

# UCSF

## UC San Francisco Previously Published Works

### Title

Genome, epigenome and RNA sequences of monozygotic twins discordant for multiple sclerosis.

### Permalink

<https://escholarship.org/uc/item/5wz869x8>

### Journal

Nature, 464(7293)

### ISSN

0028-0836

### Authors

Baranzini, Sergio E  
Mudge, Joann  
van Velkinburgh, Jennifer C  
[et al.](#)

### Publication Date

2010-04-01

### DOI

10.1038/nature08990

Peer reviewed



Published in final edited form as:

*Nature*. 2010 April 29; 464(7293): 1351–1356. doi:10.1038/nature08990.

## Genome, epigenome and RNA sequences of monozygotic twins discordant for multiple sclerosis

Sergio E. Baranzini<sup>1</sup>, Joann Mudge<sup>2</sup>, Jennifer C. van Velkinburgh<sup>2</sup>, Pouya Khankhanian<sup>1</sup>, Irina Khrebtukova<sup>3</sup>, Neil A. Miller<sup>2</sup>, Lu Zhang<sup>3</sup>, Andrew D. Farmer<sup>2</sup>, Callum J. Bell<sup>2</sup>, Ryan W. Kim<sup>2</sup>, Greg D. May<sup>2</sup>, Jimmy E. Woodward<sup>2</sup>, Stacy J. Caillier<sup>1</sup>, Joseph P. McElroy<sup>1</sup>, Refujia Gomez<sup>1</sup>, Marcelo J. Pando<sup>4</sup>, Leonda E. Clendenen<sup>2</sup>, Elena E. Ganusova<sup>2</sup>, Faye D. Schilkey<sup>2</sup>, Thiru Ramaraj<sup>2</sup>, Omar A. Khan<sup>5</sup>, Jim J. Huntley<sup>3</sup>, Shujun Luo<sup>3</sup>, Pui-yan Kwok<sup>6,7</sup>, Thomas D. Wu<sup>8</sup>, Gary P. Schroth<sup>3</sup>, Jorge R. Oksenberg<sup>1,7</sup>, Stephen L. Hauser<sup>1,7</sup>, and Stephen F. Kingsmore<sup>2</sup>

<sup>1</sup>Department of Neurology, University of California at San Francisco, San Francisco, CA 94143, USA

<sup>2</sup>National Center for Genome Resources, Santa Fe, NM 87505, USA

<sup>3</sup>Illumina Inc., Hayward, CA 94545, USA

<sup>4</sup>Stanford Medical School Blood Center, Palo Alto, CA 94303, USA

<sup>5</sup>Department of Neurology, Wayne State Medical School, Detroit, MI 48201, USA

<sup>6</sup>Cardiovascular Research Institute, University of California at San Francisco, San Francisco, CA 94143, USA

<sup>7</sup>Institute for Human Genetics, University of California at San Francisco, San Francisco, CA 94143, USA

<sup>8</sup>Department of Bioinformatics, Genentech Inc., South San Francisco, California 94080, USA

### Abstract

Monozygotic (MZ) or “identical” twins have been widely studied to dissect the relative contributions of genetics and environment in human diseases. In multiple sclerosis (MS), an autoimmune demyelinating disease and common cause of neurodegeneration and disability in young adults, disease discordance in MZ twins has been interpreted to indicate environmental importance in its pathogenesis<sup>1–8</sup>. However, genetic and epigenetic differences between MZ

Users may view, print, copy, download and text and data- mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use: [http://www.nature.com/authors/editorial\\_policies/license.html#terms](http://www.nature.com/authors/editorial_policies/license.html#terms)

Correspondence should be addressed to S.E.B. ([sebaran@cgl.ucsf.edu](mailto:sebaran@cgl.ucsf.edu)) or S.F.K. ([sfk@ncgr.org](mailto:sfk@ncgr.org)).

**Full Methods** are available in the online version of the paper at [www.nature.com/nature](http://www.nature.com/nature).

**Supplementary Information** is linked to the online version of the paper at [www.nature.com/nature](http://www.nature.com/nature).

**Author Contributions** S.E.B., G.P.S., J.R.O., S.L.H and S.F.K. designed the project. S.F.K., S.E.B., J.M. and J.R.O. wrote the paper with input from the other authors. S.E.B., J.M., J.C.V., L.Z., R.W.K., G.D.M., J.E.W., S.J.C., J.P.M., R.G., M.J.P., L.E.C., E.E.G., F.D.S., J.J.H. and S.L. performed the experiments. S.E.B., J.M., J.C.V., P.K., I.K., N.A.M., L.Z., A.D.F., C.J.B., T.R., S.L., P.K., T.D.W., G.P.S., J.R.O., S.L.H. and S.F.K. analyzed the data. S.L.H., J.R.O. and O.A.K. supervised patient recruitment.

Data is deposited at dbGaP under accession phs000239.v1.p1. Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints).

twins have been described, challenging the accepted experimental paradigm in disambiguating effects of nature and nurture.<sup>9–12</sup> Here, we report the genome sequences of one MS-discordant MZ twin pair and messenger RNA (mRNA) transcriptome and epigenome sequences of CD4<sup>+</sup> lymphocytes from three MS-discordant, MZ twin pairs. No reproducible differences were detected between co-twins among ~3.6 million single nucleotide polymorphisms (SNPs) or ~0.2 million insertion-deletion polymorphisms (indels). Nor were any reproducible differences observed between siblings of the three twin pairs in HLA haplotypes, confirmed MS-susceptibility SNPs, copy number variations, mRNA and genomic SNP and indel genotypes, or expression of ~19,000 genes in CD4<sup>+</sup> T cells. Only two to 176 differences in methylation of ~2 million CpG dinucleotides were detected between siblings of the three twin pairs, in contrast to ~800 methylation differences between T cells of unrelated individuals and several thousand differences between tissues or normal and cancerous tissues. In the first systematic effort to estimate sequence variation among MZ co-twins, we did not find evidence for genetic, epigenetic or transcriptome differences that explained disease discordance. These are the first female, twin and autoimmune disease individual genome sequences reported.

---

We sought to assess the magnitude of genetic, epigenetic and transcriptomic differences in CD4<sup>+</sup> lymphocytes from MS-affected and unaffected MZ twin sibships (Supp. Fig. 1). CD4<sup>+</sup> T cells are involved in the pathophysiology of MS (MIM 126200).<sup>1</sup> mRNA, genomic DNA (gDNA) and reduced-representation, bisulphite-treated gDNA were prepared from negatively isolated, CD4<sup>+</sup> T lymphocytes from three pairs of adult, MZ twins who were discordant for MS ('001, affected; '101, unaffected). Affected individuals fulfilled McDonald criteria for MS diagnosis.<sup>13</sup> Lack of sibling affection was assessed through clinical evaluation, and, for twin 041896-101, confirmed by magnetic resonance brain imaging, and cerebrospinal studies. MZ twin pair 041896 was female, of Ashkenazi Jewish origin and beyond the susceptibility age-range for MS at time of study (Supp. Table 1). Twin pair 230178 was female and African American, while twins 041907 were white males. Individual 041896-001 had onset of MS at age 30 years, and is currently in the secondary progressive phase; individuals 230178-001 and 041907-001 had MS onset at ages 38 and 13, respectively, and have relapsing-remitting disease. Molecular typing of HLA loci showed identical genotypes within the three twin pairs (Supp. Table 1). Only co-twins 041907 had DRB1\*1501, the strongest genetic susceptibility factor for MS.<sup>14</sup>

Nucleic acid samples were sequenced by sequencing-by-synthesis with reversible-terminator chemistry<sup>15–18</sup>. mRNA was prepared from blood samples drawn on different days from twin pair 041896 to ascertain sampling variance. Fifty to 68 million, high-quality, 36–44 nt, singleton sequences from each of eight mRNA samples were aligned to the NCBI human genome reference and read-counts-per-gene were calculated<sup>18–20</sup> (Supp. Table 2). Sequencing to this depth (median relative transcript coverage of 5.0-fold and 6.4-fold for 041896-001 and 041896-101, respectively) allowed determination of diversity of the polyadenylated transcriptome in CD4<sup>+</sup> lymphocytes: ~92% of 20,601 genes with exon annotations were expressed, as assessed both by aligned reads and the upper asymptote of the best-fit sigmoid curve (Supp. Table 2, Supp. Fig. 2). The distribution of transcript abundance was a left-skewed, bell-shaped curve with >7 log<sub>10</sub> dynamic range (Supp. Fig. 2), in agreement with a previous study<sup>17</sup>. Digital gene expression values correlated well with

exon-resolution array hybridization results (Supp. Fig. 3), in agreement with another report<sup>21</sup>. Surprisingly, diagnosis or treatment of MS accounted for only 9.4% of variance in transcript abundance in T cells of MZ twins, compared with 57.3% attributable to twin pair-to-twin pair differences, 23.6% to day-to-day variation (as assessed in twin pair 041896 alone) and 3.5% to sequencing lane-to-lane variation, Supp. Figs. 4–7). The variance in transcript abundance attributable to MS was within the range of variances obtained by random permutation of MS diagnosis labels (Supp. Fig. 8, Supp. Table 3). Thus, robust gene expression differences were not observed between MS-affected and unaffected twins in CD4<sup>+</sup> lymphocytes that were inexplicable by other variables.

One billion, high-quality, shotgun, whole genome sequences were generated from twins 041896-001 and -101, corresponding to 21.7- and 22.5-fold aligned coverage, and representing 99.6% and 99.5% of the NCBI human reference genome, respectively (Suppl. Table 4). Comparisons of genome coverage of the twins with the AK1 genome, which was determined using identical procedures, revealed no individual coverage bias<sup>15</sup> (Supp. Fig. 9 and 10).

Viral infection has been suggested to contribute to the etiology of MS. Upon re-alignment of unaligned sequences to 2864 viral genomes, ~0.02% of DNA reads from twins 041896 and 0.2% of RNA reads from the three twin pairs matched 310 viral genomes. A large majority of these alignments reflected simple sequence repeats or endogenous retroviral sequences. Upon reverse-transcription and PCR, no reproducible differences were found between sequences aligning to viral genomes in T cells from MS-affected and unaffected individuals.

Approximately 3.6 million SNPs and ~0.2 million indels were detected in subject 041896-001 and -101 genomes, using optimized criteria, which are similar to values reported for male genomes (15 and references therein; Supp. Table 5). Indels varied in size from –31 to +8 nt, with an approximately normal frequency distribution. Of 13 common risk variants previously associated with MS susceptibility<sup>14</sup>, co-twins 041896 were homozygous for five, heterozygous for five, and three were absent. This genetic load is predicted to increase risk for development of MS ~8-fold under an additive model (Supp. Table 6). Co-twins 230178 were homozygous for seven susceptibility loci and heterozygous for two, and co-twins 041907 were homozygous for eight risk alleles and heterozygous for two, conferring 14-fold and 43-fold increased risk, respectively (Supp. Table 6). These data should be interpreted cautiously since translation of genetic burden into risk for complex disorders is rudimentary. Clustering of 9.9 million SNPs in eight individual genome sequences showed close similarity of the twin 041896, female genomes and their separation from six male genomes (Supp. Fig. 11).

SNP genotype differences were sought between affected and unaffected twin siblings in genomic DNA and mRNA (Supp. Fig. 1). Firstly, stringent bioinformatic filters were trained both to call SNPs in aligned genome and mRNA sequences and to infer SNP genotypes, by comparing genotypes obtained from duplicate Affymetrix 6.0 SNP array hybridizations with those derived from genome and mRNA sequencing (Supp. Tables 7 and 8; Supp. Fig. 12)<sup>15</sup>. These filters excluded low coverage or repetitive genomic sequences (<11; –fold or >44-fold coverage, respectively), yielding high positive predictive values (PPV) to enable meaningful

co-twin comparisons. Secondly, these filters were used to determine SNP genotypes in aligned genomic sequences of twin pair 041896 and in aligned mRNA sequences of the three twin pairs. Thirdly, identities and differences in inferred SNP genotypes were sought between affected and unaffected twin siblings. Co-twin genotype differences were categorized either as changes from homozygous reference allele to heterozygote, or heterozygote to homozygous variant (Table 1). Of 1,089,550 SNP genotypes inferred in genomes 041896 using these filters, 3,241 (0.3%) differed between twins (Table 1). Of over 730,000 genomic SNP genotypes determined by duplicate array hybridizations, 126 (0.02%), 153(0.02%), and 120 (0.02%) differed between siblings in the three twin pairs, respectively, which was considerably less than ~8,500 SNPs that were discordant between repeated hybridizations of individual DNA samples (Supp. Table 9). mRNA sequencing covered ~65.6 Mb of annotated exons to a depth of ~5-fold. Three hundred and twenty two (0.6%), 1,017 and 380 SNP genotypes inferred in mRNA sequences differed between siblings of twin pairs 041896, 230178 and 041907, respectively (Table 1). Finally, replication of co-twin SNP genotype identities and differences was sought. No differences in SNP genotypes inferred by one approach were recapitulated by a second method. In contrast, >98% of SNPs that were identical in twin siblings and genotyped by two methods (array hybridization, mRNA sequencing or genomic DNA sequencing) were replicated (Table 1). Furthermore, Sanger resequencing revealed identical genotypes in twin pair 041896 for a set of 15 SNP differences well supported by at least one method.

The SNP genotyping filters were also used to infer indel genotypes in genome and mRNA sequences of the twins: 91.9% of indels detected in both genome and mRNA sequences had identical genotypes (Table 1). Of 26,908 indel genotypes inferred in the genomes of twins 041896, 213 (0.8%) differed between siblings. Of 1,322, 1,073 and 407 indel genotypes inferred in mRNA sequences from twins 041896, 230178 and 041907, 8, 39 and 10 differed between twin siblings, respectively (Table 1). No indel genotype differences identified by one approach were recapitulated by a second method. In summary, siblings in three MZ twin pairs exhibited no replicable nucleotide variation differences in non-repetitive sequences, as assessed by genome and mRNA sequencing and SNP array hybridization. Much longer reads and lower error rates will be required to evaluate variation differences in repetitive sequences comprehensively. Detection of no replicable SNP genotype differences between siblings of any of the three twin pairs in peripheral CD4<sup>+</sup> T cells accords with estimated rates of somatic mutation of  $8.4 \times 10^{-9}$  to  $4.6 \times 10^{-10}$  per nt per generation in human tumors, *Saccharomyces cerevisiae* and *Drosophila melanogaster*<sup>22–24</sup>.

Expression quantitative trait loci (eQTL) are emerging as a molecular mechanism for common SNPs that are significant in genome-wide association studies of disease<sup>25</sup>. In light of an absence of significant MS-associated genotypic or mRNA expression differences between twins, we sought allele specific differences in mRNA expression. For heterozygous coding SNPs (cSNPs), the expression of both alleles in CD4<sup>+</sup> lymphocytes was measured to address deviation from the 1:1 expected ratio (allelic imbalance). 268 heterozygous cSNPs exhibited allelic imbalance *in cis* at 188 loci in twin 041896 transcriptomes, as determined by significant deviation of aligned genomic and mRNA read counts (Supp. Table 10). Single base mismatches do not cause systematic bias in GSNAP alignments. Two imprinted genes

showed altered allelic expression in both co-twins (ZNF331 and GNAS) as did three genes that exhibit altered allelic expression in human cerebellar cortex (ABLM1, UBE2I and KIAA1267, Kingsmore et al., unpublished) and two that previously have shown altered allelic expression in CD4+ lymphocytes (the MS-associated gene CD6 and acid trehalase-like 1 [ATHL1])<sup>14,26</sup>. We used quantitative PCR to validate each of the three possible outcomes: a) where both twins showed an expected 1:1 ratio of allelic expression; b) where both twins show skewed expression of an allele in the same direction and magnitude, indicative of cis-acting eQTL or imprinting, and; c) where the direction or magnitude of the imbalance differed between the twins (Supp. Fig. 13). Interestingly, 115 (43%) cSNPs differed between twins (i.e. differential allelic expression; Fig. 1, Supp. Table 10). These results suggest that some gene expression differences between twins represent chromatid-specific alterations in transcription. Variance in allelic expression between samples mirrored that observed in overall mRNA levels, with twin pair-to-twin pair accounting for 51.2%, day-to-day variation for 27.7% and MS diagnosis for 8.0% of variance. No cSNPs showing allelic imbalance were shared among the three twin pairs. Interestingly, however, cSNPs that show allelic imbalance were significantly closer to transcription factor binding sites than random SNPs, providing a novel, potential mechanism of action.

Structural variants were identified in the six genomes by hybridization of duplicate arrays. In contrast to a recent report, we found no copy number variants or allelic gains / losses that differed between siblings in any twin pair<sup>12</sup>. Twins 041896 displayed 143 structural variants comprising 89 Mb, twins 230178 exhibited 13 variants comprising 3 Mb and twins 041907 had 58 variants encompassing 33 Mb (Supp. Fig. 1, 14, 15, Suppl. Table 11). Of note, seven structural variants were common to all three twin pairs, and changed the copy number of two genes (Late Cornified Envelope-3B (LCE3B) and T Cell Receptor Gamma Chain Alternate Reading Frame Protein (TARP)) and one pseudogene (A Disintegrin And Metalloprotease-6 (ADAM6), Supp. Table 12). LCE3B was not expressed in T cell mRNA samples from these patients. TARP was expressed at a level of  $12.9 \pm 6.1$  reads/million (mean  $\pm$  standard deviation) and did not show altered expression in MS. These genes have not previously been associated with MS.

An additional axis of heritable genetic information in human genomic DNA is cytosine methylation, which serves several functions including regulation of gene expression, silencing of retrotransposons, genomic imprinting and X-chromosome inactivation and has been implicated in several diseases<sup>27,28</sup>. We sought to compare genome-scale DNA methylation profiles between twin siblings at nucleotide resolution. We aligned 50 – 90 million, high-quality, 50 nt, reduced representation bisulphite sequences (RRBS) from ten samples – the three pairs of twin T lymphocytes, normal lung and lung cancer, and normal breast and breast cancer<sup>16</sup> (Supp. Table 13). For twins 041896, these corresponded to 45.5- and 32.7-fold coverage of 1.4 million uniquely aligning, non-repetitive *MspI* fragments, and 2,146,620 and 2,033,078 CpG dinucleotides from -001 and -101 genomes, respectively (Table 2). Bisulphite conversion of non-CpG cytosines was >99%. Almost identical numbers of CpG sites were identified in the forward and reverse strands, as expected (Supp. Table 14). As reported for mouse, methylation levels of CpG dinucleotides in human T cells displayed a bimodal distribution, with most being unmethylated or methylated (>95% of



reads in either state) [Fig. 2a, Supp. Fig. 16]16. Approximately one quarter of CpGs were methylated. Over 90% of CpG sites were common to siblings within each twin pair (Table 2). CpGs aggregated into clusters (corresponding to CpG islands)16 at a ratio of 1.58 – 1.74 CpGs per cluster. Over 92% of CpG clusters were common to siblings within each twin pair (Table 2 and Supp. Table 14). Highly congruent results were obtained with two alignment algorithms (Suppl. Table 14; Suppl. Figs. 17 and 18) and two reference genome datasets. Of ~2 million CpGs represented by  $\geq 10$  high-quality reads in twins 041896, only two showed a switch between siblings from  $\leq 20\%$  methylated to  $\geq 80\%$  by ELAND and four by GSNAP, none of which was supported by both methods (Fig. 2b, Table 2). Likewise, 10 of 1.7 million CpG sites in twins 230178 and 176 of 1.7 million CpG sites in twins 041907 showed a switch in methylation by ELAND (Fig. 2c,d, Supp. Table 15). Two CpG methylation switches between affected and unaffected siblings were common to twin pairs 230178 and 41907, albeit with opposite directions of change ( $>80\% \rightarrow <20\%$  mCpG in 041907-001 and -101, respectively, whereas  $<20\% \rightarrow >80\%$  mCpG in 230178-001 and -101, at a CpG site 9912 nt 5' of TMEM1 and 8536 and 10,659 nt 5' of PEX14). To put these findings in context, we evaluated the magnitude of methylation changes in CD4<sup>+</sup> T cells from unrelated individuals, between tissues and between normal and cancerous tissue. 586 – 827 inter-individual  $<20\% \rightarrow >80\%$  CpG methylation differences were observed (Fig. 2e,f); 4255 – 7180 CpG methylation shifts were observed between T lymphocytes, lung and breast tissues (Table 2, Fig. 2i,j). Breast and lung cancers revealed 1,557 and 16,509 CpG methylation shifts, when compared with normal breast and lung tissue, respectively (Fig. 2g,h, Table 2). A second pattern of change in CpG methylation was observed in comparison of male and female samples: 394 CpGs were  $<5\%$  methylated in 041907-001 T lymphocytes (male) but 20 – 50% methylated in 041896-001 (female). Likewise, 406 CpG sites were  $<5\%$  methylated in 041907-101 (male) and 20 – 50% methylated in 041896-101 (female). Of these, 385 and 389, respectively, mapped to Chromosome X, consistent with female X inactivation (Fig. 2e). Similarly, a very large number of CpG sites that were  $<10\%$  methylated in normal lung were 20–70% methylated in lung cancer (Fig. 2h). A previous study has shown epigenetic differences between dizygotic twins to be qualitatively greater than between MZ twins<sup>29</sup>. Here we show the magnitude of epigenetic differences between MZ twin sibling CD4<sup>+</sup> lymphocytes to be at least an order of magnitude less than those between individuals and ~3 orders less than those observed between tissues and in malignant transformation.

In summary, the recent GWAS-identification of novel risk loci is opening a broad window into genetic intricacies underpinning complex diseases. Although genetic knowledge remains incomplete, a new generation of sequencing and analytical tools may prove to hold great potential, as shown herein. Likewise, a discordant MZ twins study controls for many genetic and non-genetic confounders, enhancing the tractability of mechanisms in complex disorders. We sought genetic, epigenetic or transcriptomic differences between CD4<sup>+</sup> T cells of twin siblings that might explain MS-discordance. While MS is a neurologic disease, T cells are fundamentally involved in its pathophysiology<sup>1</sup>. However, no reproducible differences in SNPs, indels, CNVs, gene expression levels or sequences aligning to viral genomes were detected between CD4<sup>+</sup> T cells of co-twins. To provide analytical rigor, SNP and indel differences were sought using at least two different approaches and CNV

experiments were performed in duplicate. However, analysis of nucleotide variants was limited in scope by exclusion of low coverage regions and repetitive sequences (since the latter cannot be reliably interrogated by alignment of short reads or array hybridization), by moderate sensitivity for detection of structural variants of size 50 – 1500 nt (which fall between the resolution of sequencing and array hybridization), and limited feasibility to detect possible somatic mosaicism. A previous study has shown differences in selection of T-cell receptors after antigen stimulation between MZ twins discordant for MS30. Quantitative analysis of T-cell repertoire or immunoglobulin locus recombination was not possible at ~22X depth of aligned coverage. Progress in single molecule sequencing technologies with longer reads and deeper coverage should overcome many of these limitations in the future, as would examination of additional cellular compartments of innate and adaptive immunity. Additionally, deep RRBS revealed very few changes in CpG methylation between CD4<sup>+</sup> T cells of twin siblings and no differences common to two or more twin pairs. It should be noted, however, that RRBS was limited to the investigation of dramatic shifts in CpG methylation in a relatively broad population of T cells. Other epigenetic mechanisms, differences within lymphocyte subsets, mono-allelic differences or other tissues were not examined. These caveats aside, however, MZ twins lacked genetic, epigenetic or transcriptomic differences in T cells to explain MS-discordance. A number of tantalizing, novel, differences were detected that will require replication and additional studies: 43% of eQTLs had a different direction or magnitude of imbalance in twin siblings. In summary, a singular genetic, epigenetic or transcriptomic mechanism underpinning MS-discordance in MZ-twins was not detected in a study of unprecedented resolution. While disease discordant MZ twins appear to provide a framework for analysis of complex disorders that has fewer variables, additional stratification and/or concomitant measurement of multiple data types may be necessary to yield molecular mechanisms underpinning disease.

## METHODS SUMMARY

The study was approved by the UCSF Institutional Review Board. Informed, written consent was obtained from all individuals. CD4<sup>+</sup> lymphocytes were isolated from peripheral blood and nucleic acids extracted with standard methods. Two samples were obtained on different days from twins 041896 and single samples from the others. HLA typing was by AlleleSEQR (Atria Genetics) and Assign SBT software (Conexio Genomics). Genome-wide genotypes and CNVs were detected with Affymetrix 6.0 arrays in duplicate. Log-R ratios were generated with Affymetrix Genotyping Console 3.0.2 and analyzed with Nexus software (BioDiscovery Inc., El Segundo, CA). Short- and long-insert, paired-end libraries were generated from gDNA, mRNA and reduced-representation, bisulphite-treated gDNA as described<sup>15–18</sup>. Paired-end and singleton, 36–130 nt reads were generated using Illumina GAIIx instruments. Sequences were aligned principally to NCBI reference genome build 36.3, with GSNAP and tolerance of 5% mismatches<sup>15</sup>. SNPs, indels and gene expression were analyzed with Alpheus using filters trained with array results<sup>15,18–20</sup>: Genomic SNP calling filters were >20% and >4 uniquely aligning reads with average quality score (Q) ≥20 (Supp. Table 7). mRNA SNP calling filters were Q ≥20, presence in ≥20% and ≥2 reads and ≥1 uniquely aligning read. Nucleotides with coverage 11–44X and Q ≥20 were genotyped



according to frequency cutoffs in Supp. Table 8; Genotype differences were called where frequencies differed by >50%. eQTLs were detected by allelic mRNA read counts differing from equality with  $\chi^2$  p-values of  $<10^{-7}$ . Gene expression was assessed by  $\log_2$ -transformed aligned read counts. Putative SNP differences were validated by Sanger sequencing and putative gene expression differences using Affymetrix Human Exon 1.0 ST arrays. Putative eQTLs and virus alignments were validated by quantitative PCR (with allele specificity for the former). Statistical analysis used JMP-Genomics (SAS Institute, Cary, NC) or R (<http://www.R-project.org>).

## Methods

### Array-based Genotyping and CNV Detection

Genome-wide genotypes (>900,000 SNPs) and CNVs (~1.8 million probes) were detected with Affymetrix 6.0 arrays (Santa Clara, CA). Genomic DNA from each individual was tested on duplicate arrays. Log-R ratios (normalized probe intensities) were generated with Affymetrix Genotyping Console 3.0.2 and analyzed with Nexus software (BioDiscovery Inc., El Segundo, CA), which identifies CNVs with a circular binary algorithm using intensity data from all probes, and allele ratios from SNP probes.

### Alignment of mRNA and gDNA Sequences to Reference Databases

mRNA-Seq and whole genome shotgun sequences were aligned to the NCBI reference genome (build 36.3) with GSNAP and tolerance of 5% mismatches<sup>15,20</sup> (Supp. Table 2 and 4). For definition of exon boundaries, annotations from RefSeq Transcript (downloaded 9/2/2008) and from 5,224 non-redundant UCSC transcripts (downloaded 4/13/2009) were appended to Build 36.3 of the reference human genome. Long (75 – 130 nucleotides) genomic reads were found to align poorly using these criteria, due to low terminal quality scores and higher rates of mismatch. Therefore, unaligned long, genomic, paired reads were further aligned to the NCBI reference genome with GSNAP by trimming to paired 75 nucleotides (nt) and tolerance of  $\leq 10$  mismatches.

mRNA-Seq and whole genome shotgun reads not mapping to the human genome were aligned to 2,864 NCBI viral genome sequences (release 35) with GSNAP and tolerance of 5% mismatches. Alignments were visualized using Alpheus<sup>20</sup> and CMTV<sup>31</sup>. High likelihood true alignments were identified on the basis of:

1. Significant read coverage of the viral genome;
2. Elimination of reads composed primarily of simple sequence repeats;
3. Unique read alignments;
4. Paired read alignments with correct orientation and distance separating read pairs;
5. Alignments of non-clonal reads to contiguous stretches of viral genome sequence.

Putative, novel viral sequences with average quality scores (Q)  $\geq 20$  were assembled by ABySS<sup>32</sup> or by reference-guided assembly with AMOScmp-shortReads-alignmentTrimmed<sup>33</sup>. Default parameters were used. Contigs were aligned to the NCBI nr database using BLASTN 2.2.21.

## mRNA-Seq Based Measurement of Gene Expression Changes

Upon alignment of mRNA-Seq reads, read counts were calculated per gene for each lane of sequence and  $\log_2$  transformed. Distribution analysis (Supp. Fig. 4) and Mahalanobis differences (Supp. Fig. 6) were assessed for log-transformed read counts from each lane of mRNA-Seq and outlier lanes were removed. Principal component analysis (Supp. Fig. 6) and variance decomposition of principal components was undertaken for log-transformed read counts from each lane to assess sources of variability in gene expression (Supp. Fig. 7). Since Diagnosis (MS-affected versus non-affected) accounted for 9.4% of variance, all possible permutations of lanes of sequence were examined to determine whether Diagnosis-associated variance was greater than a random permutation (experimental design file in Supp. Table 3). Principal component analysis and variance decomposition of principal components were repeated with log-transformed read counts from each lane for each permutation to assess permuted diagnosis-associated variance in gene expression (Supp. Fig. 8). Since true Diagnosis-associated variance was not greater than permuted variance, genes differing between MS-affected and unaffected individuals were not assessed by weighted ANOVA.

## Eland Alignment of RRBS

Treatment of DNA with bisulfite converts cytosine residues to uracil, but leaves 5-methylcytosine residues unaffected. Thus, alignments of 50 bp, singleton, reduced representation bisulphate sequences (RRBS) to the human genome are complicated by the simplification of the genetic code from four to three bases, except at methylcytosine (mC) locations. Eland-extended performs alignments of the first 32 nt of a read with up to two substitutions, and then extends the alignment with unlimited mismatches. Alignment of 3-base reads (following conversion of residual cytosines to thymidines in the RRBS reads) to a 3-base genome (following conversion of all cytosines to thymidines) with Eland-extended resulted in many non-unique alignments. In order to circumvent this problem, we made use of the fact that all RRBS start at an *MspI* site (which comprise the majority of CpG residues and large majority of CpG islands<sup>16</sup>). Thus, 3-base reads were aligned to a 3-base version of ~3.7% of the human genome, comprising 2.3 million *MspI* fragments of up to 50 nt in length, derived from the NCBI human genome sequence, version 36.3, totaling 113 Mb in length (Supp. Table 16). The fragments were of two types: 133,609 fragments of 30 – 50 bp that were flanked by *MspI* sites on both ends and 2.2 million 50 bp fragments with a 5' flanking *MspI* site (representing genomic *MspI* fragments of greater than 50 bp in length). Only unique alignments with Phred-like scores  $>4$  (greater than 50% likelihood of being correct alignments) and only those starting with a 5' thymidine (base 1 of a converted *MspI* fragment) were retained (Supp. Table 13). Alignments to fragments of less than 50 nt terminated at the end of the fragment. Eland does not align to *MspI* fragments of less than 30 nt in length. Following alignment of converted reads, thymidine residues were corrected to their original sequence in the RRBS and reference, and C-to-T transitions were identified. Percent methylation for CpG sites was scored by the ratio of C/(C+T) calls for each C that was followed by a G. Percent conversion of C to T when followed by another base was used for estimation of bisulfite conversion rate, and was  $> 99.8\%$ .

## RRBS Alignment with GSNAP

RRBS were also aligned with GSNAP to the NCBI human genome reference sequence, version 36.3, allowing 5% mismatches and without penalizing C-to-T transitions (Supp. Table 13). Since GSNAP reports only the best alignments (those with the fewest mismatches) using the entire 50 nt alignment, unique alignments were possible using the entire genome without penalizing C-to-T transitions. % methylation was assessed for CpG sites with at least 10-fold coverage, based on all alignments (i.e. not restricted to unique). Only CpG sites within MspI fragments were considered. For identification of differences between subjects from “largely methylated” to “largely unmethylated”, we sought positions where there was at least 80% cytosine in one subject and less than 20% cytosine in the other.

GSNAP is a short-read alignment program based on GMAP that employs a hash table and a compressed version of the reference genome, which is constructed once for that genome<sup>34</sup>. The reference may include arbitrary contigs (up to 4 billion), so that one may also align to a reference transcriptome, with redundancy allowed among the contigs. The hash table contains the locations of a given 12-mer in the genome, subject to sampling. The sampling step occurs during pre-processing of the genome, so that genomic locations are stored only for every third 12-mer in the genome. Sampling is needed to reduce the memory footprint of the program below 4 gigabytes for a human-sized genome. GSNAP can handle short reads of > 24 or more nt, with each read in the input potentially having a varying length. There is theoretically no upper bound on the length of the query sequence, except that this bound is compiled into GSNAP by default at 200 nt; longer sequences can be handled simply by changing this constant at compile time.

GSNAP has specialized algorithms for identifying exact mappings, one-mismatch mappings, multiple-mismatch mappings, and indel mappings (including a user-specified number of mismatches). Exact mappings are identified by taking the intersection of genomic positions over a spanning set of 12-mers in the query sequence. The spanning set must contain 12-mers in the same phase modulo 3, to account for the sampling used in pre-processing the genome, so the program must test each of the three possible phases. For spanning set members that overhang the ends of the query sequence by 1 or 2 nt, the relevant genomic positions can be obtained by substituting 1 or 2 wildcard nt, respectively, and taking the union of genomic locations in the hash table.

Candidates for one-mismatch mappings are similarly identified by computing an incomplete intersection, in which one 12-mer in the spanning set does not contain the given genomic location. These candidate genomic mappings are then compared against a compressed version of the genome to verify that only one mismatch was present.

Candidates for multiple-mismatch mappings are determined by processing a sorted list of genomic locations from all 12-mers in the query sequence. This sorted list is computed efficiently using a heap-based priority queue. For each candidate genomic location, a floor on the number of mismatches can be computed from the pattern of query positions of the 12-mers that match the genomic location. Candidates with a sufficiently low floor (based either on a user-specified limit or on the best mapping determined so far) are then compared against the compressed genome to determine the actual number of mismatches.

For identifying indel mappings, GSNAP accumulates partial genomic alignments during the multiple-mismatch algorithm, where a partial alignment can be supported by a single 12-mer in the query sequence. These partial alignments are then scanned in genomic order to identify pairs that are sufficiently close to constitute a candidate indel, where the default distances are 30 nt for an insertion and 12 nt for a deletion. These candidate pairs are then compared against the compressed genome to determine the number of mismatches. To identify indels occurring at either end of the query sequence, the program computes floors that exclude the 12-mers on either end. Candidates with a sufficiently low floor are then compared against the compressed genome to identify a possible indel at the end and to count the actual number of mismatches.

Although GSNAP allows repetitive regions of the genome to be masked before building the genomic data structure, in typical usage (as herein) the genome is not pre-masked. Therefore, GSNAP is able to align sequences to redundant regions in the genome, including repetitive regions, and report all such alignments. In default mode (as herein), the program reports only the best alignments, those with the fewest mismatches, although the program also can be run to identify and report suboptimal alignments. GSNAP differs from ELAND in that it processes the reference genome first, constructs a hash table of the genome, and then aligns the short reads to the genome. In contrast, ELAND processes the short reads first, constructs a hash table of the short reads, and then scans the genome to find matches.

### Identification of Optimal Bioinformatic Filters for SNP Detection and Genotyping

SNP detection in Illumina GAI sequences is complicated by relatively high sequencing error rates, particularly at nucleotides 50 – 130 using the chemistry and base calling software available during the first half of 2009. SNP genotyping in Illumina GAI sequences is complicated by a continuous, albeit trimodal, distribution of frequencies of SNP- and reference sequence-containing reads at a given location (Supp. Figure 12). In order to translate SNP- and reference sequence-containing read frequencies into genotypes and to understand the sensitivity and specificity of SNP detection and genotyping, comparisons between array-based SNP genotypes and sequencing results were performed extensively. Unambiguous SNP genotypes from duplicate array hybridizations (with SNP calls and concordant genotypes in both replicates) were assessed to be true. Subsets of SNPs common to Affymetrix 6.0 arrays and sequence datasets were identified. Optimal SNP genotyping filters (those with maximal positive predictive value (PPV) and near-optimal sensitivity) for each sequence dataset were identified by determining the number of true positives, false positives and false negatives and determining the PPV and sensitivity of all combinations of the following criteria: number of reads calling the SNP, number of uniquely aligning reads calling the SNP, % reads calling the SNP, average quality score (Q), and minimum quality score. To detect changes in SNP genotype, each possible genotype in a diploid genome was modeled (homozygous reference allele, heterozygote, and homozygous variant allele) and the optimal change in allele frequency was determined. Resultant filters are shown in Supp. Tables 7 and 8. These methods represent a refinement of those used previously<sup>15</sup>, and which were extensively validated by Sanger resequencing and genotyping arrays.

## Identification of Allele-Specific Expression

Allele-specific expression in mRNA sequences was identified by methods similar to those described<sup>25</sup>. Frequencies of SNP- and reference sequence-containing reads at a given heterozygous location in mRNA sequences are continuous, albeit unimodal (Supp. Figure 12), reflecting both random reference and variant-containing read sequencing, effects of clonal reads and allele-specific expression. Unambiguous heterozygous SNP locations in each individual were determined based on duplicate array hybridizations (with SNP calls and concordant genotypes in both replicates) and by the SNP calling criteria developed above. Allele-specific expression effects were assessed by application of genome-wide p values to significance testing of deviation from 50:50 read frequencies. Artifactual allele-specific expression associated with enrichment of clonal reads was evaluated for many, putative allele-specific expression SNPs by visualization of start and stop sites of reads using Alpheus. Artifactual allele-specific expression associated with bias in GSNAP alignment of reads containing or lacking specific SNPs was evaluated as discussed above.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgements

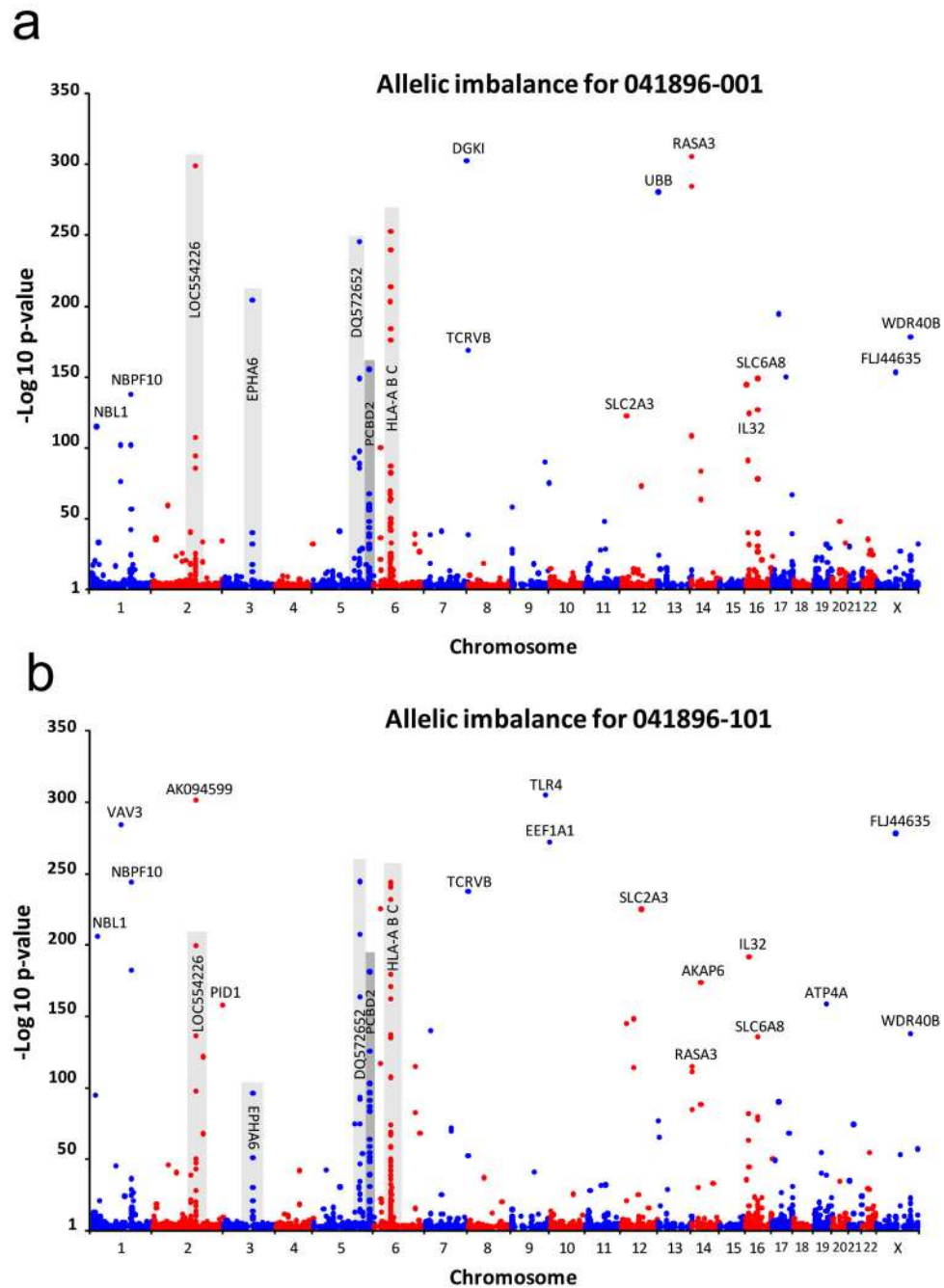
We deeply thank the study subjects. This work was supported by grants from Small Ventures USA Inc., A.J. Brass Foundation, Nancy Davis Foundation, NIH grants RR016480 to F.D.S., RO1NS26799 to S.L.H., RO1NS46297 to J.R.O. and NMSS grants RG3060C8 and RG2901D9 to J.R.O. S.E.B is a Harry Weaver Neuroscience Scholar of NMSS. *A deo lumen, ab amicis auxilium.*

## References

1. Chitnis T. The role of CD4 T cells in the pathogenesis of multiple sclerosis. *Int. Rev. Neurobiol.* 2007; 79:43–72. [PubMed: 17531837]
2. Oksenberg JR, Baranzini SE, Sawcer S, Hauser SL. The genetics of multiple sclerosis: SNPs to pathways to pathogenesis. *Nature Rev. Genet.* 2008; 9:516–526. [PubMed: 18542080]
3. Sadovnick AD, et al. Evidence for genetic basis of multiple sclerosis. *Lancet.* 1996; 347:1728–1730. [PubMed: 8656905]
4. Nielsen NM, et al. Familial risk of multiple sclerosis: A nationwide cohort study. *Am. J. Epidemiol.* 2005; 162:774–778. [PubMed: 16120694]
5. Mumford CJ, et al. The British Isles survey of multiple sclerosis in twins. *Neurology.* 1994; 44:11–15. [PubMed: 8290043]
6. Willer CJ, et al. Twin concordance and sibling recurrence rates in multiple sclerosis. *Proc. Natl Acad. Sci. USA.* 2003; 100:12877–12882. [PubMed: 14569025]
7. Islam T, et al. Differential twin concordance for multiple sclerosis by latitude of birthplace. *Ann. Neurol.* 2006; 60:56–64. [PubMed: 16685699]
8. French Research Group on Multiple Sclerosis. Multiple sclerosis in 54 twinships: Concordance rate is independent of zygosity. *Ann. Neurol.* 1992; 32:724–727. [PubMed: 1471862]
9. Machin GA. Some causes of genotypic and phenotypic discordance in monozygotic twin pairs. *Am. J. Med. Genet.* 1996; 61:216–228. [PubMed: 8741866]
10. Gringras P, Chen W. Mechanisms for differences in monozygous twins. *Early Hum. Dev.* 2001; 64:105–117. [PubMed: 11440823]
11. Fraga MF, et al. Epigenetic differences arise during the lifetime of monozygotic twins. *Proc. Natl Acad. Sci. USA.* 2005; 102:10604–10609. [PubMed: 16009939]

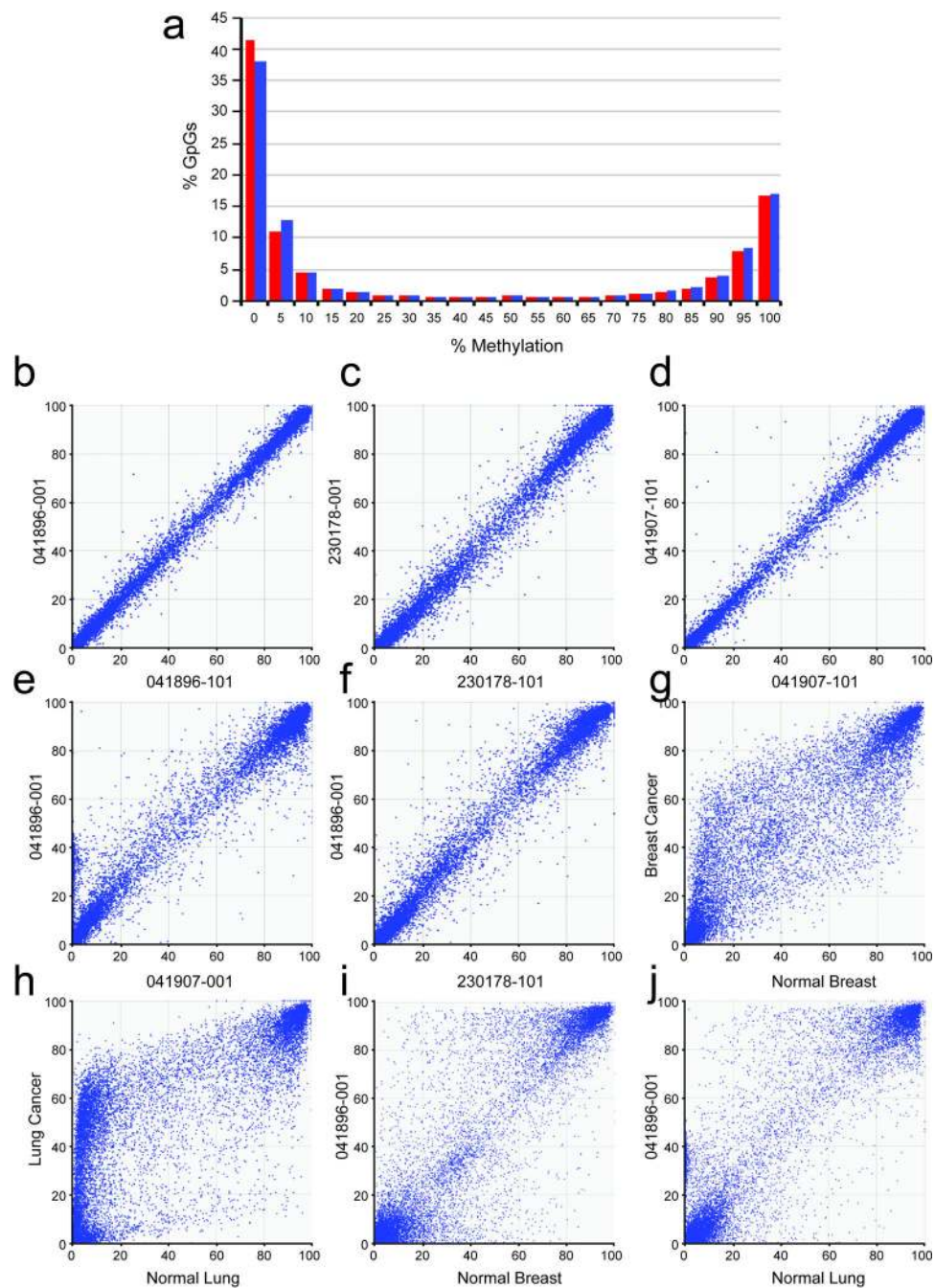
12. Bruder CE, et al. Phenotypically concordant and discordant monozygotic twins display different DNA copy-number-variation profiles. *Am. J. Hum. Genet.* 2008; 82:763–771. [PubMed: 18304490]
13. Kim J-I, et al. A highly annotated whole genome sequence of a Korean Individual. *Nature.* 2009; 460:1011–1015. [PubMed: 19587683]
14. McDonald WI, et al. Recommended diagnostic criteria for multiple sclerosis: guidelines from the International Panel on the diagnosis of multiple sclerosis. *Ann. Neurol.* 2001; 50:121–127. [PubMed: 11456302]
15. De Jager PL, et al. Meta-analysis of genome scans and replication identify CD6, IRF8 and TNFRSF1A as new multiple sclerosis susceptibility loci. *Nature Genet.* 2009; 41:776–782. [PubMed: 19525953]
16. Meissner A, et al. Genome-scale DNA methylation maps of pluripotent and differentiated cells. *Nature.* 2008; 454:766–770. [PubMed: 18600261]
17. Wang ET, et al. Alternative isoform regulation in human tissue transcriptomes. *Nature.* 2008; 456:470–476. [PubMed: 18978772]
18. Mudge J, et al. Genomic convergence analysis of schizophrenia: mRNA sequencing reveals altered synaptic vesicular transport in post-mortem cerebellum. *PLoS ONE.* 2008; 3:e3625. [PubMed: 18985160]
19. Sugarbaker DJ, et al. Transcriptome sequencing of malignant pleural mesothelioma tumors. *Proc. Natl Acad. Sci. USA.* 2008; 105:3521–3526. [PubMed: 18303113]
20. Miller NA, et al. Management of High-Throughput DNA Sequencing Projects: Alpheus. *J. Comput. Sci. Syst. Biol.* 2008; 1:132–148. [PubMed: 20151039]
21. Mane SP, et al. Transcriptome sequencing of the Microarray Quality Control (MAQC) RNA reference samples using next generation sequencing. *BMC Genomics.* 2009; 10:264. [PubMed: 19523228]
22. Jones S, et al. Comparative lesion sequencing provides insights into tumor evolution. *Proc. Natl Acad. Sci. USA.* 2008; 105:4283–4288. [PubMed: 18337506]
23. Lynch M, et al. A genome-wide view of the spectrum of spontaneous mutations in yeast. *Proc. Natl Acad. Sci. USA.* 2008; 105:9272–9277. [PubMed: 18583475]
24. Haag-Liautard C, et al. Direct estimation of per nucleotide and genomic deleterious mutation rates in *Drosophila*. *Nature.* 2007; 445:82–85. [PubMed: 17203060]
25. Moffatt MF, et al. Genetic variants regulating ORMDL3 expression contribute to the risk of childhood asthma. *Nature.* 2007; 448:470–473. [PubMed: 17611496]
26. Heap GA, et al. Genome-wide analysis of allelic expression imbalance in human primary cells by high-throughput transcriptome resequencing. *Hum. Mol. Genet.* 2010; 19:122–134. [PubMed: 19825846]
27. Jones PA, Baylin SB. The epigenomics of cancer. *Cell.* 2007; 128:683–692. [PubMed: 17320506]
28. Cedar H, Bergman Y. Linking DNA methylation and histone modification: patterns and paradigms. *Nature Rev. Genet.* 2009; 10:295–304. [PubMed: 19308066]
29. Kaminsky ZA, et al. DNA methylation profiles in monozygotic and dizygotic twins. *Nature Genet.* 2009; 41:240–245. [PubMed: 19151718]
30. Utz U, et al. Skewed T-cell receptor repertoire in genetically identical twins correlates with multiple sclerosis. *Nature.* 1993; 364:243–246. [PubMed: 7686632]
31. Sawkins MC, et al. Comparative map and trait viewer (CMTV): an integrated bioinformatic tool to construct consensus maps and compare QTL and functional genomics data across genomes and experiments. *Plant Mol. Biol.* 2004; 56:465–480. [PubMed: 15604756]
32. Simpson JT, et al. ABySS: a parallel assembler for short read sequence data. *Genome Res.* 2009; 19:1117–1123. [PubMed: 19251739]
33. Pop M, Phillippy A, Delcher AL, Salzberg SL. Comparative genome assembly. *Brief Bioinform.* 2004; 5:237–248. [PubMed: 15383210]
34. Wu TD, Nacu S. Fast and SNP-tolerant detection of complex variants and splicing in short reads. *Bioinformatics AOP.* 2010 Feb 10.





**Figure 1. Comparison of genomic locations of heterozygous cSNPs exhibiting imbalanced allelic expression in mRNA of twins 041896-001 (a) and -101 (b)**

Allelic imbalance was detected in cSNPs called by  $\geq 10$  gDNA reads with  $Q \geq 20$  and where 20–80% of uniquely aligning gDNA reads called the SNP, together with detection in  $\geq 10$  mRNA reads with  $Q \geq 20$ . 268 of 14,461 heterozygous cSNPs (1.9%) showed significant allelic imbalance in expression ( $p < 10^{-7}$ ), of which 153 (57%) were of the same magnitude and direction in both subjects.



**Figure 2. Comparisons of methylation of genomic CpG sites in CD4<sup>+</sup> lymphocytes and breast and lung tissue samples**

**a**, Frequency distribution of CpG site methylation in 041896-001 (blue) and -101 (red) using ELAND-extended. **b-j**, Pairwise comparisons of CpG site methylation using ELAND-extended in CD4<sup>+</sup> lymphocytes from MZ twin siblings 041896-001 and -101 (**b**), 230178-001 and -101 (**c**) and 041907-001 and -101 (**d**); inter-individual differences between CD4<sup>+</sup> lymphocytes from 041896-001 and 041907-001 (**e**) and 041896-001 and 230178-101 (**f**); neoplastic differences between breast tissue and breast cancer (**g**) and between normal

lung tissue and lung cancer (**h**); and between-tissue differences between CD4<sup>+</sup> lymphocytes and breast tissue (**i**) and lung tissue (**j**).

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

**Table 1**

SNP and indel genotype identities and differences between siblings in three twin pairs

Genotype Change & Individual	Platform	Twin Pair 041896				Twin Pair 230178				Twin Pair 41907			
		SNP Genotypes	Replicated SNP Genotype Differences <sup>4</sup>	Indel Genotypes	Replicated Indel Genotype Difference	SNP Genotypes	Replicated SNP Genotype Difference	Indel Genotypes	Replicated Indel Genotype Difference	SNP Genotypes	Replicated SNP Genotype Difference	Indel Genotypes	Replicated Indel Genotype Difference
No change	genome-seq <sup>1</sup> SNP array (X2)	1,086,309	79,209	26,908	91 (91.9%)	n.d.	-	n.d.	-	n.d.	-	n.d.	-
	mRNA-seq <sup>3</sup>	736,782	1,638 (98.3%)	-	-	783,189	-	-	796,870	-	-	-	-
	genome-seq <sup>1,2</sup> SNP array (X2)	51,201	8,816 (98.2%)	1,314	91 (91.9%)	39,816	888 (95.3%)	1,034	18,123	385 (98.0%)	397	397	
Ref in '001 → Het In '101	mRNA-seq <sup>2,3</sup>	202	0	3	0	n.d.	-	n.d.	n.d.	-	n.d.	n.d.	
	genome-seq <sup>1,2</sup> SNP array (X2)	32	0	-	-	36	0	-	32	0	-	-	
	mRNA-seq <sup>2,3</sup>	12	0	0	0	6	0	0	2	0	0	0	
Het in '001 → Ref in '101	genome-seq <sup>1,2</sup> SNP array (X2)	134	0	1	0	n.d.	-	n.d.	n.d.	-	n.d.	n.d.	
	mRNA-seq <sup>2,3</sup>	49	0	-	-	31	0	-	11	0	-	-	
Het in '001 het. → Hom Variant in '101	genome-seq <sup>1,2</sup> SNP array (X2)	5	0	0	0	9	0	0	16	0	0	0	
	mRNA-seq <sup>2,3</sup>	1513	0	128	0	n.d.	-	n.d.	n.d.	-	n.d.	n.d.	
	genome-seq <sup>1,2</sup> SNP array (X2)	29	0	-	-	24	0	-	17	0	-	-	
	mRNA-seq <sup>2,3</sup>	203	0	7	0	573	0	23	170	0	5	5	
Hom Variant in '001 → Het in '101	genome-seq <sup>1,2</sup> SNP array (X2)	1392	0	81	0	n.d.	-	n.d.	n.d.	-	n.d.	n.d.	
	mRNA-seq <sup>2,3</sup>	16	0	-	-	62	0	-	60	0	-	-	
	genome-seq <sup>1,2</sup> SNP array (X2)	102	0	1	0	429	0	16	192	0	5	5	

Genotype categories: homozygous reference (Ref), heterozygous variant (Het), homozygous variant (Hom Variant).

<sup>1</sup> Nucleotide genotyped if 11–44X coverage & Q ≥20;

<sup>2</sup> Genotypes determined according to frequency cutoffs in Supp. Table 8 & differences called where frequencies differed by >50%;

<sup>3</sup> Genotyped if present in >2 reads, >1 uniquely aligning read & Q ≥20;

<sup>4</sup> Detected by platform on corresponding row, replicated by platform listed on row below.

Comparison of CpG sites and CpG clusters between CD4<sup>+</sup> lymphocytes from three pairs of MZ twins, breast and lung cancer and normal tissue samples

**Table 2**

Genomic DNA Sample	CpG sites*	CpG clusters	Ratio of CpGs to Clusters	CpGs Shared	CpG Clusters Shared	mCpG Unique to One Sample**	Between sample comparison**	CpGs Shared	CpG Clusters Shared	mCpG Unique to One Sample**
041896-001 T Cell	2,146,620	1,230,241	1.74	98.1%	98.2%	2	041896- & 230178-	97.4%	97.7%	522
041896-101 T Cell	2,033,078	1,190,741	1.71			0	001 T Cell			305
230178-001 T Cell	1,636,285	1,038,787	1.58	97.8%	97.9%	3	041896-001 & 230178-101 T Cell	96.5%	96.9%	445
230178-101 T Cell	1,917,131	1,155,024	1.66			7				362
041907-001 T Cell	1,779,140	1,094,361	1.63	90.6%	92.7%	174	041896- & 041907-	97.5%	98.1%	304
041907-101 T Cell	1,642,200	1,038,090	1.58			2	001 T Cell			282
Normal Breast	1,829,855	1,086,405	1.68	96.7%	97.9%	696	041896-001 T Cell	97.3%	98.0%	5620
Breast Cancer	2,010,173	1,192,180	1.69			861	& Normal Breast			1560
Normal Lung	2,096,524	1,216,046	1.72	97.9%	98.8%	6,891	041896-001 T Cell	96.1%	97.0%	3329
Lung Cancer	1,619,178	956,760	1.69			9,618	& Normal Lung			926

\* > 10 RRBS reads aligned by ELAND-extended & Q > 20;

\*\* CpG > 80% methylated in one sample & < 20% in other;

\*\*\* Not replicated upon RRBS read alignment with GSNAP.