PLOS | GENETICS

# Functional Variants in *NFKBIE* and *RTKN2* Involved in Activation of the NF-κB Pathway Are Associated with Rheumatoid Arthritis in Japanese

Keiko Myouzen[1], Yuta Kochi[1,2]*, Yukinori Okada[1,2,3], Chikashi Terao[4,5], Akari Suzuki[1], Katsunori Ikari[6], Tatsuhiko Tsunoda[7], Atsushi Takahashi[3], Michiaki Kubo[8], Atsuo Taniguchi[6], Fumihiko Matsuda[4,9,10], Koichiro Ohmura[5], Shigeki Momohara[6], Tsuneyo Mimori[5], Hisashi Yamanaka[6], Naoyuki Kamatani[11], Ryo Yamada[12], Yusuke Nakamura[13], Kazuhiko Yamamoto[1,2]

1 Laboratory for Autoimmune Diseases, Center for Genomic Medicine (CGM), RIKEN, Yokohama, Japan, 2 Department of Allergy and Rheumatology, Graduate School of Medicine, the University of Tokyo, Tokyo, Japan, 3 Laboratory for Statistical Analysis, CGM, RIKEN, Yokohama, Japan, 4 Center for Genomic Medicine, Kyoto University Graduate School of Medicine, Kyoto, Japan, 5 Department of Rheumatology and Clinical Immunology, Graduate School of Medicine, Kyoto University, Kyoto, Japan, 6 Institute of Rheumatology, Tokyo Women's Medical University, Tokyo, Japan, 7 Laboratory for Medical Informatics, CGM, RIKEN, Yokohama, Japan, 8 Laboratory for Genotyping Development, CGM, RIKEN, Yokohama, Japan, 9 CREST Program, Japan Science and Technology Agency, Kawaguchi, Saitama, Japan, 10 Institut National de la Santé et de la Recherche Médicale (INSERM), Unité U852, Kyoto University Graduate School of Medicine, Kyoto, Japan, 11 Laboratory for International Alliance, CGM, RIKEN, Yokohama, Japan, 12 Unit of Statistical Genetics, Center for Genomic Medicine, Graduate School of Medicine, Kyoto University, Kyoto, Japan, 13 Laboratory of Molecular Medicine, Human Genome Center, Institute of Medical Science, University of Tokyo, Tokyo, Japan

## Abstract

Rheumatoid arthritis is an autoimmune disease with a complex etiology, leading to inflammation of synovial tissue and joint destruction. Through a genome-wide association study (GWAS) and two replication studies in the Japanese population (7,907 cases and 35,362 controls), we identified two gene loci associated with rheumatoid arthritis susceptibility (*NFKBIE* at 6p21.1, rs2233434, odds ratio (OR) = 1.20, $P = 1.3 \times 10^{-15}$; *RTKN2* at 10q21.2, rs3125734, OR = 1.20, $P = 4.6 \times 10^{-9}$). In addition to two functional non-synonymous SNPs in *NFKBIE*, we identified candidate causal SNPs with regulatory potential in *NFKBIE* and *RTKN2* gene regions by integrating *in silico* analysis using public genome databases and subsequent *in vitro* analysis. Both of these genes are known to regulate the NF-κB pathway, and the risk alleles of the genes were implicated in the enhancement of NF-κB activity in our analyses. These results suggest that the NF-κB pathway plays a role in pathogenesis and would be a rational target for treatment of rheumatoid arthritis.

## Introduction

Rheumatoid arthritis (RA [MIM 180300]) is an autoimmune disease [1] with a complex etiology involving several genetic factors as well as environmental factors. Previous genome-wide association studies (GWAS) for RA have discovered many genetic loci [2–6], although the causal mechanisms linking the variants in these loci and disease etiology are largely unknown, except for in a few cases [6–8]. In contrast to mutations in Mendelian, monogenic diseases, most disease-associated variants in complex diseases, including autoimmune diseases, have moderate effects on disease susceptibility. This is because the disease causal variants in complex diseases are thought to have moderate effects on gene function, while amino acid changes introduced by the mutations of monogenic diseases have critical impacts on protein function [9]. Moreover, it has been demonstrated that the majority of autoimmune disease loci are expression quantitative trait loci (eQTLs) [10,11], indicating that accumulation of quantitative, but not qualitative, changes in gene function likely predisposes individuals to the disease. This renders it difficult to pinpoint the causal variants in the GWAS loci, especially in eQTLs, because all the variations in strong linkage disequilibrium (LD) with the marker SNP in a GWAS, the majority of which are not covered by the SNP array, are possible candidates for causal variants.

In recent years, with the emergence of next-generation sequencing technologies, the way we approach disease-causing variants has dramatically changed. First, a comprehensive map of human genetic variations is now available owing to the 1000 Genome Project [12], which allows us to grasp most of the potential common variants. This also enables us to perform genotype imputation of SNPs that are not directly genotyped in the GWAS, and consequently, to test them for association. Second, genomic studies using new technologies, such as chromatin immunoprecipitation-sequencing (ChIP-seq) and DNase I hypersensitive sites sequencing (DNase-seq), have advanced our understanding of how each genomic cluster regulates gene

## Author Summary

Rheumatoid arthritis (RA) is a chronic autoimmune disease affecting approximately 1% of the general adult population. More than 30 susceptibility loci for RA have been identified through genome-wide association studies (GWAS), but the disease-causal variants at most loci remain unknown. Here, we performed replication studies of the candidate loci of our previous GWAS using Japanese cohorts and identified variants in *NFKBIE* and *RTKN2* gene loci that were associated with RA. To search for causal variants in both gene regions, we first examined non-synonymous (ns)SNPs that alter amino-acid sequences. As *NFKBIE* and *RTKN2* are known to be involved in the NF-κB pathway, we evaluated the effects of nsSNPs on NF-κB activity. Next, we screened *in silico* variants that may regulate gene transcription using publicly available epigenetic databases and subsequently evaluated their regulatory potential using *in vitro* assays. As a result, we identified multiple candidate causal variants in *NFKBIE* (2 nsSNPs and 1 regulatory SNP) and *RTKN2* (2 regulatory SNPs), indicating that our integrated *in silico* and *in vitro* approach is useful for the identification of causal variants in the post–GWAS era.

transcription. If disease-associated variants are present in a critical site for gene regulation suggested by the ChIP-seq and DNase-seq studies, the disease-associated variants might possibly influence gene transcription levels such as through altered transcription factor-DNA binding avidity.

In the present study, we first performed replication studies of candidate loci in our previous GWAS and identified two association signals with genome-wide significance ($P<5\times10^{-8}$) in nuclear factor of kappa light polypeptide gene enhancer in B-cells inhibitor, epsilon (*NFKBIE* [MIM 604548]) and rhotekin 2 (*RTKN2*) loci. By utilizing publicly available datasets yielded by the above-mentioned genomic studies, we then performed integrated *in silico* and *in vitro* analysis to identify plausible causal variants in *NFKBIE* and *RTKN2* loci.

## Results

### Identification of rheumatoid arthritis susceptibility genes

We previously performed a GWAS of RA using a Japanese case-control cohort (2,303 cases and 3,380 controls) and identified significant associations in major histocompatibility complex, class II, DR beta 1 (*HLA-DRB1* [MIM 142857]), and chemokine (C-C motif) receptor 6 (*CCR6* [MIM 601835]) loci ($P_{GWAS}<5\times10^{-8}$) [6]. To reveal additional risk loci from those showing moderate associations in the GWAS (31 loci, $5\times10^{-8}<P_{GWAS}<5\times10^{-5}$), we selected a landmark SNP from each locus and genotyped it for an additional cohort (replication-1: 2,187 cases and 28,219 controls) (Table S1, S2). Among the 31 SNPs genotyped, seven SNPs were nominally associated with RA ($P<0.05$), which included SNPs in the tumor necrosis factor, alpha-induced protein 3 (*TNFAIP3* [MIM 191163]), and signal transducer and the activator of transcription 4 (*STAT4* [MIM 600558]) gene loci that were previously reported to be associated with RA [13,14] (Table S2). In a combined analysis of the GWAS and the 1st replication study, we identified two associations with genome-wide significance ($P<5\times10^{-8}$) in *NFKBIE* (6p21.1, rs2233434, $P=4.1\times10^{-11}$, odds ratio (OR) = 1.21, 95% confidence interval (CI) = 1.14–1.28) and in *RTKN2* (10q21.2, rs3125734, $P=3.7\times10^{-8}$, OR = 1.23, 95% CI = 1.14–1.32) (Table 1 and

Figure 1). *NFKBIE* was previously reported as a novel RA susceptibility gene locus in a meta-analysis of three GWASs for RA in the Japanese population, which included the GWAS set that the present study used [15]. *RTKN2* is located in the same region (10q21) as *ARID5B*, in which a significant association signal was also reported in the meta-analysis [15]. In our GWAS set, however, two significant signals were observed at rs3125734 (*RTKN2*: $P=4.8\times10^{-5}$) and rs10821944 (*ARID5B*: $P=7.4\times10^{-4}$), the former of which was tested as a landmark in the replication study. These two SNPs were in weak LD ($r^2=0.11$) and the independent effect of each SNP was observed after conditioning on each SNP (*RTKN2*: $P=1.5\times10^{-3}$, *ARID5B*: $P=0.024$, respectively). This indicated that two independent associations existed in this region, and the association of *RTKN2* is novel. We also confirmed the association in the *STAT4* locus [14] with genome-wide significance (2q32.2, rs10168266, $P=3.2\times10^{-8}$, OR = 1.16, 95% CI = 1.10–1.22) (Table S2). The associations in *NFKBIE* and *RTKN2* were further replicated in the 2nd replication cohort (3,417 cases and 3,763 controls; rs2233434, $P=1.1\times10^{-5}$, OR = 1.19, 95% CI = 1.10–1.30 and rs3125734, $P=0.016$, OR = 1.14, 95% CI = 1.02–1.26, respectively), confirming the associations in these loci (a combined analysis of three sets; rs2233434, $P=1.3\times10^{-15}$, OR = 1.20, 95% CI = 1.15–1.26 and rs3125734, $P=4.6\times10^{-9}$, OR = 1.20, 95% CI = 1.13–1.27, respectively) (Table 1 and Figure 1). We also genotyped these SNPs for individuals with systemic lupus erythematosus (SLE [MIM 152700]) ($n=657$) and Graves' disease (GD [MIM 275000]) ($n=1,783$). We identified a significant association of *RTKN2* (rs3125734) with GD ($P=3.4\times10^{-5}$, OR = 1.24, 95% CI = 1.12–1.37), whereas no significant associations were detected in *NFKBIE* (rs2233434) with either disease or in *RTKN2* (rs3125734) with SLE (Table S3).

### Functional analysis of non-synonymous SNPs

*NFKBIE* and *RTKN2* genes are both involved in the NF-κB pathway: *NFKBIE* encodes IκB epsilon (IκBε), a member of the IκB family [16], and its binding to NF-κB inhibits the nuclear translocation of NF-κB [17]; *RTKN2* encodes a member of Rho-GTPase effector proteins highly expressed in CD4$^+$ T cells [18] and is implicated in the activation of the NF-κB pathway [19]. Considering that the NF-κB pathway is critical for the pathogenesis of RA [20], these two genes could be strong candidates in these regions. To identify disease-causing variants, we first sequenced the coding regions of the genes using DNA from patients ($n=48$) to find variants that alter amino acid sequences. We identified four non-synonymous (ns)SNPs, which were all registered in the dbSNP database: two nsSNPs in *NFKBIE* (rs2233434 (Val194Ala) and rs2233433 (Pro175Leu)) and two in *RTKN2* (rs3125734 (Arg462His) and rs61850830 (Ala288Thr)), where rs2233434 and rs3125734 were the same as the landmark SNPs in the GWAS (Figure 1 and Figure 2A). The two nsSNPs of each locus were in strong LD (Figure 2B) and were both associated with disease (Table S4). In the haplotype analysis, a single common risk haplotype with a frequency $>0.05$ was observed in each locus, and significant associations with disease risk were detected (*NFKBIE*, $P=5.3\times10^{-8}$, Table S5; *RTKN2*, $P=5.7\times10^{-5}$, Table S6).

To investigate the effect of these nsSNPs on protein function, we evaluated them by *in silico* analysis using PolyPhen and SIFT software, which predicts possible impacts of amino acid substitutions on the structure and function of proteins, but all four nsSNPs were predicted to have little effect (Table S7), contrasting with the effect of Mendelian disease mutations [9]. We next examined their influence on the NF-κB activity in cells by performing NF-κB

58

**Table 1.** Association analysis of *NFKBIE* and *RTKN2* with rheumatoid arthritis.

| Gene | dbSNP ID | Allele (1/2) | Study set | Number of subjects Case | Control | Frequency of allele 1 Case | Control | Odds ratio (95% CI) | *P*-value[a] |
|------|----------|--------------|-----------|------|---------|------|---------|---------------------|----------|
| *NFKBIE* | rs2233434 | G/A | GWAS | 2,303 | 3,380 | 0.254 | 0.216 | 1.24 (1.13–1.35) | $2.2 \times 10^{-6}$ |
| | | | Replication study-1 | 2,186 | 28,204 | 0.245 | 0.215 | 1.19 (1.10–1.27) | $4.2 \times 10^{-6}$ |
| | | | Replication study-2 | 3,396 | 3,756 | 0.239 | 0.209 | 1.19 (1.10–1.30) | $1.1 \times 10^{-5}$ |
| | | | Combined analysis | 7,885 | 35,340 | 0.245 | 0.215 | 1.20 (1.15–1.26) | $1.3 \times 10^{-15}$ |
| *RTKN2* | rs3125734 | T/C | GWAS | 2,303 | 3,380 | 0.125 | 0.101 | 1.27 (1.13–1.43) | $4.8 \times 10^{-5}$ |
| | | | Replication study-1 | 2,185 | 28,218 | 0.129 | 0.110 | 1.20 (1.09–1.31) | $1.4 \times 10^{-4}$ |
| | | | Replication study-2 | 3,402 | 3,751 | 0.115 | 0.103 | 1.14 (1.02–1.26) | 0.016 |
| | | | Combined analysis | 7,890 | 35,349 | 0.122 | 0.108 | 1.20 (1.13–1.27) | $4.6 \times 10^{-9}$ |

[a]: Cochran-Armitage trend test was used for the GWAS and replication studies. Mantel-Haenszel method was used for the combined analysis.
doi:10.1371/journal.pgen.1002949.t001

reporter assays with haplotype-specific expression vectors. In *NFKBIE*, the non-risk haplotype (A-C: rs2233434 (non-risk allele (NR))-rs2233433 (NR)) displayed an inhibitory effect on NF-κB activity compared with the mock construct, which reflected compulsorily binding of exogenous IκBε to the endogenous NF-κB, as shown in a previous study [16]. Of note, the risk haplotype (G-T: risk allele (R)-R) showed higher NF-κB activity than A-C (NR-NR) (Figure 3A), suggesting impaired inhibitory potential of G-T (R-R) products. No haplotypic difference was detected in the protein expression levels of these constructs (Figure 3C). We also examined two additional constructs of G-C (R-NR) and A-T (NR-R) haplotypes to evaluate the effect of each nsSNP (Figure S1A, S1B). Because NF-κB activity increased in the order of A-C<G-C<A-T<G-T (rs2233434-rs2233433: NR-NR<R-NR<NR-R<R-R) when cells were stimulated with TNF-α, the C>T substitution (Pro175Leu) in rs2233433 may have more impact on the protein function of IκBε compared with the A>G substitution (Val194Ala) in rs2233434. In contrast to the observations in *NFKBIE*, no clear difference was detected between the two common haplotype products of *RTKN2* in either their effect on NF-κB activity or protein expression levels, although both products enhanced NF-κB activity as reported previously (Figure 3B, 3D) [19]. These functional analyses of nsSNPs suggest that two nsSNPs (rs2233434 and rs2233433) in the *NFKBIE* region are candidates for causal SNPs.

## ASTQ analysis suggested the existence of regulatory variants

As the majority of autoimmune disease loci have been implicated as eQTL [11], we speculated that variants in the *NFKBIE* and *RTKN2* loci would influence gene function by regulating gene expression, in addition to changing the amino acid sequences. To address this possibility, we performed allele-specific transcript quantification (ASTQ) analysis by using allele-specific probes targeting the nsSNPs in exons (rs2233434 for *NFKBIE* and rs3125734 for *RTKN2*, both of which were the GWAS landmarks). The genomic DNAs and cDNAs were extracted from peripheral blood mononuclear cells (PBMCs) in individuals with heterozygous genotype (*n* = 14 for *NFKBIE* and *n* = 6 for *RTKN2*) and from lymphoblastoid B-cell lines (*n* = 9) for *NFKBIE*. As the expression levels of *RTKN2* were low in lymphoblastoid B cells, only PBMCs were used. When quantified by allele-specific probes, transcripts from the risk allele of *NFKBIE* showed 1.1-fold and 1.2-fold lower amounts (in PBMCs and lymphoblastoid B cells, respectively) than

those from non-risk alleles (*P* = 0.012 and $5.3 \times 10^{-4}$, respectively; Figure 3E and Figure S2). In contrast, 1.5-fold higher amounts of transcripts were observed in the risk allele of *RTKN2* (*P* = 0.016; Figure 3F). These allelic imbalances suggested that both gene loci were eQTL and that there existed variants with *cis*-regulatory effects. Moreover, considering the inhibitory effects of *NFKBIE* and the activating potential of *RTKN2* on NF-κB activity, which might both be dose dependent (Figure 3G, 3H), these regulatory variants in the risk alleles should enhance NF-κB activity *in vivo*.

## Integrated *in silico* and *in vitro* analysis to search for regulatory variants

To comprehensively search the two genomic regions for causal regulatory variants, we performed an integrated *in silico* and *in vitro* analysis with multiple steps (Figure 4 and Figures S3, S4). We first determined the target genomic region by selecting LD blocks containing disease-associated SNPs ($P_{GWAS} < 1.0 \times 10^{-3}$) (Step 1). We then extracted SNPs with frequencies of >0.05 from HapMap and 1000 Genome Project databases in the region (Step 2). We excluded uncommon variants (MAF<0.05) from the analysis because of their low imputation accuracy in the GWAS (93% of uncommon variants in *NFKBIE* and 76% in *RTKN2* exhibited Rsq <0.6). There is neither structural variation (>1 kbps) nor indels (100 bps to 1 kbs) that are common in the population (frequency >0.01) in these loci. To evaluate the *cis*-regulatory potential of sequences around the SNPs *in silico*, we used the regulatory potential (RP) score [21,22]. This score was calculated based on the extent of sequence conservation among species or similarity with known regulatory motifs. We selected SNPs from the genomic elements with an RP score >0.1 (Step 3a). Subsequently, we selected SNPs from sites of transcriptional regulation as demonstrated by previous ChIP-seq studies (transcription factor binding sites [23,24] and histone modification sites [25,26]) or a DNase-seq study (DNase I hypersensitivity sites) [27] (Step 3b). Finally, these SNPs with regulatory potential were further screened out by the disease-association status (*P*<0.05) using an imputed GWAS dataset (Step 4). As a consequence, we selected 14 SNPs in *NFKBIE* and 10 SNPs in *RTKN2* that had regulatory potential predicted *in silico*.

To further investigate the regulatory potential of the SNPs, we evaluated 31-bp sequences around the SNPs by *in vitro* assays. First, we examined their ability to bind nuclear proteins by EMSAs (Step 5a) using nuclear extracts from lymphoblastoid B cells (PSC cells) and Jurkat cells. Of the 24 SNPs examined, nine
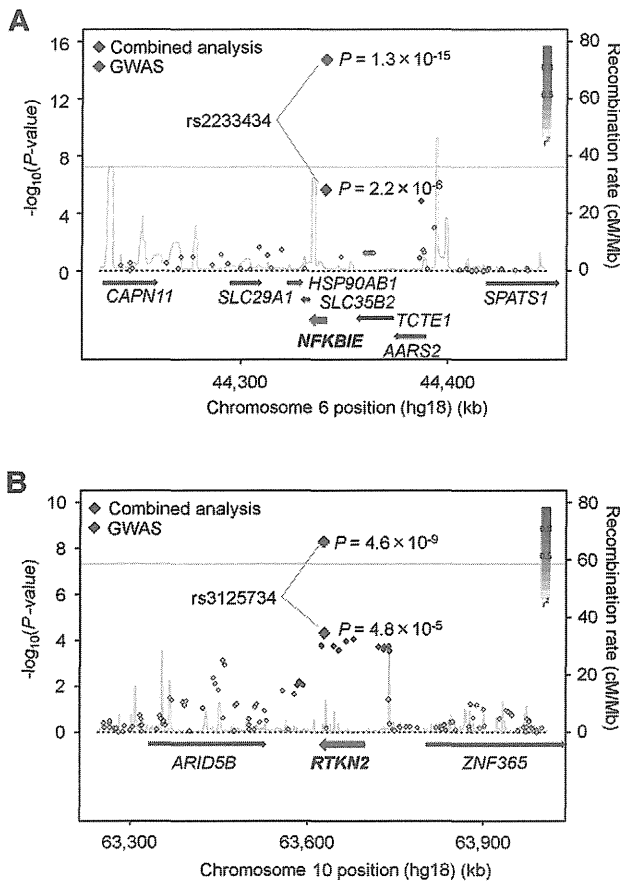
59

**A**



**B**



**Figure 1. Association plots of *NFKBIE* and *RTKN2* regions.** The diamonds represent the $-\log_{10}$ of the Cochran-Armitage trend *P*-values. Large diamonds show landmark SNPs in *NFKBIE* (rs2233434: A) and *RTKN2* (rs3125734: B). Red: GWAS, Blue: combined analysis. Red colors of each SNP indicate its $r^2$ with landmark SNP. Gray lines indicate the genome-wide significance threshold ($P < 5 \times 10^{-8}$). For each plot, the $-\log_{10}$ *P*-values (y-axis) of the SNPs are presented according to their chromosomal positions (x-axis). Physical positions are based on NCBI build 36.3 of the human genome. Genetic recombination rates, estimated using the 1000 Genome Project (JPT and CHB), are represented by the blue line.
doi:10.1371/journal.pgen.1002949.g001

SNPs displayed allelic differences, implying differential potential of transcriptional activity between these alleles (Figure 5A and Figure S5). We then evaluated the enhancing or repressing activity of the sequences by luciferase reporter assays (Step 5b). We cloned them into the pGL4.24 vector, which has minimal promoter activity, and transfected these constructs into HEK293A cells (for *NFKBIE* and *RTKN2*), lymphoblastoid B cells (for *NFKBIE*), and Jurkat cells (for *RTKN2*). Among the three SNPs examined in *NFKBIE*, the risk allele of rs2233424 (located −396 bps from the 5′ end) displayed stronger repression activity (Figure 2A and Figure 5B) than that of the non-risk allele. Among the six SNPs in *RTKN2*, the risk alleles of rs12248974 (approximately 10 kb from the 3′ end) and rs61852964 (−215 bps from the 5′ end) showed higher enhancing activity compared with the non-risk alleles (Figure 2A and Figure 5B). These results corresponded to the results of ASTQ analyses (Figure 3E, 3F). Other SNPs showed no allelic differences or had the opposite trend of transcriptional activity in the risk allele compared to the results of ASTQ analysis (Figure S6).

To confirm the regulatory potential of these SNPs, we investigated the correlation between genotypes and gene expression levels in lymphocytes utilizing the data from the previous eQTL studies. We evaluated the expression of *RTKN2* in primary T cells from Western European individuals by using Genevar software [28,29]. Though *NFKBIE* is also expressed in primary T cells, the genotypes of rs2233424 are not available. We thus evaluated gene expression data of lymphoblastoid B-cell lines obtained from HapMap individuals (Japanese (JPT) + Han Chinese in Beijing (CHB), European (CEU), and African (YRI)) [30,31] instead. The *NFKBIE* expression level decreased with the number of risk alleles of rs2233424 ($R = -0.18$, $P = 0.020$), and the *RTKN2* expression levels increased with that of rs1432411 (a proxy for rs12248974, $r^2 = 0.97$) ($R = 0.27$, $P = 0.018$) (Figure 5C), corresponding to the results of the *in vitro* assays. The data for rs61852964 in *RTKN2* was not available. Among the SNPs that displayed opposite transcriptional activities in the reporter assays compared to the results of ASTQ, the data for rs2233434, rs77986492, and rs3852694 (a proxy for rs1864836, $r^2 = 1.0$) were available (Figure S7 and S8). These SNPs displayed the opposite direction of the correlation trend as compared to the results of reporter assays, but parallel to ASTQ, implying that the regulatory effects observed in the *in vitro* assays were cancelled out by the effects of other regulatory variants on the same haplotype *in vivo*.

Finally, we validated the associations of these regulatory (r)SNPs observed in the imputed GWAS dataset. We directly genotyped them by TaqMan assay and confirmed significant associations (Table S8). As the candidate causal variants (nsSNPs and rSNPs) and the landmark SNPs of GWAS were in strong LD at each locus (Figure 2A, 2B), we evaluated the independent effect of each SNP by haplotype analysis in both loci (Table S9 and S10) and the conditional logistic regression analysis in *RTKN2* (Table S11). The conditional analysis was not performed in *NFKBIE* because three candidate causal variants were in strong LD ($r^2 > 0.9$). However, the analyses for these two loci did not demonstrate any evidence of primary or independent effects across the candidate causal variants, and it remains a possibility that all of the functional variants were involved in the pathogenesis. In addition, although the landmark nsSNP (rs3125734) in *RTKN2* did not display any influence on NF-κB activity in our *in vitro* assays, rs3125734 might influence functions of *RTKN2* other than those in the NF-κB pathway; alternatively, it is still possible that rs3125734 tags the effects of other unknown variants, such as rare variants, in addition to the other two rSNPs (rs12248974 and rs61852964).

## Discussion

In the present study, we performed a replication study of our previously reported GWAS and identified variants in *NFKBIE* and *RTKN2* loci that were associated with RA susceptibility. The associations of *NFKBIE* and *RTKN2* loci have not been reported in other populations with genome-wide significance. However, rs2233434 in *NFKBIE* showed a suggestive association (589 cases vs. 1,472 controls, $P = 0.0099$, OR = 1.57, 95% CI = 1.11–2.21) in a previous meta-analysis in European populations [32]. The weak association signal in Europeans may be partially due to the lower frequency of the risk allele (0.04 in Europeans compared to 0.22 in Japanese). On the other hand, the association of rs3125734 in *RTKN2* was not observed in a GWAS meta-analysis of European populations (cases 5,539 vs. controls 20,169, $P = 0.11$, OR = 1.04, 95% CI = 0.99–1.09). As the association of *RTKN2* locus was also implicated in Graves' disease in a Han Chinese population [33], the association in *RTKN2* locus may be unique to Asian populations.

To find the disease causal variants in disease-associated loci, target re-sequencing and variant genotyping with a large sample set followed by conditional association analysis examining the
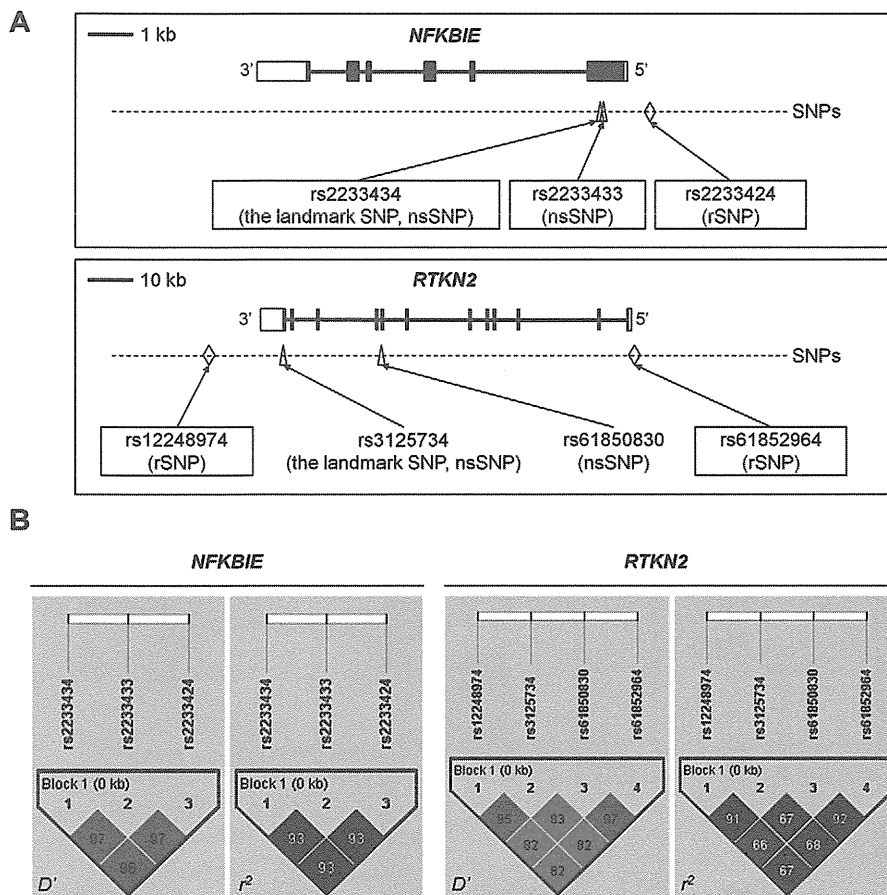
60

**Figure 2. Genomic position and LD blocks.** (A) Genomic position of non-synonymous (ns)SNPs and regulatory (r)SNPs in *NFKBIE* and *RTKN2*. *NFKBIE* (top) and *RTKN2* (bottom) correspond to transcripts NM_004556.2 and NM_145307.2, respectively. Exons are shown as boxes, where black boxes represent coding regions and open boxes represent untranslated regions. Intron sequences are drawn as lines. Open triangles represents nsSNPs and open diamond shapes indicate candidate rSNPs. dbSNP IDs of candidate causal variants were boxed in a solid line. (B) LD patterns for nsSNPs and candidate rSNPs in *NFKBIE* (left) and *RTKN2* (right) gene regions. LD blocks were constructed from genotype data of 3,290 control individuals of the GWAS. The diagrams show pairwise LD values as quantified using the $D'$ and $r^2$ values.
doi:10.1371/journal.pgen.1002949.g002

independent effects of each variant would be the first step. For this purpose, a recent attempt to fine-map the known autoimmunity risk loci in Celiac disease (MIM 212750) using an "Immunochip" brought us several insights [34]. First, no stronger signals compared to the GWAS signals were detected in most of the known loci, while additional independent signals were found in several loci. Second, none of the genome-wide significant common SNP signals could be explained by any rare highly penetrant variants. Third, although the fine-mapping strategy could localize the association signals into finer scale regions, it could not identify the actual causal variants due to strong LD among the variants, indicating that an additional approach, such as functional evaluation of candidate variants, is needed.

In the present study, we focused on common variants to find causal variants. Instead of re-sequencing additional samples, we utilized the 1000 Genome Project dataset, where the theoretically estimated cover rate for common variants (frequency of >0.05) in our population is >0.99 [12,35]. To fine-map the association signals, we performed imputation-based association analysis, where we could not find any association signals that statistically exceeded the effect of landmark SNPs (rs2233434 for *NFKBIE* and rs3125734 for *RTKN2*) in both gene regions (Figures S3 and S4).

We also performed a conditional logistic regression analysis, and found no additional independent signals of association when conditioned on each landmark SNP (data not shown). Although the imputation-based association tests may yield some bias compared to direct genotyping of the variants, these results suggested that variants in strong LD with the landmark SNPs were strong candidates for causal variants.

Following the analysis of nsSNPs, we evaluated cis-regulatory effects of variants in the two regions by ASTQ analysis using both B-cell lines and primary cells (PBMC), the majority of which consisted of T and B lymphocytes. As the mechanism of gene-regulation is substantially different between cell types [26], ASTQ analysis in more specific cell types that are relevant to the disease etiology, such as Th1 and Th17 cells, would be ideal to evaluate the cis-regulatory effects of variants. In this context, a more comprehensive catalog of the eQTL database of multiple cell types should be established for genetic study of diseases. As our ASTQ analysis demonstrated cis-regulatory effects of variants in both regions, we then performed an integrated *in silico* and *in vitro* analysis to identify candidate regulatory variants. Accumulating evidence by recent ChIP-seq and DNase-seq studies suggested that cis-regulatory variants are
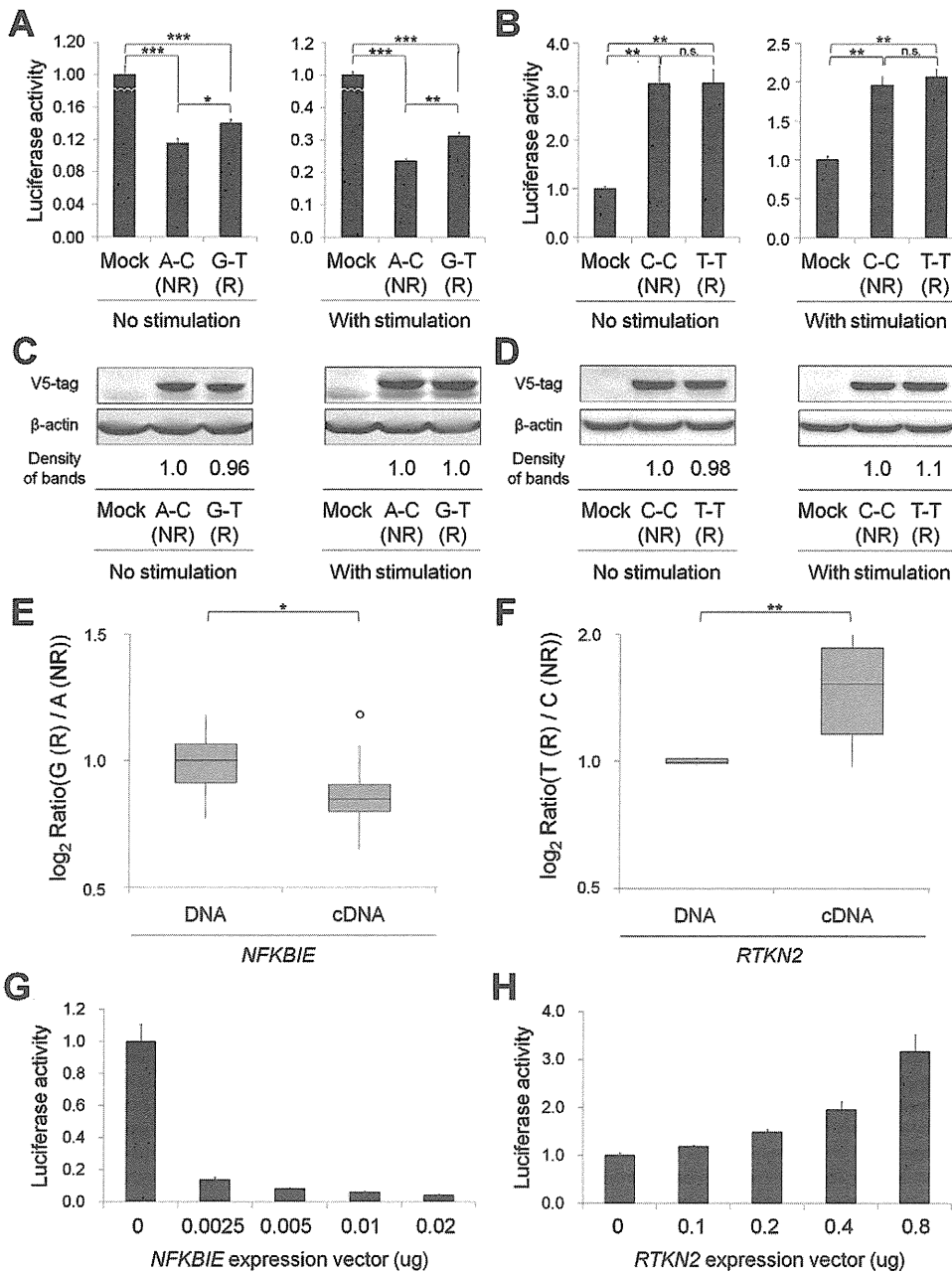
**Figure 3. Functional evaluation of nsSNPs and allelic imbalance of expression in *NFKBIE* and *RTKN2*.** (A, B) Effects of nsSNPs in *NFKBIE* (A) and *RTKN2* (B) on NF-κB activity by luciferase assays. Two haplotype constructs (A-C (rs2233434-rs2233433; non-risk (NR)) and G-T (risk (R)) for *NFKBIE* and C-C (rs3125734-rs61850830; NR) and T-T (R) for *RTKN2*) were used. The expression vector of each construct, pGL4.32[*luc2P*/NF-κB-RE] vector and pRL-TK vector were transfected into HEK293A cells. Data represent the mean ± s.d. Each experiment was performed in sextuplicate, and experiments were independently repeated three times. *$P<0.05$, **$P<1.0\times10^{-5}$, and ***$P<1.0\times10^{-10}$ by Student's *t*-test. n.s.: not significant. (C, D) Protein expression levels of each haplotype construct. Anti-V5 tag antibody was used in the Western blotting analysis to detect the expression of exogenous IκBε (C) and RTKN2 (D). Beta-actin expression was used as an internal control. The densities of the bands were quantified and normalized to that of the risk allele. (E, F) Allelic imbalance of expression in *NFKBIE* (E) and *RTKN2* (F). ASTQ was performed using samples from individuals heterozygous for rs2233434 (G/A) in *NFKBIE* and rs3125734 (T/C) in *RTKN2*. Genomic DNAs and cDNAs were extracted from PBMCs ($n=14$ for *NFKBIE* and $n=6$ for *RTKN2*). The y-axis shows the $\log_2$ ratio of the transcript amounts in target SNPs (risk allele/non-risk allele). The top bar of the box-plot represents the maximum value and the lower bar represents the minimum value. The top of box is the third quartile, the bottom of box is the first quartile, and the middle bar is the median value. The circle is an outlier. *$P=0.012$, **$P=0.016$, by Student's *t*-test. (G, H) Dose-dependent inhibition of *NFKBIE* (G) and activation of *RTKN2* (H) on NF-κB activity. Various doses of expression vectors carrying the non-risk allele of each gene were transfected into HEK293A cells with pGL4.32 and pRL-TK vectors.
doi:10.1371/journal.pgen.1002949.g003

located in the key regions of transcriptional regulation [26,36], warranting the prioritization of variants before evaluation by *in vitro* assays. This could also minimize false-positive results of the

*in vitro* assays. However, there may be additional causal variants, including rare variants, unsuccessfully selected at each step of our integrated screening. Therefore, the screening strategy
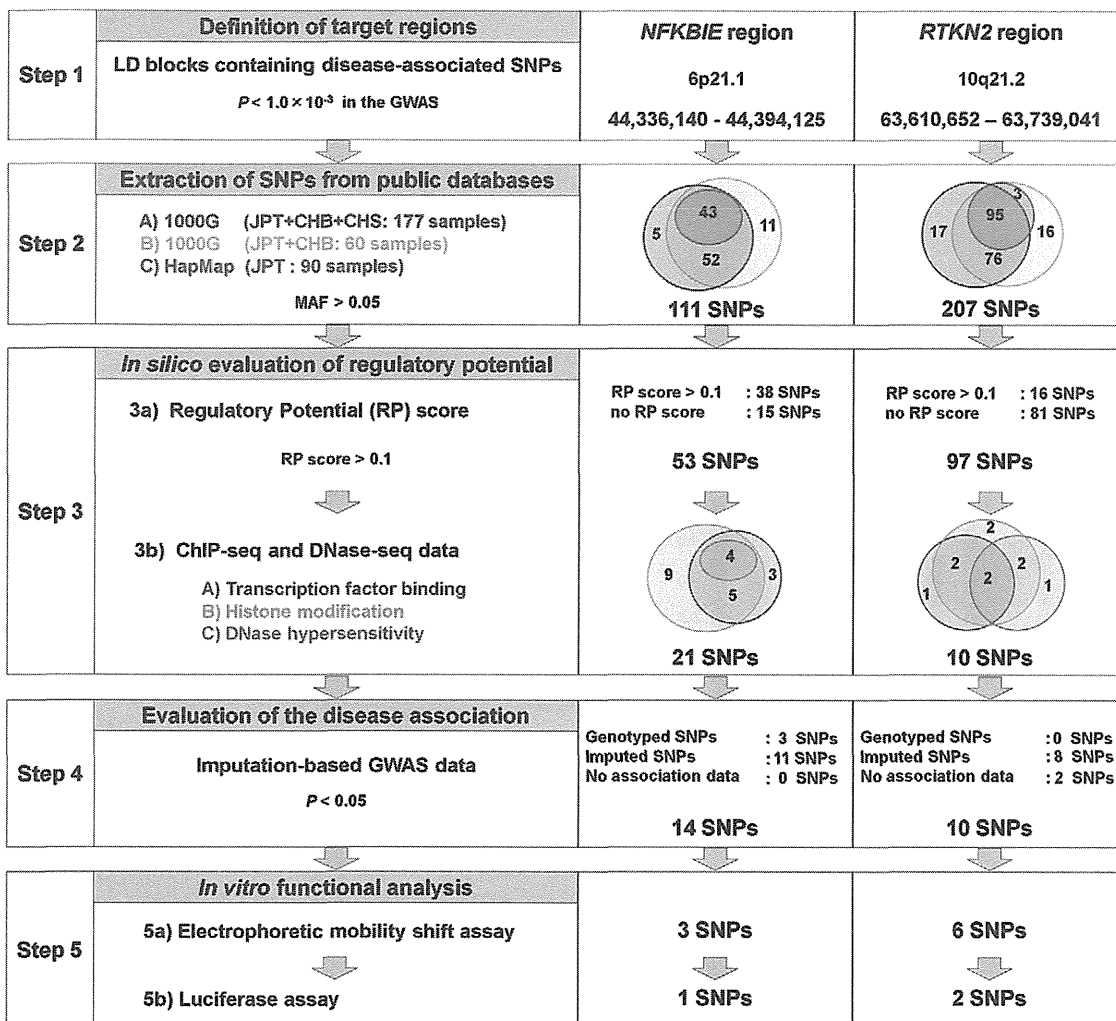
| | Definition of target regions | *NFKBIE* region | *RTKN2* region |
|---|---|---|---|
| Step 1 | LD blocks containing disease-associated SNPs<br>$P < 1.0 \times 10^{-3}$ in the GWAS | 6p21.1<br>44,336,140 - 44,394,125 | 10q21.2<br>63,610,652 - 63,739,041 |

| | Extraction of SNPs from public databases | *NFKBIE* region | *RTKN2* region |
|---|---|---|---|
| Step 2 | A) 1000G (JPT+CHB+CHS: 177 samples)<br>B) 1000G (JPT+CHB: 60 samples)<br>C) HapMap (JPT : 90 samples)<br>MAF > 0.05 | 43 / 5 / 11 / 52<br>**111 SNPs** | 95 / 17 / 16 / 76 / 3<br>**207 SNPs** |

| | In silico evaluation of regulatory potential | | |
|---|---|---|---|
| Step 3 | 3a) Regulatory Potential (RP) score<br>RP score > 0.1 | RP score > 0.1 : 38 SNPs<br>no RP score : 15 SNPs<br>**53 SNPs** | RP score > 0.1 : 16 SNPs<br>no RP score : 81 SNPs<br>**97 SNPs** |
| | 3b) ChIP-seq and DNase-seq data<br>A) Transcription factor binding<br>B) Histone modification<br>C) DNase hypersensitivity | 9 / 4 / 3 / 5<br>**21 SNPs** | 2 / 2 / 2 / 2 / 1 / 1<br>**10 SNPs** |

| | Evaluation of the disease association | | |
|---|---|---|---|
| Step 4 | Imputation-based GWAS data<br>P < 0.05 | Genotyped SNPs : 3 SNPs<br>Imputed SNPs : 11 SNPs<br>No association data : 0 SNPs<br>**14 SNPs** | Genotyped SNPs : 0 SNPs<br>Imputed SNPs : 8 SNPs<br>No association data : 2 SNPs<br>**10 SNPs** |

| | In vitro functional analysis | | |
|---|---|---|---|
| Step 5 | 5a) Electrophoretic mobility shift assay | **3 SNPs** | **6 SNPs** |
| | 5b) Luciferase assay | **1 SNPs** | **2 SNPs** |

**Figure 4. Overview of SNP selection using integrated *in silico* and *in vitro* approaches.** The figure shows the SNP selection process (left) and the results of *NFKBIE* (middle) and *RTKN2* (right). (Step 1) LD blocks that contain disease-associated SNPs ($P_{GWAS} < 1.0 \times 10^{-3}$) were selected. (Step 2) SNPs were extracted from three databases (A–C). 1000G, 1000 Genome Project; HapMap, International HapMap Project. A) JPT, CHB, and CHS samples ($n = 177$) from the 1000G (the August 2010 release). B) JPT and CHB samples ($n = 60$) from the pilot 1 low coverage study data of 1000G (the March 2010 release). C) JPT samples ($n = 90$) from HapMap phase II+III (release #27). SNPs with minor allele frequency >0.05 were selected. (Step 3) Prediction of regulatory potential *in silico*. 3a) Regulatory potential (RP) scores were used for SNP selection, where an RP score >0.1 indicated the presence of regulatory elements. SNPs without RP scores were also selected. 3b) Prediction of regulatory elements by ChIP-seq data and DNase-seq data. (A) Transcription factor binding sites, (B) histone modification sites (CTCF binding, H3K4me1, H3K4me2, H3K4me3, H3K27ac, H3K9ac), and (C) DNase I hypersensitivity sites were evaluated. ChIP-seq and DNase-seq data derived from GM12878 EBV-transformed B cells were used for *NFKBIE* and *RTKN2*. DNase-seq data of Th1, Th2, and Jurkat cells were also used for *RTKN2*. (Step 4) Association data of the imputation-based GWAS using 1000G reference genotypes were used. SNPs with a significance level of $P < 0.05$ were selected. SNPs without association data were also selected. (Step 5) EMSAs and luciferase assays were performed for evaluation of regulatory potentials *in vitro*.
doi:10.1371/journal.pgen.1002949.g004

should be refined as the quality and quantity of genomic databases improves in the future.

We identified multiple candidate causal variants in *NFKBIE* (two nsSNPs and one rSNP) and *RTKN2* (two rSNPs). We could not statistically distinguish the primary effect of each candidate causal variant, because these variants are in strong LD and on the same common haplotype. However, multiple causal variants could be involved in a single locus, which is also seen in another well-known autoimmune locus in 6q23 (*TNFAIP3* gene locus), where both an nsSNP and a regulatory variant have been shown to be functionally related to the disease [8,37]. The risk haplotype of nsSNPs in *NFKBIE* (rs2233433 and rs2233434) showed an enhancement of NF-κB activity, which might reflect an impaired

inhibitory effect of IκB-ε on nuclear translocation of NF-κB. On the other hand, down-regulated *NFKBIE* expression and up-regulated *RTKN2* expression were observed at the risk haplotypes, which may be regulated in *cis* by the rSNPs (rs2233424 in *NFKBIE*, rs12248974 and rs61852964 in *RTKN2*). As overexpression studies have also demonstrated dose-dependent attenuation of NF-κB activity by *NFKBIE*, and dose-dependent enhancement by *RTKN2*, the *cis*-regulatory effects of these rSNPs should enhance the NF-κB activity in the risk allele. Taken together with the effect of nsSNPs in *NFKBIE*, the enhancement of NF-κB activity may play a role in the pathogenesis of the disease. This is further supported by evidence that previous GWAS for RA have also identified genes related to the NF-κB pathway, such as *TNFAIP3* [13], v-rel
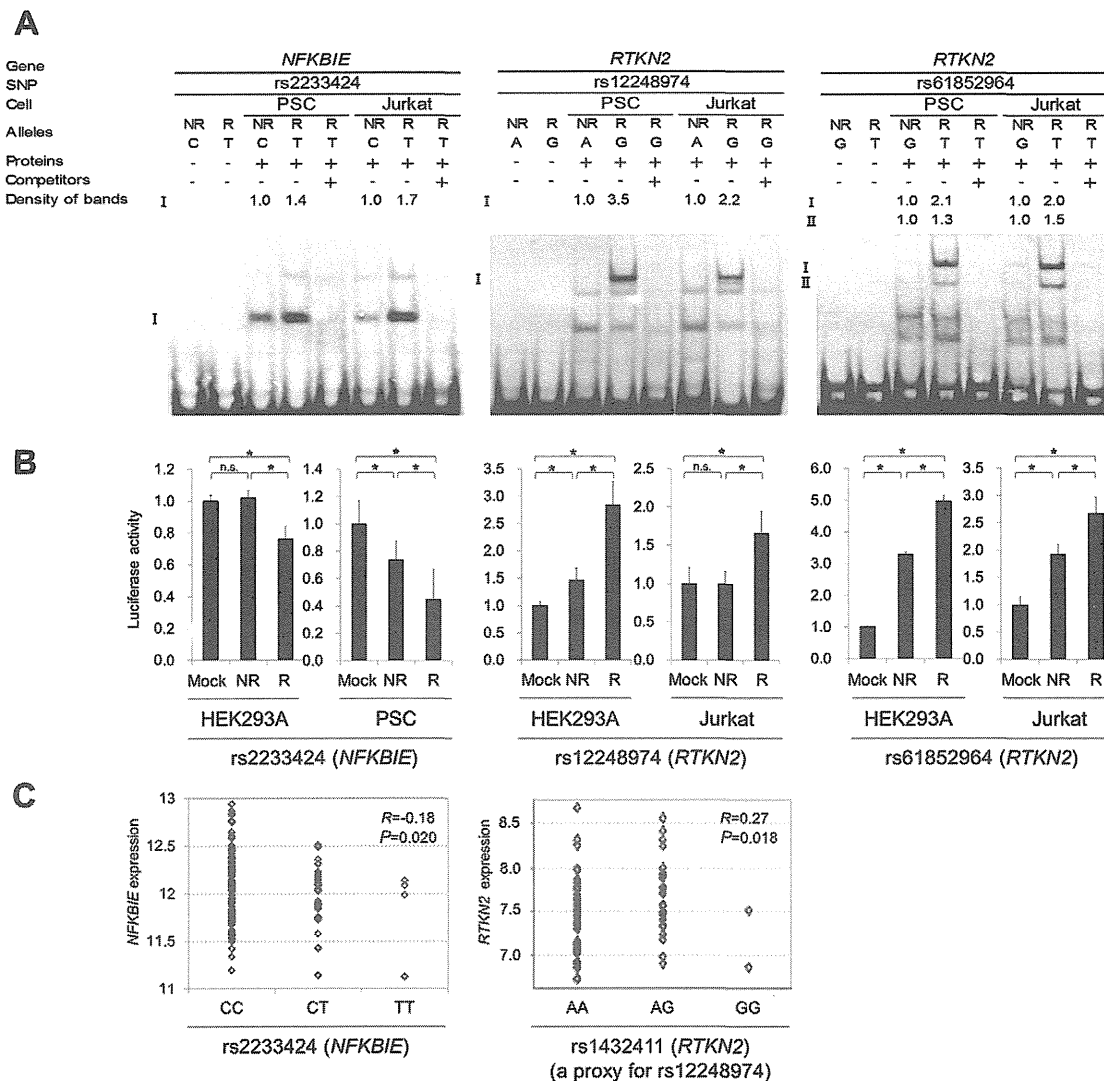
63

**Figure 5. Evaluation of candidate regulatory SNPs *in vitro*.** (A) Binding of nuclear factors from lymphoblastoid B-cells (PSC cells) and Jurkat cells to the 31-bp sequences around each SNP was evaluated by EMSA. Unlabeled probes in 200-fold excess as compared to the labeled probes were used for the competition experiment. The densities of the bands were quantified and normalized to that of the risk allele. rs2233424 in *NFKBIE* (C(NR)/ T(R)) (left), rs12248974 (A(NR)/G(R)) (middle) and rs61852964 (G(NR)/T(R)) (right) in *RTKN2*. (B) Transcriptional activities were evaluated by luciferase assays. Each 31-bp oligonucleotide was inserted into the pGL4.24[*Luc2P*/minP] vector. Luc, luciferase; minP, minimal promoter. Transcfection was performed with HEK293A (for all the SNPs), PSC cells (for rs2233424), and Jurkat cells (for rs12248974 and rs61852964). rs2233424 (left), rs12248974 (middle), and rs61852964 (right). Data represent the mean ± s.d. Each experiment was performed in sextuplicate and independently repeated three times. *$P<0.05$ by Student's *t*-test. n.s.: not significant. (C) Liner regression analysis of the relationship between SNP genotype and gene expression level. *NFKBIE* expression data in lymphoblastoid B-cell lines of HapMap individuals (JPT+CHB, CEU and YRI; $n=151$), and *RTKN2* expression data in primary T cells from umbilical cords of Western European individuals ($n=85$) were used. The x-axis shows the SNP genotypes and the y-axis represents the log₂-transformed gene expression level. R: correlation coefficient between SNP genotype and gene expression. Rs2233424 genotypes and *NFKBIE* expression level (left). The genotype classification by population: JPT+CHB, CC = 52, CT = 1; CEU, CC = 35, CT = 2; YRI, CC = 32, CT = 2, TT = 4. Rs1432411 genotypes and *RTKN2* expression level (right). Rs1432411 was used as a proxy SNP of rs12248974 ($r^2=0.97$).
doi:10.1371/journal.pgen.1002949.g005

reticuloendotheliosis viral oncogene homolog (*REL* [MIM 164910]) [5], TNF receptor-associated factor 1 (*TRAF1* [MIM 601711]) [3], and CD40 molecule TNF receptor superfamily member 5 (*CD40* [MIM 109535]) [38].

In conclusion, we identified *NFKBIE* and *RTKN2* as genetic risk factors for RA. Considering the allelic effect of both genes, enhanced NF-κB activity may play a role in the pathogenesis of the disease. Because NF-κB regulates the expression of numerous genes, including inflammatory and immune response mediators, NF-κB and its regulators identified by GWAS are promising targets for the treatment of RA.

## Materials and Methods

### Ethics statement

All subjects were of Japanese origin and provided written informed consent for participation in the study, which was approved by the ethical committees of the institutional review boards.

### Subjects

A total of 7,907 RA cases, 657 SLE cases, 1,783 GD cases, and 35,362 control subjects were enrolled in the study through medical

institutes in Japan under the support of the BioBank Japan Project, Center for Genomic Medicine at RIKEN, the University of Tokyo, Tokyo Women's Medical University, and Kyoto University. The same case and control samples were used in the previous meta-analysis of GWASs in the Japanese population (Table S1) [15] . RA and SLE subjects met the revised American College of Rheumatology (ACR) criteria for RA [39]. Diagnosis of individuals with GD was established on the basis of clinical findings and results of the routine examinations for circulating thyroid hormone and thyroid-stimulating hormone concentrations, thyroid-stimulating hormone receptors, ultrasonography, $^{[99m]}TCO_4^-$ (or $[^{123}I]$) uptake, and thyroid scintigraphy. DNAs were extracted from peripheral blood cells using a standard protocol. Total RNAs were also extracted from PBMCs of healthy individuals ($n = 20$) using an RNeasy kit (QIAGEN, Valencia, CA, USA). Details of the samples are summarized in Table S1.

## Genotyping and quality control

In the GWAS, RA cases and controls were genotyped using Illumina Human610-Quad and Illumina Human 550v3 Genotyping BeadsChips (Illumina, San Diego, CA, USA), respectively, and quality control of genotyping was performed as described previously [6]. For replication study of candidate loci, a landmark SNP was selected from each locus that satisfied $5 \times 10^{-8} < P_{GWAS} < 5 \times 10^{-5}$ in the GWAS. If multiple candidate SNPs existed within $\pm 100$ kb, the SNP with the lowest $P$-value was selected. All case subjects in the replication study and both case and control subjects in the validation study of candidate causal variants were genotyped using TaqMan SNP genotyping assays (Table S12) (Applied Biosystems, Foster City, CA, USA) with an ABI Prism 7900HT Sequence Detection System (Applied Biosystems). Because of the availability of DNA samples, only a part of the control subjects were genotyped for the validation study ($n = 3,290$, 97.3%). To enlarge the number of subjects and enhance statistical power for replication studies, we used genotype data obtained from other GWAS projects genotyped using the Illumina platforms for the replication control panels (Table S1). All SNPs were successfully genotyped with call rates $>0.98$ and were in Hardy-Weinberg equilibrium (HWE) in control subjects ($P>0.05$ as examined by $\chi^2$ test), except for rs2233434, which displayed a deviation from HWE ($P = 0.00091$). To evaluate possible genotyping biases between the platforms, we also genotyped rs2233434 and rs3125734 by TaqMan assays for randomly selected subjects genotyped using other genotyping platforms ($n = 376$), yielding high concordance rates of $\geq 0.99$.

## Association analysis

The associations of the SNPs were tested with the Cochran-Armitage trend test. Combined analysis was performed with the Mantel-Haenszel method. Haplotype association analysis and haplotype-based conditional association analysis were performed using Haploview v4.2 and the PLINK v1.07 program (see URLs) [40], respectively. The SNPs that were not genotyped in the GWAS were imputed using MACH 1.0.16 (see URLs), with genotype data from the 1000 Genome Project (JPT, CHB, and Han Chinese South (CHS): 177 individuals) as references (August 2010 release) [41]. All the imputed SNPs demonstrated *Rsq* values more than 0.60.

## DNA re-sequencing

Unknown variants in the coding sequences of *NFKBIE* and *RTKN2* were revealed by directly sequencing the DNA of 48 individuals affected with RA. DNA fragments were amplified with the appropriate primers (Table S13). Purification of PCR products

was performed with Exonuclease I (New England Biolabs, Ipswich, MA, USA) and shrimp alkaline phosphatase (Promega, Madison, WI, USA). The amplified DNAs were sequenced using the BigDye Terminator v3.1 Cycle Sequencing kit (Applied Biosystems), and signals were detected using an ABI 3700 DNA Analyzer (Applied Biosystems).

## Construction of haplotype-specific expression vectors

The full coding regions were amplified using cDNAs prepared from an Epstein-Barr virus-transfected lymphoblastoid B-cell line (Pharma SNP Consortium (PSC), Osaka, Japan) for *NFKBIE* (NM_004556.2) and from Jurkat cells (American Type Culture Collection (ATCC), Rockville, MD, USA) for *RTKN2* (NM_145307.2) with appropriate primers (Table S14) and DNA polymerases. PCR products were inserted into the pcDNA3.1D/V5-His-TOPO vector (Invitrogen, Camarillo, CA, USA) using the TaKaRa Ligation kit ver. 2.1 (Takara Bio Inc, Shiga, Japan), and mutagenized using the AMAP Multi Site-Directed Mutagenesis Kit (MBL, Nagoya, Japan). Each construct was then transformed into Jet Competent *Escherichia coli* cells (DH5α) (BioDynamics Laboratory Inc., Tokyo, Japan). These plasmids were purified using an Endofree Plasmid Maxi Kit (QIAGEN) after confirmation of the sequence.

## NF-κB reporter assay

Human embryonic kidney (HEK) 293A cells (Invitrogen) were cultured in Dulbecco's modified Eagle's medium (Sigma-Aldrich, St. Louis, MO, USA) supplemented with 10% fetal bovine serum (BioWest, Nuaillé, France), 1% penicillin/streptomycin (Invitrogen), and 0.1 mM MEM Non-Essential Amino Acids (Invitrogen). Various doses of the haplotype-specific expression vector (0.0025–0.02 μg for *NFKBIE* and 0.1–0.8 μg for *RTKN2*), pGL4.32[*luc2P*/NF-κB-RE/Hygro] vector (Promega) (0.05 μg and 0.0125 μg, respectively), and pRL-TK vector (an internal control for transfection efficiency) (0.45 μg and 0.15 μg, respectively) were transfected into the HEK293A cells using the Lipofectamine LTX transfection reagent (Invitrogen) according to the manufacturer's protocol. The total amounts of DNAs were adjusted with empty pcDNA3.1 vector. After 22 h, cells were incubated with 1 ng/ml TNF-α (Sigma) for 2 h or with medium alone. Cells were collected, and luciferase activity was measured using a Dual-Luciferase Reporter Assay system (Promega) and a GloMax-Multi+ Detection System (Promega). Each experiment was independently repeated three times, and sextuplicate samples were assayed each time.

## Western blotting

After 24 h of transfection as described for the NF-κB reporter assay, cells were lysed in NP-40 lysis buffer (150 mM NaCl, 1% NP-40, 50 mM Tris-HCl at pH 8.0, and a protease inhibitor cocktail), and incubated on ice for 30 min. After centrifugation, the supernatant fraction was collected and 4×Sodium dodecyl sulfate (SDS) sample buffer was added. After denaturation at 95°C for 5 min, proteins were analyzed by SDS-polyachrylamide gel electrophoresis (PAGE) on a 5% to 20% gradient gel (Wako, Osaka, Japan) and were transferred to polyvinylidene difluouride (PVDF) membranes (Millipore, Billerica, MA, USA). Target proteins on the membrane were probed with antibodies (mouse anti-V5 tag (Invitrogen), anti-β-actin-HRP (an internal control), and goat anti-mouse IgG2a-HRP (Santa Cruz Biotechnology, Santa Cruz, CA, USA)), visualized using enhanced chemiluminescence (ECL) detection reagent (GE Healthcare, Pollards Wood, UK), and detected using a LAS-3000 mini lumino-image analyzer

65

## Allele-specific transcript quantification (ASTQ) analysis

ASTQ analysis was performed as previously described [42]. Total RNAs and genomic DNAs were extracted from PBMCs and lymphoblastoid B-cell lines. cDNAs were synthesized using TaqMan reverse transcription reagents (Applied Biosystems). We selected SNPs (rs2233434 (A/G) for *NFKBIE* and rs3125734 (C/T) for *RTKN2*) as target SNPs. Allele-specific gene expression was measured by TaqMan SNP genotyping probes for these SNPs (Applied Biosystems). To make a standard curve, we selected two individuals that had homozygous genotypes of each target SNP. We mixed these DNAs at nine different ratios and detected the intensities. The $\log_2$ of (risk allele/non-risk allele intensity) for each SNP was plotted against the $\log_2$ of mixing homozygous DNAs. We generated a standard curve (linear regression line; $y = ax+b$), where y is the $\log_2$ of (risk allele/non-risk allele intensity) at a given mixing ratio, x is the $\log_2$ of the mixing ratio, a is the slope, and b is the intercept. We then measured the allelic ratio for each cDNA and genomic DNA from each individual by real-time TaqMan PCR. Based on a standard curve, we calculated the allelic ratio of cDNAs and genomic DNAs. Intensities were detected using an ABI Prism 7900HT Sequence Detection System (Applied Biosystems).

## Electrophoretic mobility shift assays (EMSA)

EMSA and preparation of nuclear extract from lymphoblastoid B-cell lines and Jurkat cells were performed as previously described [43]. Cells were cultured in RPMI-1640 medium (Sigma-Aldrich) supplemented with 10% fetal bovine serum and 1% penicillin/streptomycin. Following stimulation with 50 ng/ml phorbol myristate acetate (Sigma-Aldrich) for 2 h, cells were collected and suspended in buffer A (20 mM HEPES at pH 7.6, 20% glycerol, 10 mM NaCl, 1.5 mM MgCl$_2$, 0.2 mM EDTA at pH 8.0, 1 mM DTT, 0.1% NP-40, and a protease inhibitor cocktail) for 10 min on ice. After centrifugation, the pellets were resuspended in buffer B (which contains buffer A with 500 mM NaCl). Following incubation on ice for 30 min and centrifugation to remove cellular debris, the supernatant fraction containing nuclear proteins was collected. Oligonucleotides (31-bp) were designed that corresponded to genomic sequences surrounding the SNPs (Table S15). Single-stranded oligonucleotide probes were labeled using a Biotin 3' End DNA Labeling Kit (Pierce Biotechnology, Rockford, IL, USA), and sense and antisense oligonucleotides were then annealed. DNA-protein interactions were detected using a LightShift Chemiluminescent EMSA kit (Pierce Biotechnology). The DNA-protein complexes were separated on a non-denaturing 5% polyachrylamide gel in 1×TBE (Tris-borate-EDTA) running buffer for 60 min at 150 V. The DNA-protein complexes were then transferred from the gel onto a nitrocellulose membrane (Ambion, Carlsbad, CA, USA), and were cross-linked to the membrane by exposure to UV light. Signals were detected using a LAS-3000 mini lumino-image analyzer (Fujifilm). Allelic differences were analyzed using MultiGauge software (Fujifilm) by measuring the intensity of the bands.

## Luciferase assay

Oligonucleotides (31-bp) were designed as described for the EMSAs (Table S15), and complementary sense and antisense oligonucleotides were annealed. To construct luciferase reporter plasmids, pGL4.24[*luc2P*/minP] vector (Promega) was digested with restriction enzymes (XhoI and BglII) (Takara Bio Inc), and annealed oligonucleotide was ligated into a pGL4.24 vector

upstream of the minimal promoter. HEK293A ($n = 2.5 \times 10^5$), lymphoblastoid B-cell lines ($n = 2.0 \times 10^6$) and Jurkat ($n = 5.0 \times 10^5$) cells were transfected with the allele-specific constructs (0.4 μg, 1.8 μg and 2.5 μg, respectively) and the pRL-TK vector (0.1 μg, 0.2 μg and 0.25 μg, respectively) using the Lipofectamine LTX transfection reagent (for HEK293A and Jurkat cells) and Amaxa nucleofector kit (Lonza, Basel, Switzerland) (for lymphoblastoid B-cell lines). Cells were collected, and luciferase activity was measured as described for the NF-κB reporter assay. Each experiment was independently repeated three times and sextuplicate samples were assayed each time.

## Correlation analysis between gene expression and genotypes

The expression data in lymphoblastoid B-cell lines derived from HapMap individuals ($n = 210$; JPT, CHB, CEU, and YRI) and in primary T cells from umbilical cords of Western European individuals ($n = 85$) from the database of the Gene Expression Variation (Genevar) project were used. SNP genotypes were obtained from HapMap and 1000 Genome Project databases. The expression levels were regressed with the genotype in a liner model. The statistical significance of regression coefficients was tested using Student's $t$-test.

## Statistical analysis

We used $\chi^2$ contingency table tests to evaluate the significance of differences in allele frequency in the case-control subjects. We defined haplotype blocks using the solid spine of LD definition of Haploview v4.2, and estimated haplotype frequency and calculated pairwise LD indices ($r^2$) between pairs of polymorphisms using the Haploview program. Luciferase assay data and ASTQ analysis data were analyzed by Student's $t$-test.

## Web resources

The URLs for data presented herein are as follows:
PLINK, http://pngu.mgh.harvard.edu/~purcekk/plink
MACH, http://www.sph.umich.edu/csg/abecasis/mach/
UCSC Genome Browser, http://genome.ucsc.edu/;
Genevar, http://www.sanger.ac.uk/resources/software/genevar/
HapMap Project, http://www.HapMap.org/
1000 Genome Project, http://www.1000genomes.org
Online Mendelian Inheritance in Man (OMIM), http://www.omim.org/

## Supporting Information

**Figure S1** NF-kB activity was influenced by nsSNPs in *NFKBIE*. NF-κB activities were evaluated by luciferase assays. Allele specific construct, pGL4.32[*luc2P*/NF-κB-RE] luciferase vector, and pRL-TK vector were transfected into HEK293A cells. Four haplotypes (rs2233434-rs2233433; A-C, G-C, A-T, and G-T) were examined. (rs2233434: A = non-risk (NR), G = risk (R); rs2233433: C = NR, T = R). Twenty-two hours after transfection, cells were stimulated with medium alone (A) or TNF-α (B) for 2 h. Data represent the mean ± s.d. Each experiment was performed in sextuplicate, and experiments were independently repeated three times. *$P < 0.05$ and **$P < 1.0 \times 10^{-5}$ by Student's $t$-test. n.s.: not significant. (TIF)

**Figure S2** Allelic imbalance of expression in *NFKBIE*. ASTQ was performed using samples from individuals heterozygous for rs2233434 (G/A) in *NFKBIE*. Genomic DNAs and cDNAs were extracted from lymphoblastoid B cells ($n = 9$). The y-axis shows the $\log_2$ ratio of the transcript amounts in target SNPs (risk allele/non-risk

(Fujifilm, Tokyo, Japan). Band intensities were measured using MultiGauge software (Fujifilm).

66

allele). The top bar of the box-plot represents the maximum value and the lower bar represents the minimum value. The top of box is the third quartile, the bottom of box is the first quartile, and the middle bar is the median value. The circle is an outlier. $*P = 5.3 \times 10^{-4}$ by Student's *t*-test.
(TIF)

**Figure S3** SNP selection using *in silico* analysis in the *NFKBIE* region. Step 1: Definition of the target region. *P*-values of the SNPs in the GWAS (top) and genomic structure (middle), and the $D'$-based LD map (bottom). The green diamond shapes represent the $-\log_{10}$ of the Cochran-Armitage trend *P*-values. The dashed line indicates the significance threshold ($P < 1 \times 10^{-3}$). The LD map was drawn based on genotype data of the 1000 Genome Project (JPT, CHB and CHS: 177 samples) using Haploview software v4.2. LD blocks were defined by the solid spine method. The red box (top) represents the target region of the *in silico* analysis (Chr6: 44,336,140-44,394,125). Step 2: Target SNPs were extracted from public databases (HapMap and 1000 Genome Project). SNPs with MAF $>0.05$ were selected. Step 3: Evaluation of regulatory potential. Step 3a: The regulatory potential (RP) score was calculated for sequences surrounding the SNPs by ESPERR (evolutionary and sequence pattern extraction through reduced representations) method. SNPs with RP score $>0.1$ were selected. Step 3b: Subsequently, SNPs within the predicted, regulatory genomic elements were selected by using ChIP-seq data of transcription factor binding sites (Txn factor), histone modification sites (CTCF binding, H3K4me1, H3K4me2, H3K4me3, H3K27ac, H3K9ac) or DNase-seq data of DNase I hypersensitivity sites (DNase HS). ChIP-seq data and DNase-seq data used the signals derived from GM12878 EBV-transformed B cells. All these analyses of Steps 2 to 3 were performed by using the UCSC genome browser. Step 4: Evaluation of disease association. Association data of both genotyped (green diamonds) and imputed (black diamonds) SNPs in the GWAS samples were used. Red triangles represent 14 extracted SNPs *in silico*. The dashed line indicates the significance threshold ($P < 0.05$).
(TIF)

**Figure S4** SNP selection using *in silico* analysis in the *RTKN2* region. SNP selection in the *RTKN2* region was performed the same as in the case of the *NFKBIE* region as described in Figure S3, except that we used DNase-seq data derived from Th1, Th2, and Jurkat cells in addition to GM12878 EBV-transformed B cells.
(TIF)

**Figure S5** Results of EMSAs for candidate regulatory SNPs. Binding affinities of nuclear factors from lymphoblastoid B-cells (PSC cells) and Jurkat cells to the 31-bp sequences around each allele of the candidate regulatory SNPs were evaluated by EMSA. Nuclear factors from PSC cells were used for *NFKBIE*, and Jurkat cells were used for *RTKN2*. 14 SNPs in *NFKBIE* (A) and 10 SNPs in *RTKN2* (B) were tested. NR: non-risk allele; R: risk allele. Arrows indicate bands showing allelic differences in each SNP.
(TIF)

**Figure S6** Luciferase assays for regulatory SNPs. Transcriptional activities of the 31-bp genomic sequences around the SNPs were evaluated by luciferase assays. Each oligonucleotide was inserted into the pGL4.24[*luc2P*/minP] vector upstream of the minimal promoter (minP), and allele-specific constructs were transfected into HEK293A cells. Relative luciferase activity is expressed as the ratio of luciferase activity of each allele-specific construct to the luciferase activity of the mock construct. Data represent the mean $\pm$ s.d. Each experiment was independently repeated three times, and each sample was measured in sextuplicate. $*P < 1 \times 10^{-3}$ by

Student's *t*-test. n.s.: not significant. (A) rs2233434 and rs77986492 in the *NFKBIE* region. (B) rs3864793, rs1864836, rs4979765, and rs4979766 in the *RTKN2* region. NR: non-risk allele; R: risk allele.
(TIF)

**Figure S7** The correlation between *NFKBIE* expression and rs2233434 and rs77986492 genotypes. Linear regression analysis of the relationship between SNP genotypes and *NFKBIE* expression. Gene expression data from EBV-transformed lymphoblastoid B cell lines of HapMap individuals (JPT+CHB, CEU, and YRI). (A) rs2233434 ($n = 204$) and (B) rs77986492 ($n = 152$). The genotype classification by population: rs2233434 (JPT+CHB, AA = 61, AG = 28, GG = 1; CEU, AA = 52, AG = 2; YRI, AA = 53, AG = 72) and rs77986492 (JPT+CHB, CC = 52, CT = 24; CEU, CC = 35, CT = 2; YRI, CC = 38, CT = 1). The x-axis shows SNP genotypes and the y-axis represents the $\log_2$-transformed *NFKBIE* expression level. *R*: the correlation coefficient between *NFKBIE* expression and SNP genotype.
(TIF)

**Figure S8** The correlation between *RTKN2* expression and rs3852694 genotypes. Linear regression analysis of the relationship between the rs3852694 genotype and *RTKN2* expression. Rs3852694 was used as a proxy SNP of rs1864836 ($r^2 = 1.0$). Gene expression data in primary T cells from umbilical cords of Western European individuals ($n = 85$) were presented by using Genevar software. The x-axis shows the rs3852694 genotypes (AA, AG, GG) and the y-axis represents the $\log_2$-transformed *RTKN2* expression level. *R*: the correlation coefficient between *RTKN2* expression and rs3852694 genotype.
(TIF)

**Table S1** Summary of samples.
(DOC)

**Table S2** Association results of the GWAS and 1st replication study.
(DOC)

**Table S3** Association analysis of *NFKBIE* and *RTKN2* with autoimmune diseases.
(DOC)

**Table S4** Association analysis of nsSNPs with RA.
(DOC)

**Table S5** Haplotype association study of nsSNPs in *NFKBIE*.
(DOC)

**Table S6** Haplotype association study of nsSNPs in *RTKN2*.
(DOC)

**Table S7** Predicting the effects of nsSNPs on protein function.
(DOC)

**Table S8** Association analysis of candidate rSNPs with RA.
(DOC)

**Table S9** Haplotype association study of candidate causal SNPs in *NFKBIE*.
(DOC)

**Table S10** Haplotype association study of candidate causal SNPs in *RTKN2*.
(DOC)

**Table S11** The conditional haplotype-based association analysis of candidate causal SNPs in *RTKN2*.
(DOC)

**Table S12** Probes and Primers used for TaqMan assays.
(DOC)

**Table S13**   Primers used for DNA re-sequencing. (DOC)

**Table S14**   Primers used for construction of expression vectors. (DOC)

**Table S15**   Oligonucleotides used for EMSAs and Luciferase assays. (DOC)

## Acknowledgments

## Author Contributions

Conceived and designed the experiments: K Myouzen, Y Kochi, Y Okada, C Terao, K Ikari, K Ohmura, R Yamada, K Yamamoto. Performed the experiments: K Myouzen, Y Kochi, C Terao, A Suzuki, K Ikari, K Ohmura. Analyzed the data: K Myouzen, Y Kochi, Y Okada, C Terao, T Tsunoda, A Takahashi, R Yamada. Contributed reagents/materials/ analysis tools: M Kubo, A Taniguchi, F Matsuda, K Ohmura, S Momohara, T Mimori, H Yamanaka, N Kamatani, Y Nakamura. Wrote the paper: K Myouzen, Y Kochi, Y Okada, C Terao, K Yamamoto.

## References

1. Gabriel SE (2001) The epidemiology of rheumatoid arthritis. Rheum Dis Clin North Am 27: 269–281
2. Suzuki A, Yamada R, Chang X, Tokuhiro S, Sawada T, et al. (2003) Functional haplotypes of PADI4, encoding citrullinating enzyme peptidylarginine deiminase 4, are associated with rheumatoid arthritis. Nat Genet 34: 395–402
3. Plenge RM, Seielstad M, Padyukov L, Lee AT, Remmers EF, et al. (2007) TRAF1-C5 as a risk locus for rheumatoid arthritis–a genomewide study. N Engl J Med 357: 1199–1209
4. Wellcome Trust Case Control Consortium (2007) Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. Nature 447: 661–678
5. Gregersen PK, Amos CI, Lee AT, Lu Y, Remmers EF, et al. (2009) REL, encoding a member of the NF-kappaB family of transcription factors, is a newly defined risk locus for rheumatoid arthritis. Nat Genet 41: 820–823
6. Kochi Y, Okada Y, Suzuki A, Ikari K, Terao C, et al. (2010) A regulatory variant in CCR6 is associated with rheumatoid arthritis susceptibility. Nat Genet 42: 515–519
7. Begovich AB, Carlton VE, Honigberg LA, Schrodi SJ, Chokkalingam AP, et al. (2004) A missense single-nucleotide polymorphism in a gene encoding a protein tyrosine phosphatase (PTPN22) is associated with rheumatoid arthritis. Am J Hum Genet 75: 330–337
8. Adrianto I, Wen F, Templeton A, Wiley G, King JB, et al. (2011) Association of a functional variant downstream of TNFAIP3 with systemic lupus erythematosus. Nat Genet 43: 253–258
9. Thomas PD, Kejariwal A (2004) Coding single-nucleotide polymorphisms associated with complex vs. Mendelian disease: evolutionary evidence for differences in molecular effects. Proc Natl Acad Sci U S A 101: 15398–15403
10. Okada Y, Shimane K, Kochi Y, Tahira T, Suzuki A, et al. (2012) A Genome-Wide Association Study Identified AFF1 as a Susceptibility Locus for Systemic Lupus Eyrthematosus in Japanese. PLoS Genet 8: e1002455. doi:10.1371/journal.pgen.1002455
11. Dubois PC, Trynka G, Franke L, Hunt KA, Romanos J, et al. (2010) Multiple common variants for celiac disease influencing immune gene expression. Nat Genet 42: 295–302
12. 1000 Genomes Project Consortium (2010) A map of human genome variation from population-scale sequencing. Nature 467: 1061–1073
13. Plenge RM, Cotsapas C, Davies L, Price AL, de Bakker PI, et al. (2007) Two independent alleles at 6q23 associated with risk of rheumatoid arthritis. Nat Genet 39: 1477–1482
14. Remmers EF, Plenge RM, Lee AT, Graham RR, Hom G, et al. (2007) STAT4 and the risk of rheumatoid arthritis and systemic lupus erythematosus. N Engl J Med 357: 977–986
15. Okada Y, Terao C, Ikari K, Kochi Y, Ohmura K, et al. (2012) Meta-analysis identifies nine new loci associated with rheumatoid arthritis in the Japanese population. Nat Genet 45: 511–516
16. Li Z, Nabel GJ (1997) A new member of the I kappaB protein family, I kappaB epsilon, inhibits RelA (p65)-mediated NF-kappaB transcription. Mol Cell Biol 17: 6184–6190
17. Whiteside ST, Epinat JC, Rice NR, Israel A (1997) I kappa B epsilon, a novel member of the I kappa B family, controls RelA and cRel NF-kappa B activity. Embo J 16: 1413–1426
18. Collier FM, Gregorio-King CC, Gough TJ, Talbot CD, Walder K, et al. (2004) Identification and characterization of a lymphocytic Rho-GTPase effector: rhotekin-2. Biochem Biophys Res Commun 324: 1360–1369
19. Collier FM, Loving A, Baker A, J., McLeod J, Walder K, et al. (2009) RTKN2 Induces NF-kappaB Dependent Resistance to Intrinsic Apoptosis in HEK cells and Regulates BCL-2 Gene in Human CD4+ Lymphocytes. J Cell Death 2: 9–23
20. Makarov SS (2001) NF-kappa B in rheumatoid arthritis: a pivotal regulator of inflammation, hyperplasia, and tissue destruction. Arthritis Res 3: 200–206
21. Kolbe D, Taylor J, Elnitski L, Eswara P, Li J, et al. (2004) Regulatory potential scores from genome-wide three-way alignments of human, mouse, and rat. Genome Res 14: 700–707
22. Taylor J, Tyekucheva S, King DC, Hardison RC, Miller W, et al. (2006) ESPERR: learning strong and weak signals in genomic sequence alignments to identify functional elements. Genome Res 16: 1596–1604
23. Johnson DS, Mortazavi A, Myers RM, Wold B (2007) Genome-wide mapping of in vivo protein-DNA interactions. Science 316: 1497–1502
24. Valouev A, Johnson DS, Sundquist A, Medina C, Anton E, et al. (2008) Genome-wide analysis of transcription factor binding sites based on ChIP-Seq data. Nat Methods 5: 829–834
25. Mikkelsen TS, Ku M, Jaffe DB, Issac B, Lieberman E, et al. (2007) Genome-wide maps of chromatin state in pluripotent and lineage-committed cells. Nature 448: 553–560
26. Ernst J, Kheradpour P, Mikkelsen TS, Shoresh N, Ward LD, et al. (2011) Mapping and analysis of chromatin state dynamics in nine human cell types. Nature 473: 43–49
27. Sabo PJ, Kuehn MS, Thurman R, Johnson BE, Johnson EM, et al. (2006) Genome-scale mapping of DNase I sensitivity in vivo using tiling DNA microarrays. Nat Methods 3: 511–518
28. Dimas AS, Deutsch S, Stranger BE, Montgomery SB, Borel C, et al. (2009) Common regulatory variation impacts gene expression in a cell type-dependent manner. Science 325: 1246–1250
29. Yang TP, Beazley C, Montgomery SB, Dimas AS, Gutierrez-Arcelus M, et al. (2010) Genevar: a database and Java application for the analysis and visualization of SNP-gene associations in eQTL studies. Bioinformatics 26: 2474–2476
30. Stranger BE, Forrest MS, Dunning M, Ingle CE, Beazley C, et al. (2007) Relative impact of nucleotide and copy number variation on gene expression phenotypes. Science 315: 848–853
31. Stranger BE, Nica AC, Forrest MS, Dimas A, Bird CP, et al. (2007) Population genomics of human gene expression. Nat Genet 39: 1217–1224
32. Stahl EA, Raychaudhuri S, Remmers EF, Xie G, Eyre S, et al. (2010) Genome-wide association study meta-analysis identifies seven new rheumatoid arthritis risk loci. Nat Genet 42: 508–514
33. Chu X, Pan CM, Zhao SX, Liang J, Gao GQ, et al. (2011) A genome-wide association study identifies two new risk loci for Graves' disease. Nat Genet 43: 897–901
34. Trynka G, Hunt KA, Bockett NA, Romanos J, Mistry V, et al. (2011) Dense genotyping identifies and localizes multiple common and rare variant association signals in celiac disease. Nat Genet 43: 1193–1201
35. Li Y, Sidore C, Kang HM, Boehnke M, Abecasis GR (2011) Low-coverage sequencing: implications for design of complex trait association studies. Genome Res 21: 940–951
36. Degner JF, Pai AA, Pique-Regi R, Veyrieras JB, Gaffney DJ, et al. (2012) DNase I sensitivity QTLs are a major determinant of human expression variation. Nature 482: 390–394
37. Musone SL, Taylor KE, Lu TT, Nititham J, Ferreira RC, et al. (2008) Multiple polymorphisms in the TNFAIP3 region are independently associated with systemic lupus erythematosus. Nat Genet 40: 1062–1064
38. Raychaudhuri S, Remmers EF, Lee AT, Hackett R, Guiducci C, et al. (2008) Common variants at CD40 and other loci confer risk of rheumatoid arthritis. Nat Genet 40: 1216–1223
39. Arnett FC, Edworthy SM, Bloch DA, McShane DJ, Fries JF, et al. (1988) The American Rheumatism Association 1987 revised criteria for the classification of rheumatoid arthritis. Arthritis Rheum 31: 315–324
40. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, et al. (2007) PLINK: a tool set for whole-genome association and population-based linkage analyses. Am J Hum Genet 81: 559–575
41. Li Y, Willer C, Sanna S, Abecasis G (2009) Genotype imputation. Annu Rev Genomics Hum Genet 10: 387–406

68

42. Akamatsu S, Takata R, Ashikawa K, Hosono N, Kamatani N, et al. (2010) A functional variant in NKX3.1 associated with prostate cancer susceptibility down-regulates NKX3.1 expression. Hum Mol Genet 19: 4265–4272

43. Andrews NC, Faller DV (1991) A rapid micropreparation technique for extraction of DNA-binding proteins from limiting numbers of mammalian cells. Nucleic Acids Res 19: 2499

69

nature
genetics

# Genome-wide association analyses in east Asians identify new susceptibility loci for colorectal cancer

Wei-Hua Jia[1,16], Ben Zhang[2,16], Keitaro Matsuo[3], Aesun Shin[4], Yong-Bing Xiang[5], Sun Ha Jee[6], Dong-Hyun Kim[7], Zefang Ren[1], Qiuyin Cai[2], Jirong Long[2], Jiajun Shi[2], Wanqing Wen[2], Gong Yang[2], Ryan J Delahanty[2], Genetics and Epidemiology of Colorectal Cancer Consortium (GECCO)[8], Colon Cancer Family Registry (CCFR)[8], Bu-Tian Ji[9], Zhi-Zhong Pan[1], Fumihiko Matsuda[10], Yu-Tang Gao[5], Jae Hwan Oh[11], Yoon-Ok Ahn[12], Eun Jung Park[6], Hong-Lan Li[5], Ji Won Park[11], Jaeseong Jo[6], Jin-Young Jeong[7], Satoyo Hosono[3], Graham Casey[13], Ulrike Peters[14,15], Xiao-Ou Shu[2], Yi-Xin Zeng[1,17] & Wei Zheng[2,17]

To identify new genetic factors for colorectal cancer (CRC), we conducted a genome-wide association study in east Asians. By analyzing genome-wide data in 2,098 cases and 5,749 controls, we selected 64 promising SNPs for replication in an independent set of samples, including up to 5,358 cases and 5,922 controls. We identified four SNPs with association $P$ values of $8.58 \times 10^{-7}$ to $3.77 \times 10^{-10}$ in the combined analysis of all east Asian samples. Three of the four were replicated in a study conducted in 26,060 individuals of European descent, with combined $P$ values of $1.22 \times 10^{-10}$ for rs647161 (5q31.1), $6.64 \times 10^{-9}$ for rs2423279 (20p12.3) and $3.06 \times 10^{-8}$ for rs10774214 (12p13.32 near the *CCND2* gene), derived from meta-analysis of data from both east Asian and European-ancestry populations. This study identified three new CRC susceptibility loci and provides additional insight into the genetics and biology of CRC.

CRC is one of the most commonly diagnosed malignancies in east Asia and many other parts of the world[1]. Genetic factors have an important role in the etiology of both sporadic and familial CRC[2]. However, less than 6% of CRC cases can be explained by rare, high-penetrance variants in the CRC susceptibility genes identified to date, such as the *APC, SMAD4, AXIN2, BMPR1A, POLD1, STK11, MUTYH* and DNA mismatch repair genes[2]. Over the past two decades, many candidate gene studies have evaluated common genetic risk factors for CRC; only a few of these have been replicated in subsequent studies[3]. Recent genome-wide association studies (GWAS) have identified approximately 15 common genetic susceptibility loci for CRC[4–12]. However, these newly identified genetic factors, along with known high-penetrance variations in CRC susceptibility genes, explain less than 15% of the heritability for this common malignancy[10,11]. Furthermore, with the exception of a small study conducted in Japan[12], all other GWAS have been conducted in populations of European ancestry, which differ from other populations in certain features of genetic architecture. Many of the variants discovered in populations of European ancestry show only weak or no association with CRC in other ancestry groups[13]. Therefore, additional GWAS are needed, particularly in populations not of European ancestry, to fully uncover the genetic basis for CRC susceptibility.

In 2009, we initiated the Asia Colorectal Cancer Consortium (ACCC), a GWAS in east Asians, to search for previously unknown genetic risk factors for CRC. The discovery stage (stage 1) consisted of five GWAS conducted in China, Korea and Japan, including 2,293 CRC cases and 5,780 controls (**Supplementary Table 1**). Cases and controls were genotyped using several SNP arrays, including the Affymetrix Genome-Wide Human SNP Array 6.0 (906,602 SNPs), the Affymetrix Genome-Wide Human SNP Array 5.0 (443,104 SNPs), the Illumina Infinium HumanHap610 BeadChip (592,044 SNPs), the Illumina Human610-Quad BeadChip (620,901 SNPs) and the Illumina HumanOmniExpress BeadChip (729,462 SNPs) (**Supplementary Table 1**). After quality control exclusions as described previously[14–17], 2,098 cases and 5,749 controls remained for this study (**Supplementary Tables 1 and 2**). Also excluded from the analyses were SNPs with call rate of <95%, genotype concordance rate of <95%

71

**Table 1  Association of CRC risk with the top four risk variants identified in east Asian samples**

| | | | | | | Cases | | Controls | | Per-allele association | | Heterogeneity | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SNP | Alleles[a] | Chr. | Gene[b] | Location (bp)[c] | Stage | Sample size | MAF | Sample size | MAF | OR (95% CI)[d] | $P_{trend}$ | $P$[e] | $I^2$ |
| rs10774214 | T/C | 12p13.32 | CCND2 | 4238613 | GWAS | 2,098 | 0.373 | 5,749 | 0.348 | 1.20 (1.09–1.32) | $2.03 \times 10^{-4}$ | | |
| | | | | | Replication | 5,197 | 0.381 | 5,797 | 0.355 | 1.16 (1.09–1.23) | $5.80 \times 10^{-7}$ | | |
| | | | | | Overall | 7,295 | 0.379 | 11,546 | 0.352 | 1.17 (1.11–1.23) | $5.48 \times 10^{-10}$ | 0.615 | 0% |
| rs647161 | A/C | 5q31.1 | PITX1 | 134526991 | GWAS | 2,098 | 0.353 | 5,749 | 0.308 | 1.22 (1.12–1.33) | $3.29 \times 10^{-6}$ | | |
| | | | | | Replication | 5,217 | 0.344 | 5,815 | 0.319 | 1.14 (1.07–1.21) | $1.15 \times 10^{-5}$ | | |
| | | | | | Overall | 7,315 | 0.347 | 11,564 | 0.313 | 1.17 (1.11–1.22) | $3.77 \times 10^{-10}$ | 0.444 | 0% |
| rs2423279 | C/T | 20p12.3 | HAO1 | 7760350 | GWAS | 2,098 | 0.339 | 5,749 | 0.307 | 1.16 (1.07–1.26) | $4.96 \times 10^{-4}$ | | |
| | | | | | Replication | 5,227 | 0.315 | 5,811 | 0.297 | 1.13 (1.06–1.19) | $1.22 \times 10^{-4}$ | | |
| | | | | | Overall | 7,325 | 0.322 | 11,560 | 0.302 | 1.14 (1.08–1.19) | $2.29 \times 10^{-7}$ | 0.331 | 12% |
| rs1665650 | T/C | 10q26.12 | HSPA12A | 118477090 | GWAS | 2,098 | 0.346 | 5,749 | 0.310 | 1.20 (1.10–1.31) | $3.88 \times 10^{-5}$ | | |
| | | | | | Replication | 5,192 | 0.328 | 5,808 | 0.320 | 1.10 (1.04–1.17) | 0.0018 | | |
| | | | | | Overall | 7,290 | 0.333 | 11,557 | 0.315 | 1.13 (1.08–1.19) | $8.58 \times 10^{-7}$ | 0.404 | 4% |

Chr., chromosome; OR, odds ratio; CI, confidence interval.
[a]Minor/major allele for east Asians. OR was estimated for the minor allele. [b]Closest gene. [c]Location based on NCBI Human Genome Build 36.3. [d]Adjusted for age, sex, the first ten principal components (stage 1) and study site. [e]$P$ for heterogeneity across studies in GWAS and replication was calculated using Cochran's $Q$ test.

between positive control samples, minor allele frequency (MAF) of <5% or $P$ value for Hardy-Weinberg equilibrium of $<1.0 \times 10^{-5}$ in controls for each study. Imputation was conducted for each study following the MaCH algorithm[18] using phased HapMap 2 Han Chinese in Beijing, China (CHB) and Japanese in Tokyo, Japan (JPT) samples as the reference. No apparent genetic admixture was detected, except for one sample from KCPS-II (**Supplementary Fig. 1**). Associations between CRC risk and each of the genotyped and imputed SNPs were evaluated using logistic regression within each study after adjusting for age, sex and the first ten principal components using mach2dat[18]. Meta-analyses were conducted under a fixed-effects model using the METAL program[19]. There was little evidence for inflation in the association test statistics for any of the five studies (genomic inflation factor ($\lambda$) range of 1.02 to 1.04) or for all studies combined ($\lambda = 1.01$) (**Supplementary Fig. 2** and **Supplementary Table 1**). The observed number of SNPs with small $P$ values was slightly larger than that expected by chance (**Supplementary Fig. 2**).

Multiple genomic locations were found that were potentially related to CRC risk (**Supplementary Fig. 3**). Nine SNPs identified from published GWAS conducted in populations of European ancestry showed associations with CRC risk at $P < 0.05$ in stage 1 (data not shown). To improve the statistical power for evaluating these SNPs, we genotyped 6,476 additional samples to bring the total sample size to 5,252 cases and 9,071 controls. Except for the 2 SNPs that are monomorphic in east Asians (rs6691170 and rs16892766), all 16 of the other SNPs identified from published GWAS conducted in European-ancestry populations showed association with CRC risk in the same direction as reported previously (**Supplementary Table 3**). A significant association with CRC risk at $P < 0.05$ was found for 13 SNPs, including rs6687758, rs10936599, rs10505477, rs6983267, rs7014346, rs10795668, rs3802842, rs4444235, rs4779584, rs9929218, rs4939827, rs10411210 and rs961523. Except for two SNPs (rs6983267 and rs4779584), no statistically significant heterogeneity at $P < 0.05$ was observed between east Asian and European-ancestry populations (**Supplementary Table 3**).

To identify new genetic factors for CRC, we selected 64 SNPs for replication in an independent set of 5,358 cases and 5,922 controls recruited in 5 studies conducted in China, Korea and Japan (**Supplementary Table 2**). SNPs were selected from among those

that (i) had MAF of >5%; (ii) showed no heterogeneity across studies ($P_{het} > 0.05$ and $I^2 < 25\%$); (iii) were not in linkage disequilibrium (LD; $r^2 < 0.2$) with any known CRC risk variant reported from previous GWAS; (iv) had high imputation quality in each of the five studies (RSQ > 0.5); and (v) were associated at $P < 0.01$ in the combined analysis of all five studies included in stage 1. These criteria were used to prioritize SNPs for replication.

Of the 64 SNPs evaluated in stage 2, 7 showed association with CRC risk at $P < 0.05$ with a direction of association consistent with that observed in stage 1 (**Table 1** and **Supplementary Table 4**). In the combined analysis of data from stages 1 and 2, $P$ values for associations with two SNPs (rs647161 at 5q31.1, odds ratio (OR) = 1.17, $P = 3.77 \times 10^{-10}$, and rs10774214 at 12p13.32, OR = 1.17, $P = 5.48 \times 10^{-10}$) were lower than the conventional genome-wide significance level of $5.0 \times 10^{-8}$, providing convincing evidence for an association of these SNPs with CRC risk (**Table 1**). An additional SNP, rs2423279, showed a significant association in stage 2 after Bonferroni correction (corrected $P < 7.8 \times 10^{-4}$) but did not reach the conventional GWAS significance level for association with CRC risk in the combined analysis of all samples (OR = 1.14, $P = 2.29 \times 10^{-7}$). The association between CRC risk and each of these three SNPs was consistent across most studies (**Fig. 1**). Results for the other four SNPs that replicated in stage 2 at $P < 0.05$ (rs1665650, rs2850966, rs1580743 and rs4503064) are also presented (**Supplementary Table 4**), including one SNP (rs1665650) with an association $P$ value of $8.58 \times 10^{-7}$ in the combined analysis of all data from both stages (**Table 1**).

We next evaluated these top four SNPs (**Table 1**) using data from GWAS in the Genetics and Epidemiology of Colorectal Cancer Consortium (GECCO) and the Colon Cancer Family Registry (CCFR), which together include 11,870 cases and 14,190 controls of European ancestry[4,20,21]. Three of the four SNPs were replicated in the GECCO and CCFR sample, although the strength of the associations was weaker than in east Asians (**Table 2**). These results provide independent support of our findings in the east Asian population. Meta-analyses of data from both east Asian and European-ancestry populations provided strong evidence for associations of CRC risk with three SNPs, with $P$ values all below the genome-wide significance threshold of $5 \times 10^{-8}$ (**Table 2**). The weaker associations observed in European-ancestry populations could be explained, in part, by differences in LD patterns at these loci for east
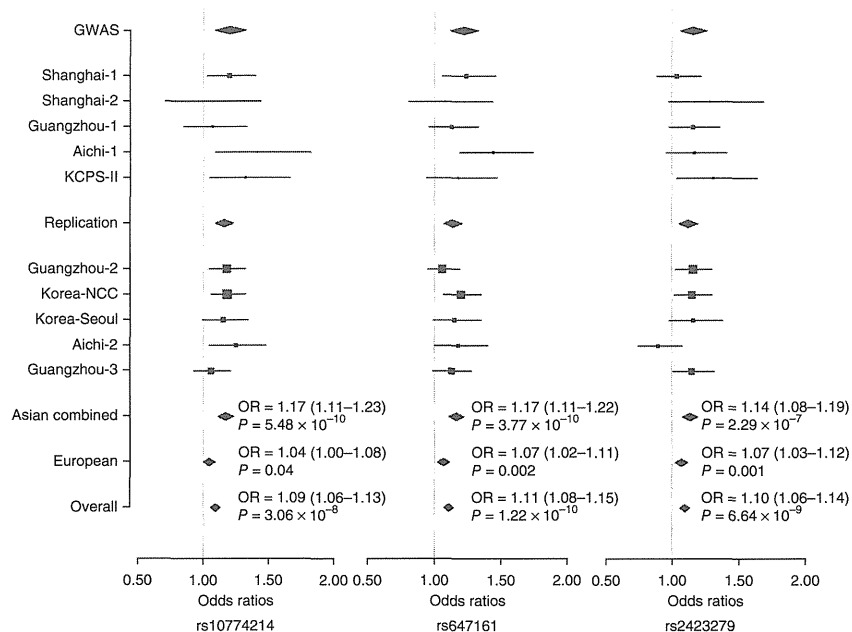
72

**Figure 1** Forest plots for the three SNPs showing evidence of an association with CRC risk. Per-allele ORs are presented, with the area of each box proportional to the inverse variance weight of the estimate. Horizontal lines represent 95% confidence intervals.

(ref. 22). CCND2 is closely related to CCND1, a well-established human oncogene[22,23]. Although CCND2 has been less well studied than CCND1, several studies, including The Cancer Genome Atlas (TCGA), have shown that CCND2 is overexpressed in a substantial proportion of human colorectal tumors[22-25]. Overexpression of this cyclin may be an independent predictor of survival in individuals with CRC[24]. Several other genes, including PARP11, FGF23, FGF6, C12orf5 and RAD51AP1, are also in close proximity to the SNP identified in our study, of which both C12orf5 (also known as TIGAR, encoding TP53-induced glycolysis and apoptosis regulator) and RAD51AP1 were found to be overexpressed in CRC tissue included in TCGA[25]. rs10774214 is in strong LD with several SNPs that are located in potential transcription factor–binding sites, as determined using the TRANSFAC database[26]. Additional research may be warranted regarding possible mechanisms by which this SNP is related to CRC risk.

The rs647161 SNP is located on chromosome 5q31.1, where a cluster of SNPs were associated with CRC risk (**Fig. 2b**). Of the genes in this region (including PITX1, CATSPER3, PCBD2, MIR4461 and H2AFY), PITX1 is the closest to rs647161 (approximately 129 kb upstream). The PITX1 gene (encoding paired-like homeodomain 1) has been described as a tumor suppressor gene and may be involved in the tumorigenesis of multiple human cancers[27-31], including CRC[27,32]. PITX1 has been reported to suppress tumorigenicity by downregulating the RAS pathway, which is frequently altered in colorectal tumors[27]. Inhibition of PITX1 induces the RAS pathway and tumorigenicity, and restoring PITX1 in colon cancer cells inhibits tumorigenicity[27]. It also has been reported that PITX1 may activate TP53 (ref. 33) and regulate telomerase activity[34]. Consistent with its possible function as a tumor suppressor gene, PITX1 has been found to be downregulated in human cancer tissue samples and cell lines[27-30,32]. CRC tissue expressing wild-type KRAS showed significantly lower expression of PITX1 than tissue with mutant KRAS[32]. Most recently, low PITX1 expression was found to be associated with poor survival in individuals with CRC[35]. In addition, rs6596201, which is in moderate LD with rs647161 ($r^2 = 0.25$), is an expression quantitative trait locus (eQTL) ($P = 2.42 \times 10^{-28}$) for the PITX1 gene[36]. Several other genes at this locus, including C5orf24, H2AFY and NEUROG1, were also found to be highly expressed in colorectal tumors included in TCGA ($P < 0.001$)[25]. Additional studies are warranted to explore a possible role for these genes in the etiology of CRC.

Asians and Europeans (**Supplementary Fig. 4**). It is possible that causal variants in these regions are tagged by different SNPs in these two populations or that there is allelic heterogeneity, in which different underlying causal variants exist in populations of Asian and European ancestry. The difference in LD structure between Asian and European descendants and possible allelic heterogeneity in these two populations might explain, in part, why these loci were not discovered in previous studies conducted in individuals of European ancestry. The fourth SNP evaluated in the GECCO and CCFR sample, rs1665650, however, was not replicated in individuals with European ancestry (OR = 0.96, P = 0.05).

Stratification analyses showed that the association of CRC risk with each of the three replicated SNPs was generally consistent in Chinese, Korean and Japanese individuals ($P_{het} > 0.05$), although the association with rs2423279 was not statistically significant in Japanese, perhaps owing to a small sample size (**Supplementary Table 5**). Associations of these three SNPs with CRC risk were similar for men and women ($P_{het} > 0.05$) (**Supplementary Table 6**).

The rs10774214 SNP is located just 15 kb upstream of CCND2, the gene encoding cyclin D2 (**Fig. 2a**), a member of the D-type cyclin family, which also includes cyclins D1 and D3. These cyclins have a critical role in cell cycle control (from G1 to S phase) through activation of cyclin-dependent kinases (CDKs), primarily CDK4 and CDK6

**Table 2** Association of CRC risk with the top three risk variants in European descendants and east Asian and European descendants combined

| SNP | Alleles[a] | MAF[b] | | European-ancestry populations[c] | | | | East Asian and European-ancestry populations combined[c] | | | |
| | | Cases | Controls | Cases | Controls | OR (95% CI) | $P_{meta}$ | Cases | Controls | OR (95% CI) | $P_{meta}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| rs10774214 | T/C | 0.385 | 0.379 | 11,870 | 14,190 | 1.04 (1.00–1.09) | 0.040 | 19,165 | 25,736 | 1.09 (1.06–1.13) | $3.06 \times 10^{-8}$ |
| rs647161 | A/C | 0.680 | 0.667 | 11,870 | 14,190 | 1.07 (1.02–1.11) | 0.002 | 19,185 | 25,754 | 1.11 (1.08–1.15) | $1.22 \times 10^{-10}$ |
| rs2423279 | C/T | 0.263 | 0.252 | 11,870 | 14,190 | 1.07 (1.03–1.12) | 0.001 | 19,195 | 25,750 | 1.10 (1.06–1.14) | $6.64 \times 10^{-9}$ |

[a]Alleles (minor/major) for east Asians. [b]MAF in European-ancestry populations. [c]Summary statistics were generated using inverse variance–weighted fixed-effects meta-analysis.
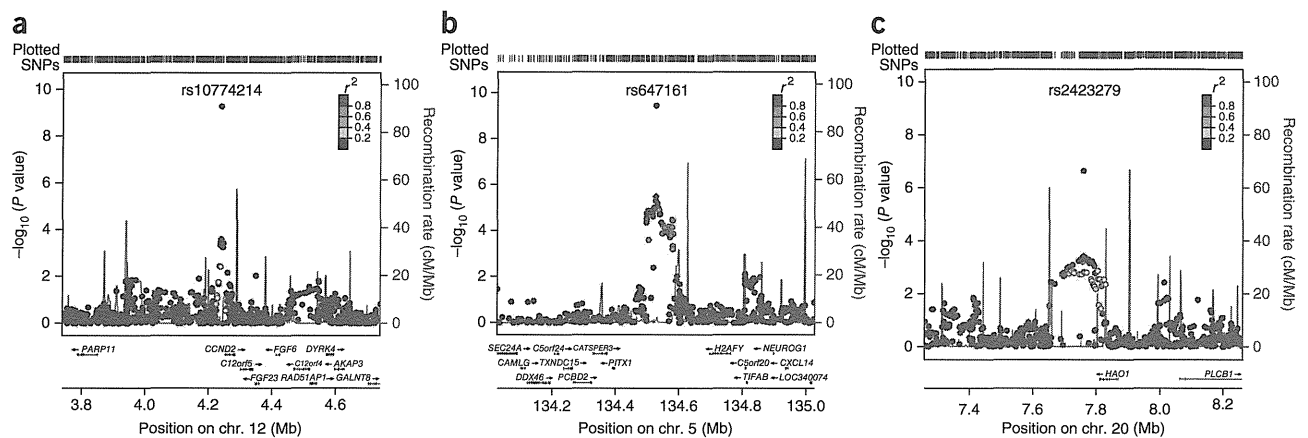
73

**Figure 2** Regional plots of association results and recombination rates for the three SNPs showing evidence of association with CRC risk. Genotyped and imputed data from GWAS samples are plotted on the basis of their chromosomal position in NCBI Human Genome Build 36.3. For each region, the SNP selected for stage 2 replication is denoted with a diamond, and the P value from the combined analysis of stage 1 and 2 data is provided. (a–c) Data are shown for rs10774214 (a), rs647161 (b) and rs2423279 (c).

The rs2423279 SNP is located on chromosome 20p12.3, close to the *HAO1* and *PLCB1* genes (**Fig. 2c**). *HAO1* encodes hydroxy-acid oxidase, which oxidizes 2-hydroxyacid. *PLCB1* encodes phospholipase C-β1, which has an important role in the intracellular transduction of many extracellular signals. Overexpression of the *PLCB1* gene has been observed in CRC tissue[25]. Possible mechanisms by which these genes are involved in CRC carcinogenesis are unknown. The rs2423279 SNP is 1,408,069 bp downstream of rs961253, a SNP previously identified in a European GWAS as being associated with CRC risk[10]. However, these two SNPs are not correlated in east Asians ($r^2 = 0$) or in Europeans ($r^2 = 0$). Adjustment for rs961253 did not change the results for rs2423279 (data not shown).

To our knowledge, this is the largest GWAS performed for CRC in east Asians, a population that differs from populations of European ancestry in CRC risk and certain aspects of genetic architecture. Results from our study, along with data from a large study conducted in a population of European ancestry, provide convincing evidence of associations with CRC risk for three new independent susceptibility loci at 5q31.1, 12p13.32 and 20p12.3. Results from this study provide new insights into the genetics and biology of CRC.

**URLs.** Cancer Genetic Markers of Susceptibility (CGEMS), http://cgems.cancer.gov/; Database of Genotypes and Phenotypes (dbGaP), http://www.ncbi.nlm.nih.gov/gap; EIGENSTRAT, http://genepath.med.harvard.edu/~reich/EIGENSTRAT.htm; eqtl.uchicago.edu, http://eqtl.uchicago.edu/Home.html; GTEx eQTL Browser, http://www.ncbi.nlm.nih.gov/gtex/GTEX2/gtex.cgi; Haploview, http://www.broad.mit.edu/mpg/haploview/; HapMap Project, http://hapmap.ncbi.nlm.nih.gov/; IntOGen, http://www.intogen.org/home; LocusZoom, http://csg.sph.umich.edu/locuszoom/; MaCH 1.0, http://www.sph.umich.edu/csg/abecasis/MACH/; mach2dat, http://genome.sph.umich.edu/wiki/Mach2dat:_Association_with_MACH_output; METAL, http://www.sph.umich.edu/csg/abecasis/Metal/; PLINK version 1.07, http://pngu.mgh.harvard.edu/~purcell/plink/; R version 2.13.0, http://www.r-project.org/; SAS version 9.2, http://www.sas.com/; SNAP, http://www.broadinstitute.org/mpg/snap/; TRANSFAC, http://www.gene-regulation.com/pub/databases.html; UCSC Genome Browser, http://genome.ucsc.edu/; WHI investigators, https://cleo.whi.org/researchers/SitePages/Write%20a%20Paper.aspx.

**METHODS**
Methods and any associated references are available in the online version of the paper.

*Note: Supplementary information is available in the online version of the paper.*

74

University of Washington, Seattle, Washington, USA), Graham Casey (Department of Preventive Medicine, University of Southern California, Los Angeles, California, USA), Andrew T. Chan (Division of Gastroenterology, Harvard Medical School and Massachusetts General Hospital, Boston, Massachusetts, USA, and Channing Division of Network Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, Massachusetts, USA), Jenny Chang-Claude (Division of Cancer Epidemiology, German Cancer Research Center, Heidelberg, Germany), Stephen J. Chanock (Division of Cancer Epidemiology and Genetics, National Cancer Institute, US NIH, Bethesda, Maryland, USA), Lin S. Chen (Department of Health Studies, University of Chicago, Chicago, Illinois, USA), Gerhard A. Coetzee (Keck School of Medicine, University of Southern California, Los Angeles, California, USA), Simon G. Coetzee (Keck School of Medicine, University of Southern California, Los Angeles, California, USA), David V. Conti (Department of Preventive Medicine, University of Southern California, Los Angeles, California, USA), Keith Curtis (Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA), David Duggan (Translational Genomics Research Institute, Phoenix, Arizona, USA), Todd L. Edwards (Division of Epidemiology, Department of Medicine, Vanderbilt University School of Medicine, Nashville, Tennessee, USA), Charles S. Fuchs (Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, Massachusetts, USA, and Channing Division of Network Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, Massachusetts, USA), Steven Gallinger (Department of Surgery, Mount Sinai Hospital, Toronto, Ontario, Canada, and Samuel Lunenfeld Research Institute, Toronto, Ontario, Canada), Edward L. Giovannucci (Channing Division of Network Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, Massachusetts, USA, and Departments of Epidemiology and Nutrition, Harvard School of Public Health, Boston, Massachusetts, USA), Stephanie M. Gogarten (School of Public Health, University of Washington, Seattle, Washington, USA), Stephen B. Gruber (Norris Comprehensive Cancer Center, University of Southern California, Los Angeles, California, USA), Robert W. Haile (Department of Preventive Medicine, University of Southern California, Los Angeles, California, USA), Tabitha A. Harrison (Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA), Richard B. Hayes (Division of Epidemiology, Department of Environmental Medicine, New York University School of Medicine, New York, New York, USA), Michael Hoffmeister (Division of Clinical Epidemiology and Aging Research, German Cancer Research Center, Heidelberg, Germany), John L. Hopper (Melbourne School of Population Health, The University of Melbourne, Melbourne, Victoria, Australia), Li Hsu (Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA, and Department of Biostatistics, University of Washington, Seattle, Washington, USA), Thomas J. Hudson (Department of Medical Biophysics, University of Toronto, Toronto, Ontario, Canada, and Department of Molecular Genetics, University of Toronto, Toronto, Ontario, Canada), David J. Hunter (Department of Epidemiology, Harvard School of Public Health, Boston, Massachusetts, USA), Carolyn M. Hutter (Division of Cancer Control and Population Sciences, National Cancer Institute, US NIH, Bethesda, Maryland, USA), Rebecca D. Jackson (Division of Endocrinology, Diabetes, and Metabolism, Ohio State University, Columbus, Ohio, USA), Mark A. Jenkins (Melbourne School of Population Health, The University of Melbourne, Melbourne, Victoria, Australia), Shuo Jiao (Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA), Charles Kooperberg (Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA), Sébastien Küry (Service de Génétique Médicale, CHU Nantes, Nantes, France), Andrea Z. LaCroix (Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA), Cathy C. Laurie (Department of Biostatistics, University of Washington, Seattle, Washington, USA), Cecelia A. Laurie (Department of Biostatistics, University of Washington, Seattle, Washington, USA), Loic Le Marchand (Cancer Epidemiology Program, University of Hawaii Cancer Center, Honolulu, Hawaii, USA), Mathieu Lemire (Ontario Institute for Cancer Research, Toronto, Ontario, Canada), David Levine (School of Public Health, University of Washington, Seattle, Washington, USA), Noralane M. Lindor (Department of Health Sciences Research, Mayo Clinic, Scottsdale, Arizona, USA), Yan Liu (Stephens and Associates, Carrollton, Texas, USA), Jing Ma (Channing Division of Network Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, Massachusetts, USA), Karen W. Makar (Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA), Polly A. Newcomb (Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA, and Department of Epidemiology, University of Washington School of Public Health, Seattle, Washington, USA), Ulrike Peters (Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA, and Department of Epidemiology, University of Washington School of Public Health, Seattle, Washington, USA), John D. Potter (Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA, Department of Epidemiology, University of Washington School of Public Health, Seattle, Washington, USA, and Centre for Public Health Research, Massey University, Palmerston North, New Zealand), Ross L. Prentice (Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA), Conghui Qu (Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA), Thomas Rohan (Department of Epidemiology and Population Health, Albert Einstein College of Medicine, Yeshiva University, Bronx, New York, USA), Robert E. Schoen (Department of Medicine and Epidemiology, University of Pittsburgh Medical Center, Pittsburgh, Pennsylvania, USA), Fredrick R. Schumacher (Department of Preventive Medicine, University of Southern California, Los Angeles, California, USA), Daniela Seminara (Division of Cancer Control and Population Sciences, National Cancer Institute, US NIH, Bethesda, Maryland, USA), Martha L. Slattery (Department of Internal Medicine, University of Utah Health Sciences Center, Salt Lake City, Utah, USA), Darin Taverna (Translational Genomics Research Institute, Phoenix, Arizona, USA), Stephen N. Thibodeau (Department of Laboratory Medicine, Mayo Clinic, Rochester, Minnesota, USA, and Department of Pathology and Laboratory Genetics, Mayo Clinic, Rochester, Minnesota, USA), Cornelia M. Ulrich (Division of Preventive Oncology, German Cancer Research Center, Heidelberg, Germany, Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA, and Department of Epidemiology, University of Washington School of Public Health, Seattle, Washington, USA), Raakhee Vijayaraghavan (Genetic Basis of Human Disease Division, Translational Genomics Research Institute, Phoenix, Arizona, USA), Bruce Weir (Department of Biostatistics, University of Washington, Seattle, Washington, USA), Emily White (Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA, and Department of Epidemiology, University of Washington School of Public Health, Seattle, Washington, USA) and Brent W. Zanke (Clinical Epidemiology Program, Ottawa Hospital Research Institute, Ottawa, Ontario, Canada).

**AUTHOR CONTRIBUTIONS**
W.Z. conceived and directed ACCC as well as the Shanghai-Vanderbilt Colorectal Cancer Genetics Project. W.-H.J., Y.-X.Z., K.M., A.S., Y.-B.X., S.H.J., D.-H.K., U.P.

and G.C. directed CRC projects at Guangzhou, Aichi, Korea-NCC, Shanghai, KCPS-II, Korea-Seoul, GECCO and CCFR, respectively. B.Z., Q.C. and W.W. coordinated the project. Q.C. directed laboratory operations. J.S. performed genotyping experiments. B.Z., J.L. and W.W. performed statistical analyses. W.Z. wrote the manuscript with substantial contributions from B.Z., Q.C., J.L., X.-O.S. and R.J.D. Z.R., G.Y., B.-T.J., Z.-Z.P., F.M., Y.-T.G., J.H.O., Y.-O.A., E.J.P., H.-L.L., J.W.P., J.J., J.-Y.J. and S.H. contributed to data and biological sample collection in the original studies included in ACCC and contributed to manuscript revision. Members of GECCO and CCFR contributed to data and biological sample collection in studies included in these consortia.

1. Jemal, A. *et al.* Global cancer statistics. *CA Cancer J. Clin.* **61**, 69–90 (2011).
2. de la Chapelle, A. Genetic predisposition to colorectal cancer. *Nat. Rev. Cancer* **4**, 769–780 (2004).
3. Dong, L.M. *et al.* Genetic susceptibility to cancer: the role of polymorphisms in candidate genes. *J. Am. Med. Assoc.* **299**, 2423–2436 (2008).
4. Zanke, B.W. *et al.* Genome-wide association scan identifies a colorectal cancer susceptibility locus on chromosome 8q24. *Nat. Genet.* **39**, 989–994 (2007).
5. Tomlinson, I. *et al.* A genome-wide association scan of tag SNPs identifies a susceptibility variant for colorectal cancer at 8q24.21. *Nat. Genet.* **39**, 984–988 (2007).
6. Broderick, P. *et al.* A genome-wide association study shows that common alleles of *SMAD7* influence colorectal cancer risk. *Nat. Genet.* **39**, 1315–1317 (2007).
7. Jaeger, E. *et al.* Common genetic variants at the *CRAC1* (*HMPS*) locus on chromosome 15q13.3 influence colorectal cancer risk. *Nat. Genet.* **40**, 26–28 (2008).
8. Tenesa, A. *et al.* Genome-wide association scan identifies a colorectal cancer susceptibility locus on 11q23 and replicates risk loci at 8q24 and 18q21. *Nat. Genet.* **40**, 631–637 (2008).
9. Tomlinson, I.P. *et al.* A genome-wide association study identifies colorectal cancer susceptibility loci on chromosomes 10p14 and 8q23.3. *Nat. Genet.* **40**, 623–630 (2008).
10. Houlston, R.S. *et al.* Meta-analysis of genome-wide association data identifies four new susceptibility loci for colorectal cancer. *Nat. Genet.* **40**, 1426–1435 (2008).
11. Houlston, R.S. *et al.* Meta-analysis of three genome-wide association studies identifies susceptibility loci for colorectal cancer at 1q41, 3q26.2, 12q13.13 and 20q13.33. *Nat. Genet.* **42**, 973–977 (2010).
12. Cui, R. *et al.* Common variant in 6q26-q27 is associated with distal colon cancer in an Asian population. *Gut* **60**, 799–805 (2011).
13. He, J. *et al.* Generalizability and epidemiologic characterization of eleven colorectal cancer GWAS hits in multiple populations. *Cancer Epidemiol. Biomarkers Prev.* **20**, 70–81 (2011).
14. Zheng, W. *et al.* Genome-wide association study identifies a new breast cancer susceptibility locus at 6q25.1. *Nat. Genet.* **41**, 324–328 (2009).
15. Bei, J.X. *et al.* A genome-wide association study of nasopharyngeal carcinoma identifies three new susceptibility loci. *Nat. Genet.* **42**, 599–603 (2010).
16. Jee, S.H. *et al.* Adiponectin concentrations: a genome-wide association study. *Am. J. Hum. Genet.* **87**, 545–552 (2010).
17. Nakata, I. *et al.* Association between the *SERPING1* gene and age-related macular degeneration and polypoidal choroidal vasculopathy in Japanese. *PLoS ONE* **6**, e19108 (2011).
18. Li, Y., Willer, C.J., Ding, J., Scheet, P. & Abecasis, G.R. MaCH: using sequence and genotype data to estimate haplotypes and unobserved genotypes. *Genet. Epidemiol.* **34**, 816–834 (2010).
19. Willer, C.J., Li, Y. & Abecasis, G.R. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* **26**, 2190–2191 (2010).
20. Peters, U. *et al.* Meta-analysis of new genome-wide association studies of colorectal cancer risk. *Hum. Genet.* **131**, 217–234 (2012).
21. Figueiredo, J.C. *et al.* Genotype-environment interactions in microsatellite stable/microsatellite instability–low colorectal cancer: results from a genome-wide association study. *Cancer Epidemiol. Biomarkers Prev.* **20**, 758–766 (2011).
22. Musgrove, E.A., Caldon, C.E., Barraclough, J., Stone, A. & Sutherland, R.L. Cyclin D as a therapeutic target in cancer. *Nat. Rev. Cancer* **11**, 558–572 (2011).
23. Mermelshtein, A. *et al.* Expression of D-type cyclins in colon cancer and in cell lines from colon carcinomas. *Br. J. Cancer* **93**, 338–345 (2005).
24. Sarkar, R. *et al.* Expression of cyclin D2 is an independent predictor of the development of hepatic metastasis in colorectal cancer. *Colorectal Dis.* **12**, 316–323 (2010).
25. Gundem, G. *et al.* IntOGen: integration and data mining of multidimensional oncogenomic data. *Nat. Methods* **7**, 92–93 (2010).
26. Matys, V. *et al.* TRANSFAC and its module TRANSCompel: transcriptional gene regulation in eukaryotes. *Nucleic Acids Res.* **34**, D108–D110 (2006).
27. Kolfschoten, I.G. *et al.* A genetic screen identifies PITX1 as a suppressor of RAS activity and tumorigenicity. *Cell* **121**, 849–858 (2005).
28. Chen, Y. *et al.* Decreased *PITX1* homeobox gene expression in human lung cancer. *Lung Cancer* **55**, 287–294 (2007).
29. Chen, Y.N., Chen, H., Xu, Y., Zhang, X. & Luo, Y. Expression of pituitary homeobox 1 gene in human gastric carcinogenesis and its clinicopathological significance. *World J. Gastroenterol.* **14**, 292–297 (2008).
30. Lord, R.V. *et al.* Increased *CDX2* and decreased *PITX1* homeobox gene expression in Barrett's esophagus and Barrett's-associated adenocarcinoma. *Surgery* **138**, 924–931 (2005).
31. Nagel, S. *et al.* Activation of paired-homeobox gene *PITX1* by del(5)(q31) in T-cell acute lymphoblastic leukemia. *Leuk. Lymphoma* **52**, 1348–1359 (2011).
32. Watanabe, T. *et al.* Differential gene expression signatures between colorectal cancers with and without *KRAS* mutations: crosstalk between the KRAS pathway and other signalling pathways. *Eur. J. Cancer* **47**, 1946–1954 (2011).
33. Liu, D.X. & Lobie, P.E. Transcriptional activation of p53 by Pitx1. *Cell Death Differ.* **14**, 1893–1907 (2007).
34. Qi, D.L. *et al.* Identification of *PITX1* as a *TERT* suppressor gene located on human chromosome 5. *Mol. Cell Biol.* **31**, 1624–1636 (2011).
35. Knösel, T. *et al.* Loss of desmocollin 1-3 and homeobox genes *PITX1* and *CDX2* are associated with tumor progression and survival in colorectal carcinoma. *Int. J. Colorectal Dis.* **27**, 1391–1399 (2012).
36. Zeller, T. *et al.* Genetics and beyond--the transcriptome of human monocytes and disease susceptibility. *PLoS ONE* **5**, e10693 (2010).

76

# ONLINE METHODS

**Study populations.** After quality control filtering, 7,456 cases and 11,671 controls from 10 studies were included in the consortium (**Supplementary Table 2**). Detailed descriptions of participating studies and demographic characteristics of study participants are provided in the **Supplementary Note**. Briefly, the consortium included 10,730 Chinese participants, 5,544 Korean participants and 2,853 Japanese participants. Chinese participants were from five studies: the Shanghai Study 1 (Shanghai-1, $n = 3,102$), the Shanghai Study 2 (Shanghai-2, $n = 485$), the Guangzhou Study 1 (Guangzhou-1, $n = 1,613$), the Guangzhou Study 2 (Guangzhou-2, $n = 2,892$) and the Guangzhou Study 3 (Guangzhou-3, $n = 2,638$). Korean participants were from three studies: the Korean Cancer Prevention Study-II (KCPS-II, $n = 1,301$), the Seoul Study ($n = 1,522$) and the Korea–National Cancer Center (Korea-NCC) Study ($n = 2,721$). Japanese participants were from two studies: the Aichi Study 1 (Aichi-1, $n = 1,346$) and the Aichi Study 2 (Aichi-2, $n = 1,507$). We also evaluated associations for the top 4 SNPs using data from 11,870 CRC cases and 14,190 controls of European ancestry included in GECCO and CCFR, which included 14 studies from the United States, Europe, Canada and Australia[4,20,21]. Approval was granted from the relevant institutional review boards at all study sites, and all included participants gave informed consent.

**Genotyping and quality control procedures.** Detailed descriptions of genotyping and quality control procedures as well as design of plates and control samples are given in the **Supplementary Note**. Briefly, in stage 1, 481 cases and 2,632 controls from Shanghai-1 were genotyped using the Affymetrix Genome-Wide Human SNP Array 6.0 as described previously[14]. The average concordance percentage of quality control samples was 99.7%, with a median value of 100% in Shanghai-1 (refs. 14,37,38). Stage 1 genotyping for 296 cases and 257 controls in Shanghai-2 was performed using Illumina HumanOmniExpress BeadChips. The same method was used to genotype cases from the Guangzhou-1 ($n = 694$) and Aichi-1 ($n = 497$) studies in stage 1. The positive quality control samples in these studies had an average concordance percentage of 99.41% and a median value of 99.97%. Cases and controls in KCPS-II were genotyped using the Affymetrix Genome-Wide Human SNP Array 5.0 (ref. 16). Controls for the Guangzhou-1 and Aichi-1 studies were genotyped previously using the Illumina Human610-Quad BeadChip[15] and Illumina Infinium HumanHap610 BeadChip[17] platforms, respectively. Details of quality control procedures for these samples have been described previously[15–17]. We excluded from the analysis samples that were genetically identical or duplicated, had a genotype-determined sex that was inconsistent with self-reported data, had unclear population structure, had close relatives with a PI-HAT estimate greater than 0.25 or had a call rate of <95%. Within each study, SNPs were excluded if (i) MAF was <5%, (ii) the call rate was <95%; (iii) the genotyping concordance percentage was <95% in quality control samples; (iv) the $P$ value for Hardy-Weinberg equilibrium was $<1.0 \times 10^{-5}$ in controls; or (v) SNPs were not on the 22 autosomes. The final numbers of cases, controls and SNPs remaining for analysis in each participating study are presented in **Supplementary Table 1**.

Genotyping for stage 2 was completed using the iPLEX Sequenom MassARRAY platform as described previously[14,39]. With the exception of samples from the Guangzhou-3 study, which were genotyped at Fudan University (Shanghai), all other samples were genotyped at the Vanderbilt Molecular Epidemiology Laboratory. The average concordance percentage of the genotyping data for positive control samples was >99% with a median value of 100% for each of the five studies. SNPs were excluded from the analysis if (i) the call rate was <95%, (ii) the genotyping concordance percentage was <95% in control samples, (iii) there was an unclear genotype call or (iv) the $P$ value for Hardy-Weinberg equilibrium was $<7.8 \times 10^{-4}$. The numbers of SNPs remaining for analysis in each participating study in stage 2 are presented in the **Supplementary Note**.

Genotyping for samples included in the GECCO and CCFR GWAS was conducted using Illumina BeadChip arrays, with the exception of the Ontario Familial Colorectal Cancer Registry study, for which Affymetrix arrays were used[4,20,21]. Details of the quality control procedures for these samples are presented in the **Supplementary Note**.

**SNP selection for replication.** SNPs were selected for stage 2 replication if (i) data were available in each of the five stage 1 studies; (ii) MAF was >5% in each stage 1 study; (iii) no heterogeneity was detected across the five studies included in stage 1 ($P_{het} > 0.05$ and $I^2 < 25\%$); (iv) there was no LD ($r^2 < 0.2$) with any known risk variant reported from previous GWAS; (v) there was no LD ($r^2 < 0.2$) with the other SNPs identified in this study; (vi) there was high imputation quality in each of the five studies (RSQ > 0.5); and (vii) $P < 0.01$ in combined analysis of all stage 1 studies.

**Evaluation of population structure.** We evaluated population structure in each of the five participating studies included in stage 1 by using principal-components analysis (PCA). Genotyping data for uncorrelated genome-wide SNPs were pooled with data from HapMap to generate the first ten principal components using EIGENSTRAT software[40] (see URLs). The first two principal components for each sample were plotted using R (see URLs). We identified and excluded one participant of KCPS-II who was more than 6 s.d. away from the means of principal components 1 and 2 (**Supplementary Fig. 1**). The remaining 7,847 samples showed clear east Asian origin, and these samples were included in the final genome-wide association analysis. Cases and controls in each of the five studies were in the same cluster as HapMap Asian samples. The estimated inflation factor $\lambda$ ranged from 1.02 to 1.04 in these studies after adjusting for age, sex and the first ten principal components, with a $\lambda$ of 1.01 for combined stage 1 data (**Supplementary Fig. 2** and **Supplementary Table 1**).

**Imputation.** We used the MaCH 1.0 program[18] (see URLs) to impute genotypes for autosomal SNPs that were present in HapMap Phase 2 release 22 separately for each of the five studies included in stage 1. Genotype data from the 90 Asian subjects from HapMap were used as the reference. For Guangzhou-1 and Aichi-1, cases and controls were genotyped using different platforms. To improve imputation quality[41], we identified SNPs for which data were available in both cases and controls (250,612 SNPs in Guangzhou-1 and 232,426 SNPs in Aichi-1) and used them to impute genotyping data. A total of 1,636,380 genotyped SNPs or imputed SNPs with high imputation quality (RSQ > 0.50) in all five studies were tested for association with CRC. To directly evaluate the imputation quality for the top four SNPs identified in our study, we genotyped them in approximately 2,500 samples included in stage 1. The agreement of genotype calls derived from direct genotyping and imputation was very high, with mean concordance rates of 98.05%, 95.61%, 99.84% and 97.90% for rs647161, rs10774214, rs2423279 and rs1665650, respectively (**Supplementary Table 7**).

**Statistical analyses.** Dosage data for genotyped and imputed SNPs for participants in each stage 1 study were analyzed using the program mach2dat[18] (see URLs). We coded 0, 1 or 2 copies of the effect allele as the dosage for genotyped SNPs, and, for imputed SNPs, we used the expected number of copies of the effect allele as the dosage score. This approach has been shown to give unbiased estimates in meta-analyses[42]. Associations between SNPs and CRC risk were assessed using ORs and 95% CIs derived from logistic regression models. ORs were estimated on the basis of the log-additive model and adjusted for age, sex and the first ten principal components. PLINK version 1.07 (see URLs) also was used to analyze genotype data[43] and yielded results virtually identical to those derived from dosage data using mach2dat[18]. Meta-analyses were performed using the inverse-variance method, assuming a fixed-effects model, and calculations were implemented in the METAL package[19] (see URLs).

Similar to stage 1, we used logistic regression models to derive ORs and 95% CIs for the 64 selected SNPs in stage 2, assuming a log-additive model with adjustment for age and sex. We performed joint analyses to generate summary results for combined samples from all studies, with additional adjustment for study site. We also conducted stratification analysis for the top four SNPs by population ancestry (Chinese, Korean or Japanese) and by sex. We used Cochran's $Q$ statistic to test for heterogeneity[44] and the $I^2$ statistic to quantify heterogeneity[45] across studies as described elsewhere in detail[46]. Analyses for stage 2, as well as combined stage 1 and 2 data, were conducted using SAS, version 9.2 (see URLs), with the use of two-tailed tests. $P$ values of $<5 \times 10^{-8}$ in the combined analysis was considered statistically significant.

We used Haploview version 4.2 (see URLs; ref. 47) to generate a genome-wide Manhattan plot for results from the stage 1 meta-analysis. Forest plots