNA methylation variation. These vSNPs are predicted by a model we proposed in which genetic variants might

increase epigenetic plasticity, providing a non-Lamarckian basis for an epigenetic role in natural selection²⁶. They included five of the nine DMPs identified in the MHC region.

This research also makes a prediction that is beyond the scope of the current experiments. Given that genetic association in the MHC cluster with rheumatoid arthritis has already been shown to be linked to specific HLA protein epitopes, the methylation mediation we observe implies an additional complementary mechanism for rheumatoid arthritis, for example, basal levels of gene expression, expression in response to antigen provocation or alternative splicing, as both gene expression and splicing can be regulated by DNA methylation. Given that there are >10 genes whose promoters are within 100-kb distance of the identified DMPs (**Supplementary Table 11**) and >50 genes within the region defined by the SNP-DMP-phenotype associations reported here, at least one of these genes should show altered regulation related to DNA methylation in rheumatoid arthritis, in addition to the linear gene-protein relationships already known.

METHODS

Methods and any associated references are available in the [online version of the paper](#).

Accession code. The Illumina 450K array data, GEO: [GSE42861](#).

Note: Supplementary information is available in the [online version of the paper](#).

ACKNOWLEDGMENTS

We thank M. Rosenblum for helpful discussions on the application of statistical mediation methodology. We thank R.A. Irizarry for his contributions to the concepts in this work, his statistical insights on batch effects and helpful comments. We thank E.A. Houseman for codes used for estimating cell proportions. We also thank the EIRA study group¹⁸ for contributing invaluable clinical information and biological samples. This work was supported by the US National Institutes of Health Centers of Excellence in Genomic Science, 5P50HG003233 to A.P.F., and by the Swedish Research Council, the Swedish Combine project, the Swedish Strategic Foundations, the AFA Insurance and the European Research Council (ERC).

AUTHOR CONTRIBUTIONS

Y.L. performed the experiments. Y.L. and M.J.A. analyzed data. L.P. and L.A. performed epidemiological data collection and evaluation. L.P. did sample genotyping and genotype imputation. M.D.F. performed epidemiology analysis and data interpretation, and assisted in experimental design. L.P., M.R., K.S. and E.H. prepared nucleic acids and/or cell sorting. A.R. performed the 450K arrays. M.T. assisted in statistical analysis. J.K., L.R., N.A. and A.S. provided reference normal 450K data from sorted cells for estimating cell proportions. Y.L., M.J.A., L.P., M.D.F., L.K., T.J.E. and A.P.F. conceived the experiments and wrote the paper.

COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Published online at <http://www.nature.com/doi/10.1038/nbt.2487>.

Reprints and permissions information is available online at <http://www.nature.com/reprints/index.html>.

- Feinberg, A.P. & Tycko, B. The history of cancer epigenetics. *Nat. Rev. Cancer* **4**, 143–153 (2004).
- Kaminsky, Z.A. *et al.* DNA methylation profiles in monozygotic and dizygotic twins. *Nat. Genet.* **41**, 240–245 (2009).
- Feinberg, A.P. *et al.* Personalized epigenomic signatures that are stable over time and covary with body mass index. *Sci. Transl. Med.* **2**, 49ra67 (2010).
- Javierre, B.M. *et al.* Changes in the pattern of DNA methylation associate with twin discordance in systemic lupus erythematosus. *Genome Res.* **20**, 170–179 (2010).
- Rakyan, V.K. *et al.* Identification of type 1 diabetes-associated DNA methylation variable positions that precede disease diagnosis. *PLoS Genet.* **7**, e1002300 (2011).
- Bjornsson, H.T., Fallin, M.D. & Feinberg, A.P. An integrated epigenetic and genetic approach to common human disease. *Trends Genet.* **20**, 350–358 (2004).
- Bjornsson, H.T. *et al.* Intra-individual change over time in DNA methylation with familial clustering. *J. Am. Med. Assoc.* **299**, 2877–2883 (2008).
- Rakyan, V.K., Down, T.A., Balding, D.J. & Beck, S. Epigenome-wide association studies for common human diseases. *Nat. Rev. Genet.* **12**, 529–541 (2011).
- Klareskog, L., Catrina, A.I. & Paget, S. Rheumatoid arthritis. *Lancet* **373**, 659–672 (2009).
- Scott, D.L., Wolfe, F. & Huizinga, T.W. Rheumatoid arthritis. *Lancet* **376**, 1094–1108 (2010).
- Padyukov, L. *et al.* A genome-wide association study suggests contrasting associations in ACPA-positive versus ACPA-negative rheumatoid arthritis. *Ann. Rheum. Dis.* **70**, 259–265 (2011).
- Raychaudhuri, S. *et al.* Genetic variants at CD28, PRDM1 and CD2/CD58 are associated with rheumatoid arthritis risk. *Nat. Genet.* **41**, 1313–1318 (2009).
- Stahl, E.A. *et al.* Genome-wide association study meta-analysis identifies seven new rheumatoid arthritis risk loci. *Nat. Genet.* **42**, 508–514 (2010).
- Raychaudhuri, S. *et al.* Five amino acids in three HLA proteins explain most of the association between MHC and seropositive rheumatoid arthritis. *Nat. Genet.* **44**, 291–296 (2012).
- Klareskog, L., Ronnelid, J., Lundberg, K., Padyukov, L. & Alfredsson, L. Immunity to citrullinated proteins in rheumatoid arthritis. *Annu. Rev. Immunol.* **26**, 651–675 (2008).
- Padyukov, L., Silva, C., Stolt, P., Alfredsson, L. & Klareskog, L. A gene-environment interaction between smoking and shared epitope genes in HLA-DR provides a high risk of seropositive rheumatoid arthritis. *Arthritis Rheum.* **50**, 3085–3092 (2004).
- Mahdi, H. *et al.* Specific interaction between genotype, smoking and autoimmunity to citrullinated alpha-enolase in the etiology of rheumatoid arthritis. *Nat. Genet.* **41**, 1319–1324 (2009).
- Klareskog, L. *et al.* 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 (2006).
- Reinius, L.E. *et al.* Differential DNA methylation in purified human blood cells: implications for cell lineage and studies on disease susceptibility. *PLoS ONE* **7**, e41361 (2012).
- Houseman, E.A. *et al.* DNA methylation arrays as surrogate measures of cell mixture distribution. *BMC Bioinformatics* **13**, 86 (2012).
- Schadt, E.E. *et al.* An integrative genomics approach to infer causal associations between gene expression and disease. *Nat. Genet.* **37**, 710–717 (2005).
- Chen, L.S., Emmert-Streib, F. & Storey, J.D. Harnessing naturally randomized transcription to infer regulatory relationships among genes. *Genome Biol.* **8**, R219 (2007).
- MacKinnon, D.P. & MacKinnon, D. *Introduction to Statistical Mediation Analysis* (Routledge Academic, 2008).
- Millstein, J., Zhang, B., Zhu, J. & Schadt, E.E. Disentangling molecular relationships with a causal inference test. *BMC Genet.* **10**, 23 (2009).
- van der Laan, M.J., Dudoit, S. & Pollard, K.S. Multiple testing. Part II. Step-down procedures for control of the family-wise error rate. *Stat. Appl. Genet. Mol. Biol.* **3** Article14 (2004).
- Feinberg, A.P. & Irizarry, R.A. Evolution in health and medicine Sackler colloquium: Stochastic epigenetic variation as a driving force of development, evolutionary adaptation, and disease. *Proc. Natl. Acad. Sci. USA* **107** (suppl. 1), 1757–1764 (2010).
- Hayes, J.D. & Strange, R.C. Glutathione S-transferase polymorphisms and their biological consequences. *Pharmacology* **61**, 154–166 (2000).
- Strange, R.C., Jones, P.W. & Fryer, A.A. Glutathione S-transferase: genetics and role in toxicology. *Toxicol. Lett.* **112–113**, 357–363 (2000).
- Strange, R.C., Spiteri, M.A., Ramachandran, S. & Fryer, A.A. Glutathione-S-transferase family of enzymes. *Mutat. Res.* **482**, 21–26 (2001).
- Bohanec Grabar, P., Logar, D., Tomsic, M., Rozman, B. & Dolzan, V. Genetic polymorphisms of glutathione S-transferases and disease activity of rheumatoid arthritis. *Clin. Exp. Rheumatol.* **27**, 229–236 (2009).
- Yun, B.R., El-Sohemy, A., Cornelis, M.C. & Bae, S.C. Glutathione S-transferase M1, T1, and P1 genotypes and rheumatoid arthritis. *J. Rheumatol.* **32**, 992–997 (2005).
- Keenan, B.T. *et al.* Effect of interactions of glutathione S-transferase T1, M1, and P1 and HMOX1 gene promoter polymorphisms with heavy smoking on the risk of rheumatoid arthritis. *Arthritis Rheum.* **62**, 3196–3210 (2010).
- Lundstrom, E. *et al.* Effects of GSTM1 in rheumatoid arthritis; results from the Swedish EIRA study. *PLoS ONE* **6**, e17880 (2011).
- Novembre, J. *et al.* Genes mirror geography within Europe. *Nature* **456**, 98–101 (2008).
- Hao, K., Chudin, E., Greenawald, D. & Schadt, E.E. Magnitude of stratification in human populations and impacts on genome wide association studies. *PLoS ONE* **5**, e8695 (2010).
- Leek, J.T. *et al.* Tackling the widespread and critical impact of batch effects in high-throughput data. *Nat. Rev. Genet.* **11**, 733–739 (2010).
- Leek, J.T. & Storey, J.D. Capturing heterogeneity in gene expression studies by surrogate variable analysis. *PLoS Genet.* **3**, e161 (2007).
- Kang, E.Y., Ye, C., Shpitser, I. & Eskin, E. Detecting the presence and absence of causal relationships between expression of yeast genes with very few samples. *J. Comput. Biol.* **17**, 533–546 (2010).
- Jaffe, A.E. *et al.* Bump hunting to identify differentially methylated regions in epigenetic epidemiology studies. *Int. J. Epidemiol.* **41**, 200–209 (2012).
- Thurlings, R.M. *et al.* Monocyte scintigraphy in rheumatoid arthritis: the dynamics of monocyte migration in immune-mediated inflammatory disease. *PLoS ONE* **4**, e7865 (2009).
- Mueller, R.B. *et al.* MHC class II expression on myeloid cells inversely correlates with disease progression in early rheumatoid arthritis. *Rheumatology (Oxford)* **46**, 931–933 (2007).

ONLINE METHODS

Sample preparation. Recruitment of rheumatoid arthritis patients in the EIRA study was described previously¹¹. Briefly, only incident cases of rheumatoid arthritis were invited for the study within the years 1996–2009, from 18 clinics in Middle Sweden. Individuals were examined by rheumatologists and all patients correspond to ACR1987 criteria. The controls from the same population were matched by sex, age, smoking status and residence area. DNA was extracted from EDTA-treated blood and kept at -80°C until use. The cell purification protocol was described previously⁴².

Illumina genome-wide genotyping. The genotyping and quality control (QC) procedures have been described previously¹¹. Briefly, the EIRA samples were genotyped with the Illumina Human Hap300 v1.0 chip, Hap370CNVduo chip or Hap550duo chip. Samples included for analysis had call rates $>95\%$ and inferred gender consistent with clinical records. SNP filtering was done based on chip type, eliminating SNPs with call-rates below 95%, monomorphic SNPs, SNPs with a minor allele frequency < 0.005 , SNPs with a Hardy-Weinberg equilibrium $P < 1.0 \times 10^{-7}$ in controls, and SNPs mapping to multiple locations and nonautosomal chromosomal SNPs. This resulted in 306,994 autosomal SNPs in 1,966 samples in Hap300, 324,981 autosomal SNPs in Hap370CNVduo on 674 samples and 527,434 autosomal SNPs on 520 samples in Hap550duo passing the QC filters. Closely related individuals were identified by RELPAIR and PLINK. The member of each pair with the lower call rate was dropped from further analysis. To quantify and control for population stratification, we used a principal components approach implemented in the EIGENSTRAT software. EIGENSTRAT identified a total of 141 significant outliers, which were removed from further analysis. This resulted in a data set of 1,934 rheumatoid arthritis cases and 1,079 controls on 297,393 SNPs.

Genomic imputation. Imputation was done using the MACH algorithm based on HapMap 3. The cleaned EIRA GWAS data set (3,000 individuals) was used for imputation. The genotype calls are based on a QC cutoff of 0.9. Amino acids imputation within *HLA-DRB1*, *HLA-DPBI* and *HLA-B* was performed previously¹⁴.

Rheumatoid arthritis genetic risk genome wide association analysis. SNPs ($n = 1,196,263$) were tested for association with rheumatoid arthritis case-control status using an additive minor-allele dosage model in the cohort of 354 ACPA-rheumatoid arthritis cases and 335 population-matched controls selected for Illumina 450K methylation assay. No non-MHC SNPs were significantly associated with rheumatoid arthritis phenotype after adjusting for multiple testing using a Bonferroni-adjusted $\alpha = 0.05$ significance level.

Illumina 450K methylation assay. For each sample, 1 μg of genomic DNA was bisulfite-converted using an EZ DNA methylation Kit (ZYMO research) according to the manufacturer's recommendations. Converted genomic DNA was eluted in 22 μl of elution buffer. DNA methylation level was measured using the Illumina Infinium HD Methylation Assay (Illumina) according to the manufacturer's instructions. Briefly, 4 μl of bisulfite-converted DNA was isothermally amplified overnight (20–24 h) and fragmented enzymatically. Precipitated DNA was resuspended in hybridization buffer and dispensed onto the Infinium HumanMethylation450 BeadChips (12 samples/chip) using a Freedom EVO robot (Tecan). The hybridization procedure was performed at 48°C overnight (16–20 h) using an Illumina Hybridization oven. After hybridization, free DNA was washed away and the BeadChips were processed through a single nucleotide extension followed by immunohistochemistry staining using a Freedom EVO robot (Tecan). Finally, the BeadChips were imaged using an Illumina iScan.

Illumina 450K microarray data preprocessing. Detection P -values were calculated to identify failed probes as per Illumina's recommendations. No arrays exceeded our quality threshold of $>5\%$ failed probes. Probes on sex chromosomes or containing SNPs (dbSNP v132) in the probe sequence were excluded. Raw data were normalized using Illumina's control probe scaling procedure and converted to methylation values on the 0–1 scale ($M/(M + U + 100)$), where M and U represent the methylated and unmethylated signal intensities respectively).

Estimate differential cell counts. Differential cell counts for each individual were estimated using a published algorithm developed²⁰ with a slight modification. Briefly, the distribution of cell types for each sample was inferred based on DNA methylation signatures of the constituent cell types. A total of five different cell types, including T cells, NK cells, B cells, monocytes and granulocytes, were included in the estimation. DNA methylation signatures on sorted human cells from the Illumina 450K arrays¹⁹ were used as validation data. Among the 500 most informative CpG probes for distinguishing cell types chosen from the Illumina Infinium 27K array²⁰, all 473 probes also present in the Illumina 450K array are included in the analysis.

Identify rheumatoid arthritis-associated DMPs. To identify the DMPs associated with the rheumatoid arthritis phenotype, we fit a linear regression model predicting methylation at each CpG sites as a function of rheumatoid arthritis status, adjusted for age, sex, smoking status and estimated differential cell counts. Rheumatoid arthritis-DMP associations were corrected for multiple testing using a stringent Bonferroni-adjusted threshold of $0.05/(298,109 \text{ CpGs}) = 1.68 \times 10^{-7}$.

Identify genotype-dependent DMPs. All genome-wide significant (Bonferroni-adjusted $P < 0.05$) rheumatoid arthritis-associated DMPs were subsequently tested for association with genotype (1,196,263 SNPs) using an additive minor-allele dosage model. Genotype-DMP associations were corrected for multiple testing using a stringent Bonferroni-adjusted threshold of $0.05/(51,476 \text{ DMPs} \times 1,196,263 \text{ SNPs}) = 8.12 \times 10^{-13}$. SNPs associated with methylation variance were identified by fitting an additive minor-allele dosage model to absolute methylation residuals, calculated as the difference between a subject's methylation value and the genotype-specific mean. A Bonferroni-adjusted $\alpha = 0.05$ cutoff was used to determine significance.

CIT. Each of the genotype (G)-methylation (M)-phenotype (Y) relationships were assessed using the CIT²⁴ to classify them as methylation mediated, methylation consequential and independent. Note that the corresponding original terms describing the CIT are causal, reactive and independent²¹. Briefly, the CIT performs statistical tests for four conditions, all of which must be met for the methylation-mediated (causal) classification: (i) G and Y are associated, (ii) G is associated with M after adjusting for Y, (iii) M is associated with Y after adjusting for G and (iv) G is independent of Y after adjusting for M. Because all component conditions must be satisfied, the CIT P -value is defined using the intersection-union test framework⁴³ as the maximum of the four component test P -values. Because the CIT was designed for continuous phenotypes rather than case-control studies, we developed a modified version based on logistic regression and confirmed that all ten reported SNP-CpG pairs retained significant causal P -values.

Post hoc statistical analysis. Within each batch, we calculated the first principal component of cell type composition using a balanced sample of cases and controls. We then used a linear model to fit methylation values to this cell type composition proxy, calculated residuals representing batch and cell type-corrected estimates, and used these in epigenotype-rheumatoid arthritis association analyses.

DNA methylation analysis from flow cytometry-sorted cells. 30–50 mL of whole blood with heparin as a preservative was collected by a clinician from rheumatoid arthritis patients and controls with consent, and separated with Ficoll within 24 h of collection. Cells were sorted using AutoMACS (Miltenyi Biotec) into four populations: CD4⁺, CD8⁺, CD14⁺ and CD19⁺, and frozen as cell pellets in PBS at -80°C . Genomic DNA was extracted with salting-out method. DNA methylation level was measured using the Illumina 450K methylation array as described before.

Analysis software. All analysis was performed in R 2.14 and Bioconductor 2.9. Illumina 450K microarray data were analyzed with the *minfi* package.

42. Ronninger, M. *et al.* The balance of expression of PTPN22 splice forms is significantly different in rheumatoid arthritis patients compared with controls. *Genome Med.* **4**, 2 (2012).

43. Casella, G. & Berger, R.L. *Statistical Inference* (Duxbury Press, 2001).