# Common epigenetic variation in a European population of mentally healthy young adults

Annette Milnik [a, b, c, *], Christian Vogler [a, b, c], Philippe Demougin [b, d], Tobias Egli [a, b],
Virginie Freytag [a, b], Francina Hartmann [a, b], Angela Heck [a, b, c], Fabian Peter [a, b, d],
Klara Spalek [b, e], Attila Stetak [a, b, c, d], Dominique J.-F. de Quervain [b, c, e],
Andreas Papassotiropoulos [a, b, c, d], Vanja Vukojevic [a, b, d]

[a] Division of Molecular Neuroscience, Department of Psychology, University of Basel, CH-4055, Basel, Switzerland
[b] Transfaculty Research Platform Molecular and Cognitive Neurosciences, University of Basel, CH-4055, Basel, Switzerland
[c] Psychiatric University Clinics, University of Basel, CH-4055, Basel, Switzerland
[d] Department Biozentrum, Life Sciences Training Facility, University of Basel, CH-4056, Basel, Switzerland
[e] Division of Cognitive Neuroscience, Department of Psychology, University of Basel, CH-4055, Basel, Switzerland

## ARTICLE INFO

## ABSTRACT

DNA methylation represents an important link between structural genetic variation and complex phenotypes. The study of genome-wide CpG methylation and its relation to traits relevant to psychiatry has become increasingly important. Here, we analyzed quality metrics of 394,043 CpG sites in two samples of 568 and 319 mentally healthy young adults. For 25% of all CpGs we observed medium to large common epigenetic variation. These CpGs were overrepresented in open sea and shore regions, as well as in intergenic regions. They also showed a strong enrichment of significant hits in association analyses. Furthermore, a significant proportion of common DNA methylation is at least partially genetically driven and thus may be observed similarly across tissues. These findings could be of particular relevance for studies of complex neuropsychiatric traits, which often rely on proxy tissues.

© 2016 The Author(s). Published by Elsevier Ltd. This is an open access article under the CC BY-NC-ND license (http://creativecommons.org/licenses/by-nc-nd/4.0/).

## 1. Introduction

Phenotypic differences between individuals are only partially explained by genetic differences. Additional sources of variation, including epigenetic regulation, are also centrally involved in trait variability and disease etiology (Petronis, 2010). Among epigenetic mechanisms, DNA methylation is being studied most extensively, because the available technology allows for the investigation of methylation patterns at both high resolution and throughput (Bibikova et al., 2011). DNA methylation represents an important connection between structural genetic variation and complex phenotypes (Tan et al., 2015), a link that can be investigated by epigenome-wide association studies (EWAS) (Jones, 2012; Rakyan et al., 2011). However, EWAS are challenged by the high complexity of the methylation signal, which displays variability related to such factors as the specifics of the population being studied, cell type, as well as temporal dynamics Davies et al., 2012; Horvath et al., 2012; Lokk et al., 2014; Pidsley et al., 2014. Furthermore, such signals consist of varying amounts of measurement error and systematic variance of no interest (e.g. technical artifacts).

Common epigenetic variation can be broadly defined as the signal's variability based on biologically or environmentally driven factors but not based on technical confounders or random noise (Altman, 2005; Flanagan, 2015). Given that the power of EWAS to identify robust biomarkers is greater for common variants in comparison to rare epimutations (Flanagan, 2015), an accurate estimation of populations' common epigenetic variation is of great importance. Replication analysis can be used to estimate the amount of a trait's naturally occurring variation (i.e. signals' variance exceeding the technical variance and random noise (Altman, 2005)) in a given population. Importantly, within the scope of

* Corresponding author. University of Basel, Division of Molecular Neuroscience, Birmansgasse 8, CH-4055, Basel, Switzerland.
E-mail address: annette.milnik@unibas.ch (A. Milnik).

high-density DNA methylation data, replication analyses can be done on different levels. Typically, analysis of technical replicates for CpG methylation data has been done on a methylome-wide level, by comparing the signal of all CpGs, e.g. by correlating two technical replicates of one DNA sample across measured CpGs (Tylee et al., 2013). Such analysis refers to comparing whole methylation profiles. On this methylome-wide level, high reproducibility within and between technologies has been shown repeatedly (Bibikova et al., 2011; Sandoval et al., 2011), as expected from the bimodal distribution of DNA methylation across the methylome.

However, association studies such as EWAS are done on the single-CpG level, where we often see an unimodal signal distribution (i.e. either high- or low- methylation sites). Therefore, we aimed to estimate the amount of detectable common epigenetic variation for each single CpG by evaluating the technical reliability on the single-CpG level, i.e. by correlating two repeated measurements of a single CpG site across samples. We performed a comprehensive reliability analysis of the single-CpG DNA methylation signal from the Illumina Infinium Human Methylation 450 K array (450 K array), in a large sample of healthy young adults. We used this reliability analysis to derive a lower bound estimate of the common epigenetic variation, given its presumed strong association with complex phenotypes. Finally we tested and replicated this assumption by conducting a series of association studies with CpG methylation in two independent samples.

## 2. Materials and methods

### 2.1. Subjects

The subjects included in this study (main sample $N = 568$; mean age 23.8 y, 18.3–36.8 y; 59% females; independent replication sample $N = 319$; mean age 24.1 y, 18.3–36.5 y; 70% females; the age information refers to the time-point of the blood sampling, see below) represent subsets of two ongoing studies, which were previously described (Heck et al., 2014; Spalek et al., 2015). The purpose of both studies is to identify biological correlates of cognitive performance by using genetics, electroencephalography and imaging techniques in healthy young adults from the general population. Saliva samples were collected at the time-point of the main investigation. Subjects were re-invited for an additional saliva and blood sampling, which took place on average 360 days (min 1 day; max 1384 days; median 341 days) after the main investigation. Samples were collected between midday and evening (mean time = 2:30 p.m., range 1:00 p.m. − 8:00 p.m.). Hematological analysis, including blood cell counts, was performed with Sysmex pocH-100i™ Automated Hematology Analyzer (Sysmex Co, Kobe, Japan). Subjects were of good general health, free of any self-reported neurological or psychiatric illness and did not take any medication (apart from oral contraception) at both time points. The investigation was carried out in accordance with the latest version of the Declaration of Helsinki. The ethics committee of the Cantons of Basel-Stadt and Basel-Landschaft approved the studies. All participants received general information about the study and gave written informed consent.

### 2.2. Affymetrix SNP 6.0 based genotyping and imputation

DNA isolated from saliva was investigated with Affymetrix SNP 6.0 array as described in the Genome-Wide Human SNP Nsp/Sty 6.0 User Guide (Affymetrix, Santa Clara, CA USA; see Supplementary Material). The mean call-rate per subject was 98.5% (90.1%– 99.7%). $N = 35$ subjects out of the main sample were identified as outliers and excluded from the association analyses (see Supplementary Material). After basic SNP-QC (in both samples MAF > 0.02; HWE > 0.001; missing rate per SNP < 5%) $N = 659,944$ SNPs entered the association analyses. Additionally autosome-wide genotype imputation was performed (see Supplementary Material).

### 2.3. HumanMethylation Infinium Infinium 450 K BeadChip based methylation analyses

Array processing (see Supplementary Fig. 1A): DNA isolated from peripheral blood was investigated with the 450 K array (Illumina, Inc., San Diego, CA, U.S.A; see Supplementary Material). The $N = 568$ subjects of the main sample were processed in two batches (2 plates and 4 plates, respectively). For the independent replication sample, $N = 319$ subjects were processed within a single batch (4 plates). Within a batch samples were processed with a randomized plate assignment and with a single bisulfite conversion.

For $N = 145$ subjects of the main sample a technical replication was performed and processed on a 450 K array in parallel with the two batches of the main sample (2 plates, 1 additional plate per batch; technical replicates). Of note, the technical replicates of the identical DNA were not processed within the same time-point, and hence were bisulfite-converted independently (see below). The technical replication measurements were performed starting from the identical DNA material as in the main sample (single DNA isolation), with randomized assignment.

Data preparation on batch-level (see Supplementary Fig. 1B): Preprocessing was done separately for each batch and also separately for the technical replicates (two batches for the main sample; two batches for the technical replicates; one batch for the independent replication sample). Data were extracted and analyzed from the generated idat files using the R package RnBeads version 0.99.9 (Assenov et al., 2014). CpG annotation was based on the manufacturer's annotation file (Human-Methylation450_15017482_v.1.2). During preprocessing, the background was subtracted using the "noob" method in the methylumi package (Davis et al., 2014), and the signal was further normalized using the SWAN algorithm (Maksimovic et al., 2012). The following probe categories were excluded from the final data sets, based on the annotation provided within the RnBeads package: non-CpG context probes (due to underrepresentation on the 450 K array, 0.6%, (Bibikova et al., 2011), functional differences when compared to CpG context as well as very low abundance of non-CpG methylation in somatic tissues (Ziller et al., 2011); $N = 3091$), probes with a SNP mapping directly to the target CpG site, as well as probes with three and more SNPs mapping within the 50mer probe (see Supplementary Fig. 2; MAF threshold was set to 0.01; $N = 18,998$ CpGs), gonosomal probes ($N = 11,473$ CpGs), non-specific probes. Using the Greedycut algorithm, we iteratively removed the probes and data sets of the highest impurity (rows and columns in the detection $p$-value table that contain the largest fraction of unreliable measurements; $p < 0.05$ (Assenov et al., 2014)).

Data preparation on sample-level (see Supplementary Fig. 1C): Postprocessing was further done for each sample separately, combining the β-values of the preprocessed data of all batches per sample (see Supplementary Fig. 3A,C, as well as Supplementary Figs. 4 and 5 for diagnostic plots of the data). The β-values were further processed step-by-step in order to correct for further influential and putative confounding factors: 1) using logit-transformation (M-value, (Du et al., 2010), done with the R-package car (Fox and Weisberg, 2011)); 2) z-transformation per plate (correcting for plate and batch effects); 3) regressing out the first 8 (main sample and technical replicates) or 7 (independent

replication sample) axes of a principal component analysis (PCA, done with the R-package pcaMethods (Stacklies et al., 2007); see Supplementary Fig. 3B,D, as well as Supplementary Figs. 6 and 7 for diagnostic plots of the data after applying steps 1–3 of post-processing). The PCA was based on CpGs with no missing values (>95% of the included CpGs). The PCA-based approach corrected for technical biases as well as for part of the variability induced by blood cell composition (see Supplementary Tables 1 and 2, Supplementary Figs. 9 and 10); 4) regressing out the effects of sex and age; 5) regressing out the effects of variants in the 50mer probe sequence, if the total variance explained by these variants exceeded 0.1% (see Supplementary Material and Supplementary Fig. 2).

The accepted missing rate per CpG was set to <1%. We further excluded cross-hybridizing probes and polymorphic CpG sites (Chen et al., 2013; Price et al., 2013) ($N_{max}$ = 63,974). Only CpGs surviving all filtering steps in all samples were used for the downstream analyses ($N$ = 394,043).

Cell-count estimates of cell types (CD 8 T helper cells, CD 4 T helper cells, natural killer cells, B-cell, Monocytes, Granulocytes) were done with the minfi-package (Aryee et al., 2014) in R, based on the algorithm provided by (Houseman et al., 2012), adapted for the Illumina 450 K array (Jaffe and Irizarry, 2014).

To validate the applied batch and plate correction, we also used ComBat (Johnson et al., 2007; Leek et al., 2012) instead of z-trans-formation per plate, which lead to a nearly identical signal per CpG ($r_{mean}$ = 0.999). Optionally, we applied SQN (stratified quantile normalization, with and without outlier exclusion; 5 samples were identified as outlier) (Touleimat and Tost, 2012) as an alternative data processing method (see Supplementary Fig. 8). When compared with SWAN normalization, the two processing methods produced on average a similar signal per CpG ($r_{mean}$ = 0.75).

## 2.4. Technical replication analyses

All analyses of technical replicates were based on the $N$ = 145 individuals from the main sample for which technical replicates were available. Pearson correlation coefficients ($r$) were used for all correlation analyses. The dataset used was analysis-dependent (see below).

On the methylome-wide level, we calculated $r$ between methylation profiles for all CpGs based on the β-values of the main sample and its technical replicates (see Supplementary Fig. 1D): either for the identical DNA (within-subject comparison, $N$ = 145 in total) or for DNA of different subjects (between-subject compari-sons; $N$ = 144 comparisons for each subject of the main sample with all other subjects from the main sample and $N$ = 144 com-parisons for each subject of the main sample with all non-identical technical replicates).

On single-CpG level we calculated $r$ separately for each CpG (see Supplementary Fig. 1E), which resulted in $N$ = 394,043 correlation coefficients. This analysis was done after applying the logit trans-formation (M-values), the correction for plate effects and the PCA correction to the datasets. To further evaluate the distribution of these $r$-values, we applied a Gaussian fit allowing up to 5 over-lapping Gaussian distributions by using the optimx function in R (settings: method L-BFGS-B, ndeps 0.0001, maxit 40,000). When allowing more than 5 Gaussian distributions, the contribution of the additional sub-distribution reached 0% (see Supplementary Fig. 11). The minimum and maximum values were restricted to $-1 \geq m \leq 1$, $0 \geq sd \leq 1$ and $0 \geq p \leq 1$. The starting values were as follows: m = c(0, 0.5, 0.6, 0.7, 0.8); sd = 0.12; p = c(0.6, 0.1, 0.1, 0.1, 0.1).

To obtain a random $r$-distribution for a sample size of $N$ = 145 subjects, we repeatedly ($N$ = 400,000 times) generated two standard normal random variables (length of 145 each) and calculated $r$ between these two variables.

## 2.5. Association analyses

The SNP-association studies were performed on a genome-wide scale for each CpG separately assuming an additive genetic model and applying an epi-genome-wide Bonferroni correction ($\alpha$ = 5%, correcting for 659,944 × 394,043 tests, resulting in $p_{bonf} < 1.9 \times 10^{-13}$). We used M-values that were additionally cor-rected for plate and batch effects, 7–8 axes of a PCA, sex and age as well as effects of variants in the 50mer probe sequence (see Supplementary methods). For the cis-analyses (defined as ± 3.5 Mbp), we used a less stringent significance threshold $p_{cis} < 1.7 \times 10^{-5}$ ($\alpha$ = 5%, correcting for at least 3000 independent tests per CpG).

To evaluate the associations with sex and age we used the PCA corrected dataset (M-values that were additionally corrected for plate and batch effects as well as 7–8 axes of a PCA). A linear model was calculated for each CpG, including sex and age as independent variables (per independent variable $\alpha$ = 5%, 394,043 independent tests, $p < 1.3 \times 10^{-7}$). Because sex and age were profoundly asso-ciated with cell-count estimates (main sample $R^2_{sex}$ = 16.32%, $p < 2.2 \times 10^{-16}$, $R^2_{age}$ = 6.21%, $p = 7.28 \times 10^{-6}$; independent replication sample $R^2_{sex}$ = 26.01%, $p < 2.2 \times 10^{-16}$, $R^2_{age}$ = 7.76%, $p = 2.94 \times 10^{-4}$), we additionally used cell-count estimates as covariates in the model. We estimated standardized betas for sex and age by z-transforming the input variables.

For the association with smoking, we used the data that was also corrected for sex and age (after the PCA correction) and calculated $r$ between single CpGs and smoking frequency ($\alpha$ = 5%, 394,043 in-dependent tests, $p < 1.3 \times 10^{-7}$). Smoking frequency was assessed on a 4-point Likert scale (0 = never, 1 = occasionally, 2 = 1–5 cigarettes/day, 3 = 6–20 cigarettes/day, 4 = 20 or more cigarettes/ day; see Supplementary Fig. 12) during the main investigation. There was no significant association between cell-count estimates and smoking (main sample $R^2_{smoking}$ = 1.2%, $p = 0.38$; independent replication sample $R^2_{smoking}$ = 0.45%, $p = 0.97$).

We compared the association results of sex, age and smoking with the following external datasets: sex (Xu et al., 2014), age (Horvath, 2013) and smoking (Guida et al., 2015; Shenker et al., 2013).

To assess the effects of batch and plate, we used the uncorrected β-values as dependent variable and calculated one linear model per CpG with one combined factor for batch and plate, since batch and plate were largely confounded (see Supplementary Fig. 1). We report the overall $R^2$ of this model as effect-size estimate.

## 2.6. Genomic representation analyses

We used the manufacturer's annotation file (Human-Methylation450_15017482_v.1.2) to classify CpGs based on their CpG density (Island, N_Shelf, N_Shore, Open_Sea, S_Shore, S_Shelf) as well as based on functional regions (1st Exon, 3′UTR, 5′UTR, Body, Intergenic, TSS1500, TSS200). We used the genomic hg19 database (genome-mysql.cse.ucsc.edu) to retrieve data about DNase I hypersensitivity cluster (table wgEncode-RegDnaseClusteredV3) and data about H3K27ac histone modifica-tion marks (data derived from wgEncodeBroadHistone-XXX-cell-type-XXX-H3k27acStdSig.bigWig for the cell-types Gm12878, H1hesc, Hsmm, Huvec, K562, Nhek, Nhlf; data has been log-transformed). We compared CpGs showing high and low natural occurring variation with respect to these genomic representations either by using the Chi$^2$-test to compare frequencies or the Kolmogorov-Smirnov test to compare distributions.

### 2.7. Software

If not mentioned differently, analyses were conducted in R (version: 2.15.1 and higher, R Development Core Team, 2012) or PLINK (Purcell et al., 2007).

## 3. Results

### 3.1. Methylome-wide reliability analyses

These analyses were done with the untransformed β-values. Based on the β-distribution on methylome-wide level (Fig. 1A), we first compared DNA methylation profiles of the identical DNA (within subjects comparison). Here, we observed high signal reproducibility of technical replicates (see Table 1; see Fig. 1B for one example). This was in agreement with previous reports (Bibikova et al., 2011; Sandoval et al., 2011). Next, we compared DNA methylation profiles between subjects (see Fig. 1C for one example). The estimated signal similarity between DNA of different subjects also suggested very high signal consistency on methylome-wide level (see Table 1). Subsequently, for each dataset of the main sample we tested, which of the remaining datasets of the main sample (including technical replicates) showed the highest similarity on methylome-wide level (top-hit). For 89% of these comparisons, the highest similarity was obtained with its technical replicate, i.e. identical DNA that has been processed independently.

### 3.2. Single CpG level reliability analyses — estimation of natural occurring variation

With the technical replication analysis on the single-CpG level (Fig. 1D—F) we further aimed to detect CpGs that comprise a significant amount of common epigenetic variation, within those samples with a technical replicate. We therefore focused on separating the technical variation or random noise from the overall detected signal. For this analysis we used data after logit-transformation (M-values) that was additionally corrected for plate and batch effects as well as the first 8 axes of a PCA, i.e. that is corrected for additional confounders (see Supplementary Table 1).

Taking into account all $N = 394,043$ CpGs, the average $r$ per single-CpG was 0.191 ($r_{min} = -0.361$, $r_{max} = 0.988$). Fig. 2A shows the distribution of $r$-values across all CpGs. Infinium type II probes in comparison to type I probes showed on average higher $r$-values (type I $r_{mean} = 0.17$; type II $r_{mean} = 0.20$; Supplementary Fig. 13). After applying an FDR-correction for multiple testing, 134,385 CpGs showed a significant positive correlation ($\alpha = 5\%$, $r > 0.198$). To further examine this distribution, we applied a Gaussian mixture model allowing up to five Gaussian distributions (Fig. 2A), which indicated that for approximately 53% of the CpGs in our sample the observed signal variability was likely based on random or technical variance (red sub-distribution Fig. 2A, mean = 0.03, sd = 0.09; alternative Gaussian mixture models suggest similar conclusion, see Supplementary Fig. 11). The majority of these CpGs' $r$-values were below 0.3. We compared this fitted distribution with a
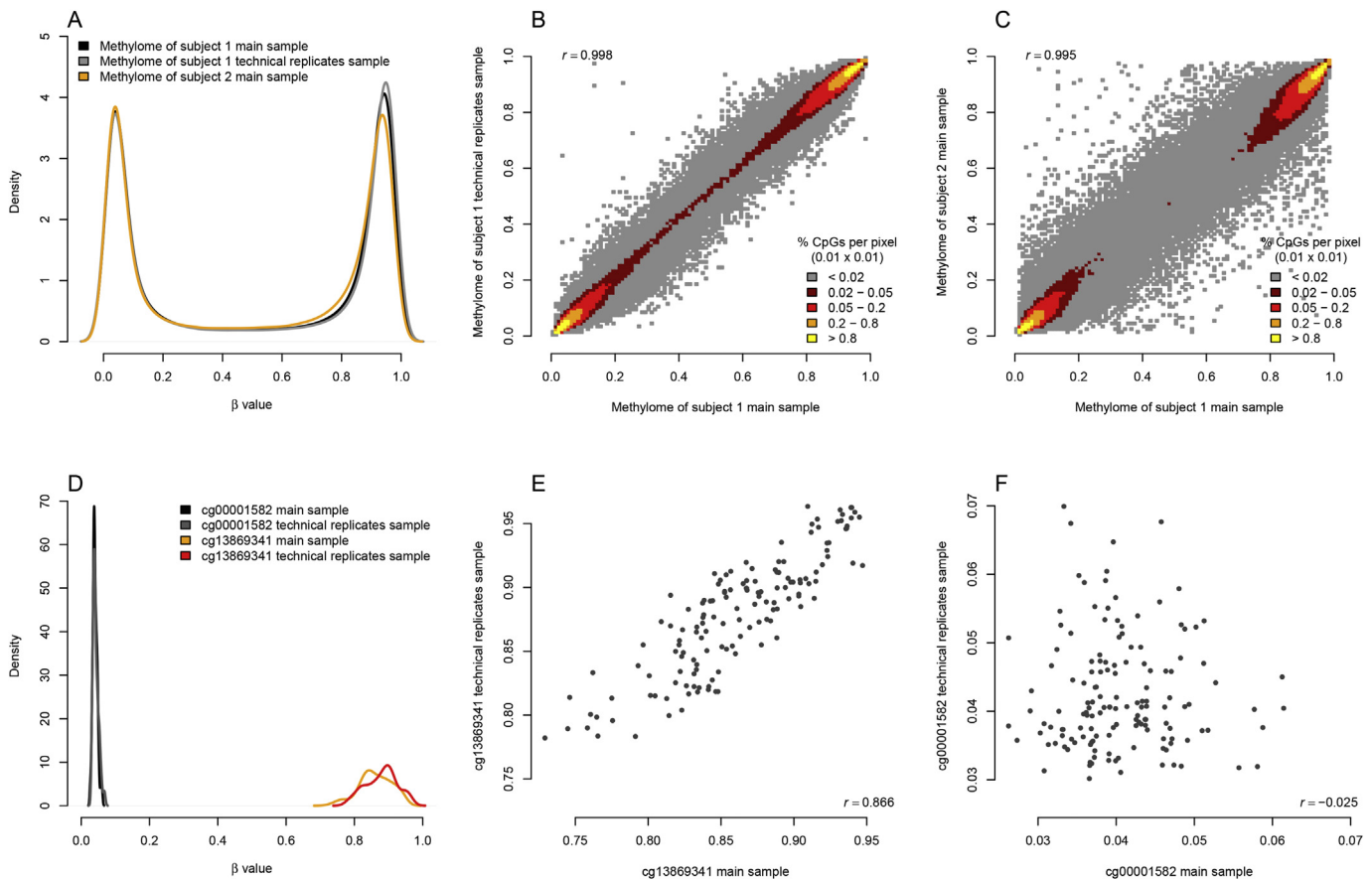


Fig. 1. **Methylation signal distributions and replication analyses.** (**A**) β-distribution on methylome level, shown for the DNA of two different subjects. For the first subject, both methylome data sets, from the main sample and from its technical replicate, are shown. β-values on methylome level plotted against each other from (**B**) the same subject and (**C**) two different subjects. (**D**) β-distribution shown for two distinct CpGs, separately for the main sample and its technical replicate ($N = 145$ pairs with technical replicates). (**E** and **F**) β-values of two distinct CpGs plotted against each other ($N = 145$ samples with technical replicates). $r$: Pearson correlation coefficient.

**Table 1**
Methylome-wide reliability analyses. For the analytical schema see also Supplementary Fig. 1D.

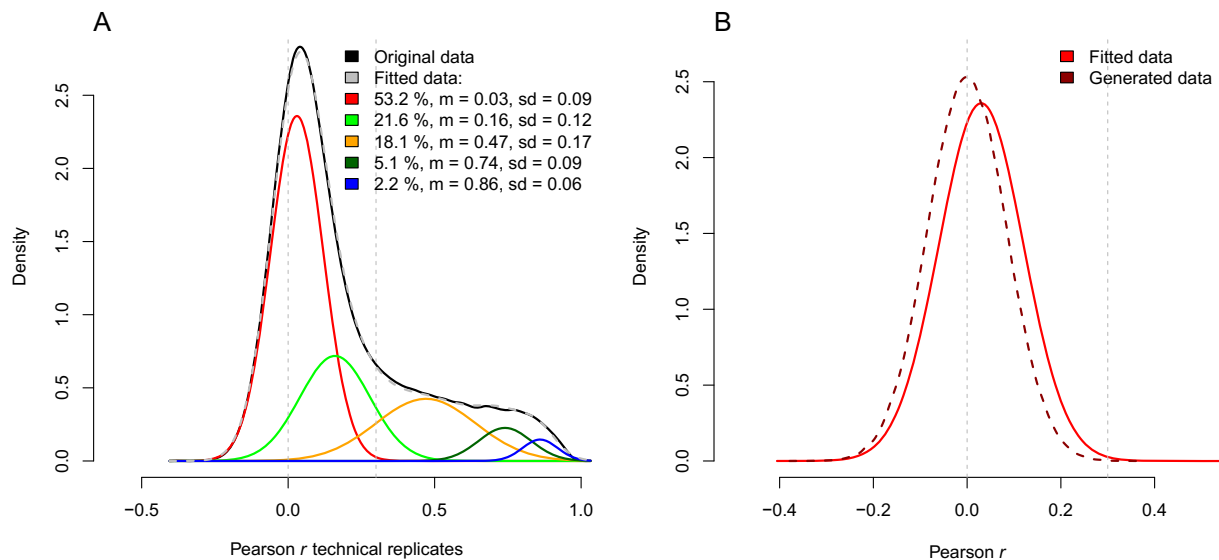| Comparison between: | Mean $r$ | Min $r$ | Max $r$ | % of top-hit |
|---|---|---|---|---|
| Technical replicates (identical DNA) | 0.997 | 0.990 | 0.999 | 89 |
| Different DNA within the main sample | 0.995 | 0.988 | 0.997 | 11 |
| Different DNA between the main sample and non-identical technical replicates | 0.994 | 0.985 | 0.997 | 0 |

$r$: Pearson correlation coefficient.



**Fig. 2. Replication analysis on single-CpG level. (A)** The $r$ distribution based on all $N = 394{,}043$ CpGs is depicted in black. Results of the Gaussian fit are depicted in grey and colored lines. The legend shows for each estimated sub-distribution the percentage of CpGs, the center (m) and the standard deviation (sd). **(B)** Superimposition of an $r$ distribution based on random signals to the fitted Gaussian distribution with the mean closest to zero (as represented by the red curve in panel **A**). The randomly generated distribution is based on standard normal random variables (length $N = 145$) and was adjusted by the corresponding probability of random probes from the Gaussian fit (53%). Vertical dotted grey lines depict the center of a random distribution (m = 0) and an $r$ of 0.3. Above $r = 0.3$ it is unlikely that a CpG shows a signal variability based on random signals only. $r$: Pearson correlation coefficient. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

distribution of $r$-values, which were based on generated random data sets (Fig. 2B; mean = 0, sd = 0.08). The fitted distribution showed a slight shift to the right and was slightly wider than the generated random distribution (mean = 0.03 in comparison to 0; sd = 0.09 in comparison to 0.08; Kolmogorov-Smirnov test $D = 0.14$, $p < 2.2 \times 10^{-16}$), indicating that within these CpGs, there is still some non-random variation. Based on the visual inspection of the plots, and specifically the boundaries of the red sub-distribution, we defined an $r$-threshold of 0.3 to separate CpGs that show a considerable part of common epigenetic variation, from CpGs showing very low amount of common epigenetic variation. When applying this threshold, type II probes in comparison to type I probes showed a significantly higher proportion of reliable CpGs (20.8% of all type I probes versus 27.1% of all type II probes; $\chi^2 = 1628.56$, df = 1, $p < 2.2 \times 10^{-16}$).

For the CpGs with $r \geq 0.3$ (reliable CpGs, $N = 99{,}839$) we assume that they robustly capture the underlying common epigenetic variation. These CpGs exhibited a higher inter-individual variability in comparison to CpGs with an $r < 0.3$ and had less extreme β-values (Figs. 3 and 1D-F; Supplementary Figs. 14 and 15). On the other hand, the signal coming from the subset of CpGs with an $r < 0.3$ most likely comprised a mixture of two different categories: random signals or signals with very low natural epigenetic variation (~2/3 of all CpGs with $r < 0.3$; red sub-distribution in Fig. 2A) and signals with a medium to low natural epigenetic variation (~1/3 of all CpGs with $r < 0.3$; light-green sub-distribution in Fig. 2A).

Differences in data post-processing could potentially bias the technical replication analysis. Therefore, we performed the

reliability analysis on the CpG level upon different post-processing steps (Supplementary Fig. 16). Correcting for the technical influences of plate and batch improved the reproducibility on the single-CpG level. To additionally evaluate the impact of the technical factors batch and plate on the single-CpG signal distribution, we compared the signals mean and standard deviation with the estimated variance explained ($R^2$) by these technical factors: CpGs with lower variability and more-extreme β-values showed more-prominent batch and plate-effects than CpGs with higher variability and less-extreme β-values (Supplementary Fig. 17). A possible explanation is that more-extreme and less-variable CpGs are more-often monomorphic, and that the detected variability is primary due to technical artifacts. The PCA-correction applied could not entirely correct for cell-type composition (see Supplementary Figs. 9 and 10). As a consequence, there was still enrichment for association signals with cell-count estimates after PCA-correction, which was more pronounced for reliable CpGs especially in our samples (Supplementary Fig. 18).

### 3.3. Association analyses

Having detected CpGs that are likely to show considerable amounts of naturally occurring variation, we next assessed the phenotypic relevance of such variation by means of significant hits in association studies using CpG methylation as dependent variable. Previously it has been demonstrated that DNA methylation can be influenced by genetic variation (Schalkwyk et al., 2010; Shoemaker et al., 2010), sex and age (Horvath, 2013; Xu et al.,
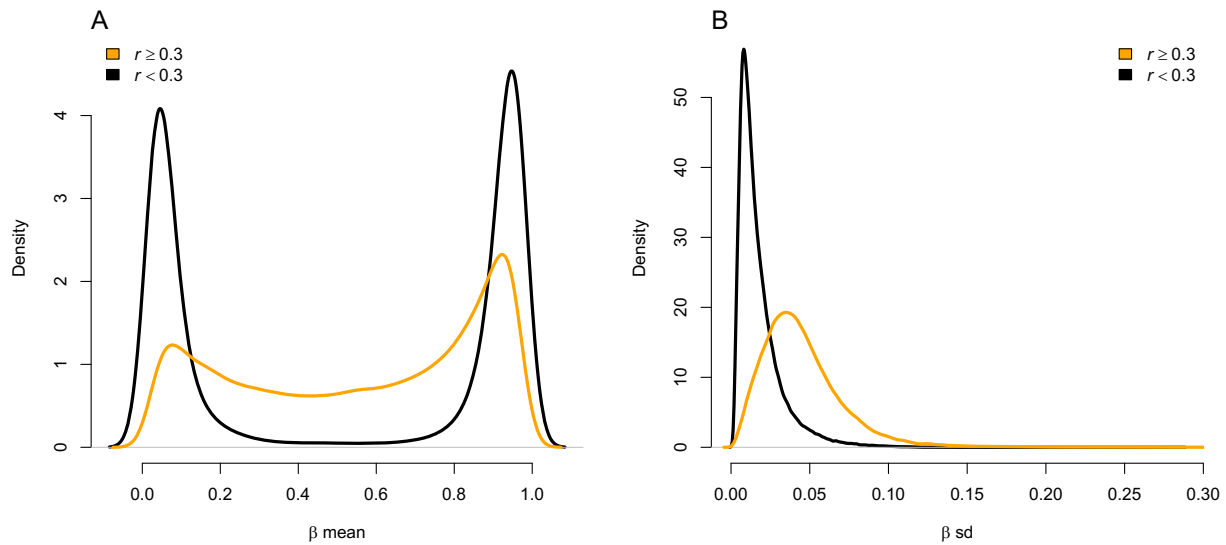
**Fig. 3. Density distributions of the beta's mean and sd, depending on the estimated common variation of the CpGs.** CpGs for which common variation could be robustly detected were less likely to show extreme β-values **(A)** and showed a wider distribution **(B)**. sd: standard deviation. r: Pearson correlation coefficient.

2014), as well as tobacco smoke exposure (Guida et al., 2015; Shenker et al., 2013). Hence, we calculated EWAS for genetic variants, sex, age and smoking in the main sample ($N = 533$, after exclusion of genetic outliers) as well as in the independent replication sample ($N = 319$). For these analyses we used M-values that were at least corrected for plate effects and for the first 7–8 axes of a PCA (for more specific information see methods).

### 3.3.1. meQTL analysis

To investigate the effect of SNPs on single-CpG methylation levels (meQTL), we first performed genome-wide association analyses on all 394,043 CpGs. Based on the main sample, this analysis identified for 8.2% of all CpGs at least one significant meQTL (see Table 2). As expected (Schalkwyk et al., 2010; Shoemaker et al., 2010), most of the top meQTLs were in *cis* of the investigated CpG (±3.5 Mbp, 95.5%).

Hence, we restricted the main analysis to a ±3.5 Mbp window surrounding the investigated CpG (*cis* analysis) and applied a less stringent *p*-value threshold ($p < 1.7 \times 10^{-5}$). Based on the main sample, we identified for 20% of all CpGs at least one significant meQTL in *cis* (Table 3). The results of the two samples showed high concordance rates for the meQTL association analysis (Supplementary Fig. 19A). CpGs that were more likely to capture common epigenetic variation showed a pronounced enrichment of significant association signals in both of the tested samples (Fig. 4A; Table 4).

Finally, we estimated the impact of the sequentially applied post-processing steps on the detection of meQTLs by repeating the same analysis for all post-processing steps. The PCA correction lead to an increase in power to detect meQTLs (Supplementary Table 3) whereas the correction of variants in the 50mer probe sequence sufficiently reduced biases of these variants in the detection of meQTLs (about 10% more hits without 50mer probe correction; Supplementary Table 4).

### 3.3.2. Association with sex, age and smoking

The phenotypic analyses revealed largely overlapping significant associations in both samples (Supplementary Fig. 19; Supplementary Table 5) with sex ($N = 5106$ in total: main sample $N = 4958$ CpGs, independent replication sample $N = 1551$ CpGs, overlap between samples $N = 1403$; $N = 892$ have been reported previously (Xu et al., 2014)), age ($N = 409$ in total: main sample $N = 352$ CpGs, independent replication sample $N = 171$ CpGs, overlap between samples $N = 114$; $N = 13$ have been reported previously (Horvath, 2013)) and smoking ($N = 19$ in total: main sample $N = 14$ CpGs, independent replication sample $N = 16$ CpGs, overlap between samples $N = 11$; $N = 18$ have been reported previously (Guida et al., 2015; Shenker et al., 2013)). Again, CpGs that were more likely to capture common epigenetic variation were more likely to show significant association results in both samples (Fig. 4B–D and Table 4). Of note, only autosomal CpGs were used for the association analysis. For all phenotypes the effect-size estimates from our study were in concordance ($|r| > 0.4$; Supplementary Figs. 20 and 21) with previously reported effect-size estimates (Guida et al., 2015; Horvath, 2013; Shenker et al., 2013; Xu et al., 2014).

### 3.4. Genomic representation of CpGs with respect to the common epigenetic variation

CpGs were classified with respect to the genomic location in the following categories: open sea, shore, shelf and island (Bibikova et al., 2011). CpGs comprising higher common epigenetic variation were overrepresented in open sea and shores, and underrepresented in CpG islands ($\chi^2 = 6707.08$, df = 5, $p < 2.2 \times 10^{-16}$; see Table 5). Accordingly, these CpGs were also overrepresented in intergenic regions ($\chi^2 = 7324.52$, df = 6, $p < 2.2 \times 10^{-16}$; see Table 6), and showed less enrichment for DNase I hypersensitivity ($D = 0.14$, $p < 2.2 \times 10^{-16}$) and for the histone modification mark

**Table 2**
Genome-wide meQTL-analyses.

| Analysis | N CpG-SNP pairs | N unique CpGs | % of CpGs | N unique SNPs |
|---|---|---|---|---|
| Main sample | 304,706 | 32,216 | 8.2 | 130,208 |
| Independent Replication sample | 170,940 | 22,308 | 5.7 | 85,301 |

DNA methylation data was corrected for variants in the 50mer probe sequence; significance threshold was set to $p < 1.9 \times 10^{-13}$.

**Table 3**
meQTL-analyses in CIS (±3.5 Mb window).

| Analysis | N CpG SNP pairs | N unique CpGs | % of CpGs | N unique SNPs |
|---|---|---|---|---|
| Main sample | 956,378 | 78,725 | 20.0 | 276,739 |
| Independent replication sample | 641,288 | 65,388 | 16.6 | 221,528 |

DNA methylation data was corrected for variants in the 50mer probe sequence. $N = 49,626$ of the unique CpGs showed a significant SNP-hit in *cis* in both studies; significance threshold was set to $p_{cis} < 1.7 \times 10^{-5}$.
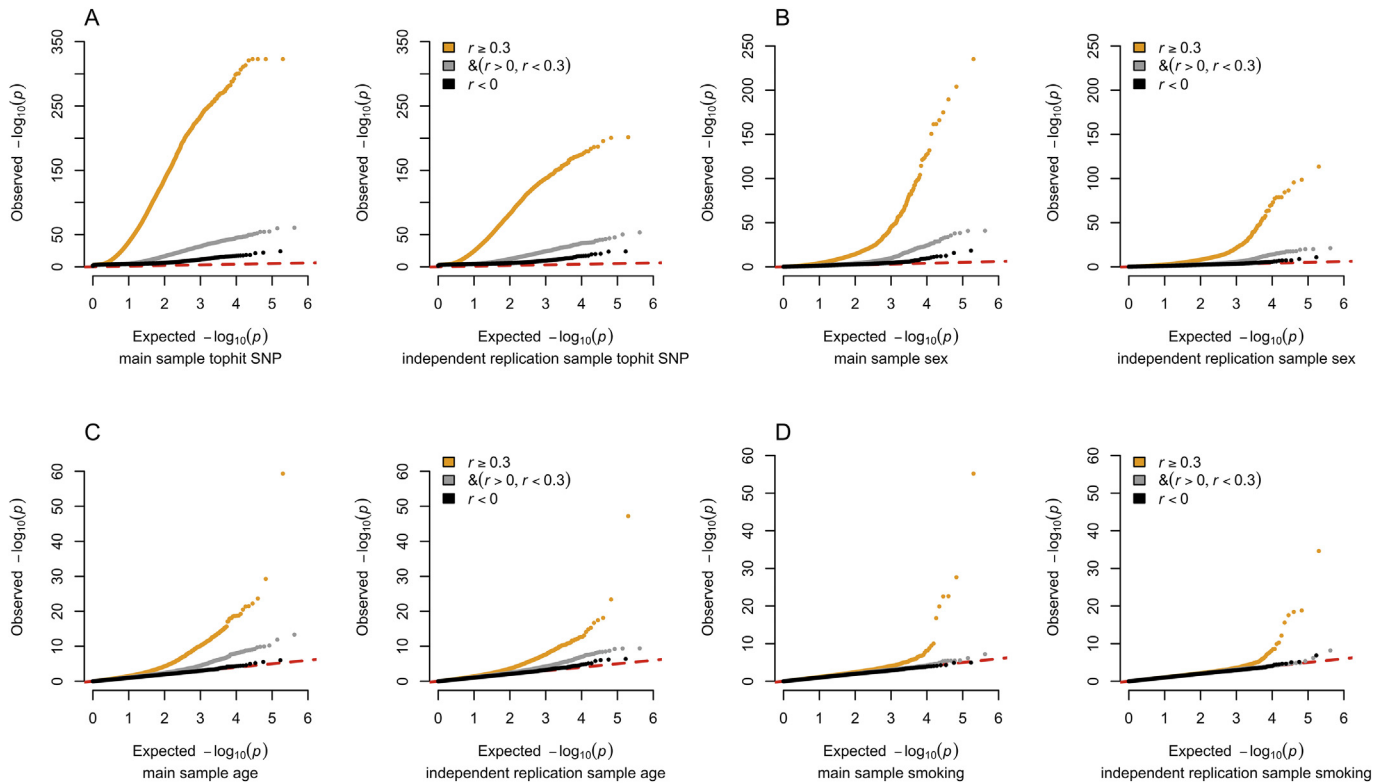


**Fig. 4. QQ-plot of the association analyses.** CpGs were classified in 3 groups, depending on the technical replication analyses and estimated common variation (orange $r \geq 0.3$; grey $0 > r < 0.3$; black $r < 0$). For each group separately, the results of a QQ-plot for the main sample (left) and the independent replication sample (right) are shown. **(A)** depicts the results of the meQTL analyses in *cis* (*p*-value of the tophit SNP), **(B)** of the association with sex, **(C)** of the association with age and **(D)** of the association with smoking. *r*: Pearson correlation coefficient. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

**Table 4**
Number of hits in the association analysis, depending on the technical reliability r per CpG.

| Study | Analysis | $r \leq 0$ ($N = 84,972$) | $0 < r < 0.3$ ($N = 209,232$) | $r \geq 0.3$ ($N = 99'839$) |
|---|---|---|---|---|
| Main sample | meQTL in *cis* | 4.93% ($N = 4,192$) | 11.99% ($N = 25,083$) | 49.53% ($N = 49,450$) |
| Ind. repl. sample | meQTL in *cis* | 4.57% ($N = 3,880$) | 9.53% ($N = 19,947$) | 41.63% ($N = 41,561$) |
| Main sample | Sex | 0.02% ($N = 19$) | 0.23% ($N = 487$) | 4.46% ($N = 4,452$) |
| Ind. repl. sample | Sex | 0.01% ($N = 5$) | 0.06% ($N = 119$) | 1.43% ($N = 1,427$) |
| Main sample | Age | 0% ($N = 0$) | 0.02% ($N = 37$) | 0.32% ($N = 315$) |
| Ind. repl. sample | Age | 0% ($N = 0$) | 0.01% ($N = 22$) | 0.15% ($N = 149$) |
| Main sample | Smoking | 0% ($N = 0$) | 0% ($N = 1$) | 0.01% ($N = 13$) |
| Ind. repl. sample | Smoking | 0% ($N = 0$) | 0% ($N = 1$) | 0.02% ($N = 15$) |

Ind. repl. sample: independent replication sample.

H3K27Ac cluster sites ($D > 0.0985$, $p < 2.2 \times 10^{-16}$; Supplementary Fig. 22).

## 4. Discussion

In the present study we performed comprehensive reliability analysis of the 450 K array in a large sample of healthy young adults. Based on the reliability analysis on the single-CpG level we estimated a medium to large common epigenetic variation ($r > 0.3$;

$r^2 > 9\%$) for 25% of the examined CpGs. These CpGs were less often hyper- and hypomethylated and showed a higher variability in the investigated population. Additionally, these CpGs were over-represented in low CpG density genomic regions like open sea and shores, as well as in intergenic regions, which was in agreement with a previous report (Gibbs et al., 2010). Accordingly, these CpGs were depleted in DNase I hypersensitivity and H3ak27ac sites. Furthermore, CpGs with an estimated higher epigenetic variability showed an enrichment of significant association signals for meQTLs

**Table 5**
Percentage of CpGs within different genomic regions, in relation to CpG density.

| %        | $r < 0.3$ | $r \geq 0.3$ |
| -------- | --------- | ------------ |
| Island   | 36.15     | 23.57        |
| N_Shelf  | 4.86      | 4.66         |
| N_Shore  | 11.95     | 17.08        |
| Open_Sea | 33.41     | 37.00        |
| S_Shore  | 9.28      | 13.62        |
| S_Shelf  | 4.35      | 4.06         |

N_Shore and S_Shore are 2 kb regions flanking CpG islands, upstream and downstream respectively; N_Shelf and S_Shelf are 2 kb regions flanking CpG island shores, upstream and downstream respectively. The remaining CpG were assigned to Open _Sea.

**Table 6**
Percentage of CpGs within different genomic regions, in relation to the transcripts "functional regions".

| %          | $r < 0.3$ | $r \geq 0.3$ |
| ---------- | --------- | ------------ |
| 1st Exon   | 5.38      | 3.54         |
| 3′UTR      | 3.90      | 2.96         |
| 5′UTR      | 9.70      | 7.54         |
| Body       | 33.97     | 33.17        |
| Intergenic | 19.61     | 30.78        |
| TSS1500    | 14.50     | 14.74        |
| TSS200     | 12.94     | 7.26         |

TSS200 is the region from Transcription start site (TSS) to −200 nt upstream of TSS; TSS1500 covers −200 to −1500 nt upstream of TSS.

as well as for sex, age and smoking, as expected for CpGs that show variability across a population (Flanagan, 2015). Importantly, we could robustly confirm the association findings in an independent replication sample, and could also replicate previously reported findings for all phenotypes (Guida et al., 2015; Horvath, 2013; Shenker et al., 2013; Xu and Taylor, 2014).

Technical replication analysis on a single-CpG level can help to identify CpGs that show common epigenetic variation in the investigated population. However, it is important to stress that we cannot draw final conclusions with regard to CpGs that showed no or very low common epigenetic variation. Variability of such CpGs may change as a function of sample size, tissue selection, samples' environmental background, genetic background, or disease status. By analyzing a large sample of healthy young adults in one tissue (i.e. blood) we most likely estimated a lower-bound variability for single-CpGs. Yet, these results may serve as a baseline reference for the naturally occurring epigenetic variation of specific CpGs, in a population of a similar origin and structure.

The results of the meQTL analyses are in line with previous studies that examined genotype-epitype associations in cell lines (Bell et al., 2011; Heyn et al., 2013), peripheral blood (McClay et al., 2015) and brain samples (Gibbs et al., 2010; Zhang et al., 2010). Of note, by focusing on the CpGs that show common epigenetic variation we could detect a profound enrichment of meQTLs, when compared to previous studies. Therefore, the subset of reliably measured CpGs could be valuable for studies of complex traits since it additionally delineates the common epigenetic variation determined by the underlying genetic architecture. Given that cross-tissue studies suggest a high level of meQTL conservation across tissues (Gibbs et al., 2010; Smith et al., 2014), it is also likely that the captured common epigenetic variation shows similar conservation across tissues. This is important, especially in the context of studies where obtaining material from the target tissue is not feasible or results in low sample size. By focusing on meQTLs findings in proxy tissues it may be possible to reflect changes in the remote target tissues, e.g. brain.

Lack of reliability on a single-CpG level in our sample may reflect a truly invariant signal (e.g. the CpG is highly methylated in all subjects), for which the measured variability is only due to chance or technical influences. However, failure to detect systematic variance can also be caused by the specifics of the microarray technology, including a fixed and limited signal resolution (Bibikova et al., 2011). Sequencing technologies can bypass this issue via customized signal resolution, e.g. by optimizing sequencing depth (Bock, 2012). Increasing sequencing depth, however, increases costs, while still being faced with the challenges of signal reliability.

Taken together, our results indicate that for the 450 K array, a relevant percentage (>25%) of single CpGs shows a medium to strong common epigenetic variation in a homogenous sample of healthy young adults. These findings could serve as a baseline-determination of CpGs that show natural epigenetic variation in a population of similar origin and structure. These CpGs also show a strong enrichment of significant hits in association analyses. The strong genetic component for CpGs comprising common epigenetic variation additionally suggests that a significant proportion of common DNA methylation variation may be shared across tissues. These findings could be of special relevance for studies of complex phenotypes, as in the case of neuropsychiatric disorders that often rely on proxy tissue.

### Contributors

AM and VV designed the study, performed experiments, analyzed data, and wrote the paper.
AP, CV and DQ designed the study.
FP, KS, PD, FH and TE performed the experiments.
AH, CV, VF, TE and AS contributed to the data analysis or results interpretation.
All authors contributed to paper writing.

### Funding

### Conflict of interest

All authors reported no biomedical financial interests or potential conflicts of interest.

### Acknowledgment

### Appendix A. Supplementary data

Supplementary data related to this article can be found at http://dx.doi.org/10.1016/j.jpsychires.2016.08.012.

### References

Altman, N., 2005. Replication, variation and normalisation in microarray experiments. Appl. Bioinforma. 4, 33—44. http://dx.doi.org/10.2165/00822942-200504010-00004.

Aryee, M.J., Jaffe, A.E., Corrada-Bravo, H., Ladd-Acosta, C., Feinberg, A.P., Hansen, K.D., Irizarry, R.A., 2014. Minfi: a flexible and comprehensive Bioconductor package for the analysis of Infinium DNA methylation microarrays. Bioinformatics 30, 1363—1369. http://dx.doi.org/10.1093/bioinformatics/btu049.

Assenov, Y., Müller, F., Lutsik, P., Walter, J., Lengauer, T., Bock, C., 2014. Comprehensive analysis of DNA methylation data with RnBeads. Nat. Methods. http://

dx.doi.org/10.1038/nmeth.3115.

Bell, J.T., Pai, A.A., Pickrell, J.K., Gaffney, D.J., Pique-Regi, R., Degner, J.F., Gilad, Y., Pritchard, J.K., 2011. DNA methylation patterns associate with genetic and gene expression variation in HapMap cell lines. Genome Biol. 12, R10. http://dx.doi.org/10.1186/gb-2011-12-1-r10.

Bibikova, M., Barnes, B., Tsan, C., Ho, V., Klotzle, B., Le, J.M., Delano, D., Zhang, L., Schroth, G.P., Gunderson, K.L., Fan, J.-B., Shen, R., 2011. High density DNA methylation array with single CpG site resolution. Genomics 98, 288—295. http://dx.doi.org/10.1016/j.ygeno.2011.07.007.

Bock, C., 2012. Analysing and interpreting DNA methylation data. Nat. Rev. Genet. 13, 705—719. http://dx.doi.org/10.1038/nrg3273.

Chen, Y., Lemire, M., Choufani, S., Butcher, D.T., Grafodatskaya, D., Zanke, B.W., Gallinger, S., Hudson, T.J., Weksberg, R., 2013. Discovery of cross-reactive probes and polymorphic CpGs in the illumina infinium HumanMethylation450 microarray. Epigenetics 8, 203—209. http://dx.doi.org/10.4161/epi.23470.

Davies, M.N., Volta, M., Pidsley, R., Lunnon, K., Dixit, A., Lovestone, S., Coarfa, C., Harris, R.A., Milosavljevic, A., Troakes, C., Al-Sarraj, S., Dobson, R., Schalkwyk, L.C., Mill, J., 2012. Functional annotation of the human brain methylome identifies tissue-specific epigenetic variation across brain and blood. Genome Biol. 13, R43. http://dx.doi.org/10.1186/gb-2012-13-6-r43.

Davis, S., Du, P., Bilke, S., Triche Jr., T., Bootwalla, M., 2014. Methylumi: Handle Illumina Methylation Data. R package version 2.12.0.

Du, P., Zhang, X., Huang, C.-C., Jafari, N., Kibbe, W.A., Hou, L., Lin, S.M., 2010. Comparison of Beta-value and M-value methods for quantifying methylation levels by microarray analysis. BMC Bioinforma. 11, 587. http://dx.doi.org/10.1186/1471-2105-11-587.

Flanagan, J.M., 2015. Epigenome-wide association studies (EWAS): past, present, and future. Methods Mol. Biol. 1238, 51—63. http://dx.doi.org/10.1007/978-1-4939-1804-1_3.

Fox, J., Weisberg, S., 2011. An R Companion to Applied Regression, Second. ed. Sage, Thousand Oaks {CA}.

Gibbs, J.R., van der Brug, M.P., Hernandez, D.G., Traynor, B.J., Nalls, M.A., Lai, S.-L., Arepalli, S., Dillman, A., Rafferty, I.P., Troncoso, J., Johnson, R., Zielke, H.R., Ferrucci, L., Longo, D.L., Cookson, M.R., Singleton, A.B., 2010. Abundant quantitative trait loci exist for DNA methylation and gene expression in human brain. PLoS Genet. 6, e1000952. http://dx.doi.org/10.1371/journal.pgen.1000952.

Guida, F., Sandanger, T.M., Castagné, R., Campanella, G., Polidoro, S., Palli, D., Krogh, V., Tumino, R., Sacerdote, C., Panico, S., Severi, G., Kyrtopoulos, S.A., Georgiadis, P., Vermeulen, R.C.H., Lund, E., Vineis, P., Chadeau-Hyam, M., 2015. Dynamics of smoking-induced genome-wide methylation changes with time since smoking cessation. Hum. Mol. Genet. 24, 2349—2359. http://dx.doi.org/10.1093/hmg/ddu751.

Heck, A., Fastenrath, M., Ackermann, S., Auschra, B., Bickel, H., Coynel, D., Gschwind, L., Jessen, F., Kaduszkiewicz, H., Maier, W., Milnik, A., Pentzek, M., Riedel-Heller, S.G., Ripke, S., Spalek, K., Sullivan, P., Vogler, C., Wagner, M., Weyerer, S., Wolfsgruber, S., de Quervain, D.J.-F., Papassotiropoulos, A., 2014. Converging genetic and functional brain imaging evidence links neuronal excitability to working memory, psychiatric disease, and brain activity. Neuron 81, 1203—1213. http://dx.doi.org/10.1016/j.neuron.2014.01.010.

Heyn, H., Moran, S., Hernando-Herraez, I., Sayols, S., Gomez, A., Sandoval, J., Monk, D., Hata, K., Marques-Bonet, T., Wang, L., Esteller, M., 2013. DNA methylation contributes to natural human variation. Genome Res. 23, 1363—1372. http://dx.doi.org/10.1101/gr.154187.112.

Horvath, S., 2013. DNA methylation age of human tissues and cell types. Genome Biol. 14, R115. http://dx.doi.org/10.1186/gb-2013-14-10-r115.

Horvath, S., Zhang, Y., Langfelder, P., Kahn, R.S., Boks, M.P., van Eijk, K., van den Berg, L.H., Ophoff, R.A., 2012. Aging effects on DNA methylation modules in human brain and blood tissue. Genome Biol. 13, R97. http://dx.doi.org/10.1186/gb-2012-13-10-r97.

Houseman, E.A., Accomando, W.P., Koestler, D.C., Christensen, B.C., Marsit, C.J., Nelson, H.H., Wiencke, J.K., Kelsey, K.T., 2012. DNA methylation arrays as surrogate measures of cell mixture distribution. BMC Bioinforma. 13, 86. http://dx.doi.org/10.1186/1471-2105-13-86.

Jaffe, A.E., Irizarry, R.A., 2014. Accounting for cellular heterogeneity is critical in epigenome-wide association studies. Genome Biol. 15, R31. http://dx.doi.org/10.1186/gb-2014-15-2-r31.

Johnson, W.E., Li, C., Rabinovic, A., 2007. Adjusting batch effects in microarray expression data using empirical Bayes methods. Biostatistics 8, 118—127. http://dx.doi.org/10.1093/biostatistics/kxj037.

Jones, P.A., 2012. Functions of DNA methylation: islands, start sites, gene bodies and beyond. Nat. Rev. Genet. 13, 484—492. http://dx.doi.org/10.1038/nrg3230.

Leek, J.T., Johnson, W.E., Parker, H.S., Jaffe, A.E., Storey, J.D., 2012. The sva package for removing batch effects and other unwanted variation in high-throughput experiments. Bioinformatics 28, 882—883. http://dx.doi.org/10.1093/bioinformatics/bts034.

Lokk, K., Modhukur, V., Rajashekar, B., Märtens, K., Mägi, R., Kolde, R., Kolt Ina, M., Nilsson, T.K., Vilo, J., Salumets, A., Tõnisson, N., 2014. DNA methylome profiling of human tissues identifies global and tissue-specific methylation patterns. Genome Biol. 15, R54. http://dx.doi.org/10.1186/gb-2014-15-4-r54.

Maksimovic, J., Gordon, L., Oshlack, A., 2012. SWAN: subset-quantile within array

normalization for illumina infinium HumanMethylation450 BeadChips. Genome Biol. 13, R44. http://dx.doi.org/10.1186/gb-2012-13-6-r44.

McClay, J.L., Shabalin, A.A., Dozmorov, M.G., Adkins, D.E., Kumar, G., Nerella, S., Clark, S.L., Bergen, S.E., Hultman, C.M., Magnusson, P.K.E., Sullivan, P.F., Aberg, K.A., van den Oord, E.J.C.G., 2015. High density methylation QTL analysis in human blood via next-generation sequencing of the methylated genomic DNA fraction. Genome Biol. 16, 291. http://dx.doi.org/10.1186/s13059-015-0842-7.

Petronis, A., 2010. Epigenetics as a unifying principle in the aetiology of complex traits and diseases. Nature 465, 721—727. http://dx.doi.org/10.1038/nature09230.

Pidsley, R., Viana, J., Hannon, E., Spiers, H.H., Troakes, C., Al-Saraj, S., Mechawar, N., Turecki, G., Schalkwyk, L.C., Bray, N.J., Mill, J., 2014. Methylomic profiling of human brain tissue supports a neurodevelopmental origin for schizophrenia. Genome Biol. 15, 483. http://dx.doi.org/10.1186/PREACCEPT-1621721621132088.

Price, M.E., Cotton, A.M., Lam, L.L., Farré, P., Emberly, E., Brown, C.J., Robinson, W.P., Kobor, M.S., 2013. Additional annotation enhances potential for biologically-relevant analysis of the Illumina Infinium HumanMethylation450 BeadChip array. Epigenet. Chromatin 6, 4. http://dx.doi.org/10.1186/1756-8935-6-4.

Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M.A.R., Bender, D., Maller, J., Sklar, P., de Bakker, P.I.W., Daly, M.J., Sham, P.C., 2007. PLINK: a tool set for whole-genome association and population-based linkage analyses. Am. J. Hum. Genet. 81, 559—575. http://dx.doi.org/10.1086/519795.

Rakyan, V.K., Down, T.A., Balding, D.J., Beck, S., 2011. Epigenome-wide association studies for common human diseases. Nat. Rev. Genet. 12, 529—541. http://dx.doi.org/10.1038/nrg3000.

R Development Core Team, 2012. R: a Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna.

Sandoval, J., Heyn, H., Moran, S., Serra-Musach, J., Pujana, M.A., Bibikova, M., Esteller, M., 2011. Validation of a DNA methylation microarray for 450,000 CpG sites in the human genome. Epigenetics 6, 692—702.

Schalkwyk, L.C., Meaburn, E.L., Smith, R., Dempster, E.L., Jeffries, A.R., Davies, M.N., Plomin, R., Mill, J., 2010. Allelic skewing of DNA methylation is widespread across the genome. Am. J. Hum. Genet. 86, 196—212. http://dx.doi.org/10.1016/j.ajhg.2010.01.014.

Shenker, N.S., Polidoro, S., van Veldhoven, K., Sacerdote, C., Ricceri, F., Birrell, M.A., Belvisi, M.G., Brown, R., Vineis, P., Flanagan, J.M., 2013. Epigenome-wide association study in the European prospective investigation into Cancer and nutrition (EPIC-Turin) identifies novel genetic loci associated with smoking. Hum. Mol. Genet. 22, 843—851. http://dx.doi.org/10.1093/hmg/dds488.

Shoemaker, R., Deng, J., Wang, W., Zhang, K., 2010. Allele-specific methylation is prevalent and is contributed by CpG-SNPs in the human genome. Genome Res. 20, 883—889. http://dx.doi.org/10.1101/gr.104695.109.

Smith, A.K., Kilaru, V., Kocak, M., Almli, L.M., Mercer, K.B., Ressler, K.J., Tylavsky, F.A., Conneely, K.N., 2014. Methylation quantitative trait loci (meQTLs) are consistently detected across ancestry, developmental stage, and tissue type. BMC Genom. 15, 145. http://dx.doi.org/10.1186/1471-2164-15-145.

Spalek, K., Fastenrath, M., Ackermann, S., Auschra, B., Coynel, D., Frey, J., Gschwind, L., Hartmann, F., van der Maarel, N., Papassotiropoulos, A., de Quervain, D., Milnik, A., 2015. Sex-dependent dissociation between emotional appraisal and memory: a large-scale behavioral and fMRI study. J. Neurosci. 35, 920—935. http://dx.doi.org/10.1523/JNEUROSCI.2384-14.2015.

Stacklies, W., Redestig, H., Scholz, M., Walther, D., Selbig, J., 2007. pcaMethods — a Bioconductor package providing PCA methods for incomplete data. Bioinformatics 23, 1164—1167.

Tan, Q., Christiansen, L., von Bornemann Hjelmborg, J., Christensen, K., 2015. Twin methodology in epigenetic studies. J. Exp. Biol. 218, 134—139. http://dx.doi.org/10.1242/jeb.107151.

Touleimat, N., Tost, J., 2012. Complete pipeline for Infinium($^{®}$) Human Methylation 450K BeadChip data processing using subset quantile normalization for accurate DNA methylation estimation. Epigenomics 4, 325—341. http://dx.doi.org/10.2217/epi.12.21.

Tylee, D.S., Kawaguchi, D.M., Glatt, S.J., 2013. On the outside, looking in: a review and evaluation of the comparability of blood and brain "-omes". Am. J. Med. Genet. B. Neuropsychiatr. Genet. 162B, 595—603. http://dx.doi.org/10.1002/ajmg.b.32150.

Xu, H., Wang, F., Liu, Y., Yu, Y., Gelernter, J., Zhang, H., 2014. Sex-biased methylome and transcriptome in human prefrontal cortex. Hum. Mol. Genet. 23, 1260—1270. http://dx.doi.org/10.1093/hmg/ddt516.

Xu, Z., Taylor, J.A., 2014. Genome-wide age-related DNA methylation changes in blood and other tissues relate to histone modification, expression and cancer. Carcinogenesis 35, 356—364. http://dx.doi.org/10.1093/carcin/bgt391.

Zhang, D., Cheng, L., Badner, J.A., Chen, C., Chen, Q., Luo, W., Craig, D.W., Redman, M., Gershon, E.S., Liu, C., 2010. Genetic control of individual differences in gene-specific methylation in human brain. Am. J. Hum. Genet. 86, 411—419. http://dx.doi.org/10.1016/j.ajhg.2010.02.005.

Ziller, M.J., Müller, F., Liao, J., Zhang, Y., Gu, H., Bock, C., Boyle, P., Epstein, C.B., Bernstein, B.E., Lengauer, T., Gnirke, A., Meissner, A., 2011. Genomic distribution and inter-sample variation of non-CpG methylation across human cell types. PLoS Genet. 7, e1002389. http://dx.doi.org/10.1371/journal.pgen.1002389.