



# The correlation of methylation levels measured using Illumina 450K and EPIC BeadChips in blood samples

Mark W Logue<sup>\*,1,2,3,4</sup>, Alicia K Smith<sup>5,6</sup>, Erika J Wolf<sup>1,2</sup>, Hannah Maniates<sup>1</sup>, Annjanette Stone<sup>7</sup>, Steven A Schichman<sup>7</sup>, Regina E McGlinchey<sup>8,9</sup>, William Milberg<sup>8,9</sup> & Mark W Miller<sup>1,2</sup>

<sup>1</sup>National Center for PTSD, VA Boston Healthcare System, USA

<sup>2</sup>Department of Psychiatry, Boston University School of Medicine, USA

<sup>3</sup>Biomedical Genetics, Boston University School of Medicine, USA

<sup>4</sup>Biostatistics, Boston University School of Public Health, USA

<sup>5</sup>Department of Gynecology & Obstetrics, Emory University, USA

<sup>6</sup>Department of Psychiatry & Behavioral Sciences, Emory University, USA

<sup>7</sup>Pharmacogenomics Analysis Laboratory, Research Service, Central Arkansas Veterans Healthcare System, USA

<sup>8</sup>Geriatric Research Educational & Clinical Center & Translational Research Center for TBI & Stress Disorders, VA Boston Healthcare System, USA

<sup>9</sup>Department of Psychiatry, Harvard Medical School, USA

\* Author for correspondence: Tel.: +1 857 364 5665; Fax: +1 857 364 4501; [mark.logue@va.gov](mailto:mark.logue@va.gov)

**Aim:** We examined concordance of methylation levels across the Illumina Infinium HumanMethylation450 BeadChip and the Infinium MethylationEPIC BeadChip. **Methods:** We computed the correlation for 145 whole blood DNA samples at each of the 422,524 CpG sites measured by both chips. **Results:** The correlation at some sites was high (up to  $r = 0.95$ ), but many sites had low correlation (55% had  $r < 0.20$ ). The low correspondence between 450K and EPIC measured methylation values at many loci was largely due to the low variability in methylation values for the majority of the CpG sites in blood. **Conclusion:** Filtering out probes based on the observed correlation or low variability may increase reproducibility of BeadChip-based epidemiological studies.

First draft submitted: 23 June 2017; Accepted for publication: 2 August 2017; Published online: 15 August 2017

**Keywords:** DNA methylation • HumanMethylation450 BeadChip • Illumina Infinium • Infinium MethylationEPIC BeadChip • intraclass correlation

The Infinium MethylationEPIC BeadChip (EPIC chip) was released in December 2015. The EPIC chip yields an estimate of the proportion of methylated DNA (denoted  $\beta$ ) for more than 850,000 CpG sites. These cover 99% of RefSeq genes and 95% of CpG islands with an average number of six probes per island. The EPIC BeadChip includes 90% of the sites from the previous-generation chip: the Illumina Infinium HumanMethylation450 BeadChip (450k chip). Several studies have appeared in the literature assessing the performance of the EPIC chip in different tissues [1–3]. These studies examined methylation of DNA samples from oncological cell lines (e.g., LNCaP and PrEC), or in pairs of tissues preserved using a variety of methods (e.g., fresh frozen vs formalin-fixed paraffin embedded) and compared methylation values of matched samples assessed on both 450k and EPIC chips. They calculated what we will here refer to as the ‘overall correlation’. That is, they ran a single DNA sample twice, or a matched pair of samples preserved using different methods, and then looked at the matched pairs of data represented by all CpG sites on chip 1 and chip 2 and evaluated the likelihood that the methylation pattern observed for that same sample on the two chips were similar. These studies reported overall correlations of matched samples run on the EPIC and 450k chip that were quite high ( $r > 0.90$  for all samples assessed). Pidsley *et al.* [2] additionally examined the ability of the platforms to detect differentially methylated probes between three cancer-associated fibroblast and three nonmalignant fibroblast cell lines. They found that 94% of the overlapping probes that were differentially methylated at the  $p < 0.01$  significance level in the EPIC chip data were also differentially methylated

at the  $p < 0.01$  significance level in the 450k data. Estimated effect sizes for differentially methylated sites were also similar.

These reliability studies suggest that patterns of methylation across the two chips were consistent for oncological samples, based on the overall methylation patterns. However, these studies did not directly examine the correlation of methylation at specific CpG sites for a collection of blood-based DNA samples, as is often used in studies examining methylation differences between individuals rather than tissues, for example, large scale case–control studies (for brevity we will dub these epidemiological studies). Here, we investigated the site-specific correlation in 145 DNA samples with methylation evaluated on both the 450k and the EPIC BeadChip. This gives us an indication of how well each CpG site would perform in an epidemiological-study setting, and the degree to which associations observed with a specific CpG site might replicate across chips. We found that, as previously reported in oncological tissue, overall correlation for methylation values from a single blood-derived DNA sample assessed on both chips were quite high. However, for the majority of measured CpG sites, the correlation of methylation values across the cohort of DNA samples was quite low. Furthermore, we present evidence that the weak correspondence in methylation values at most sites is largely attributable to the fact that a large proportion of sites measured by the Illumina BeadChips have relatively low variation in human blood samples. That is, the majority of CpGs are either almost entirely all methylated or all unmethylated in human blood samples. For CpG sites that have a higher range of observable variability, which are arguably the most interesting sites to study, the correlation between probes on the two chips was substantially higher.

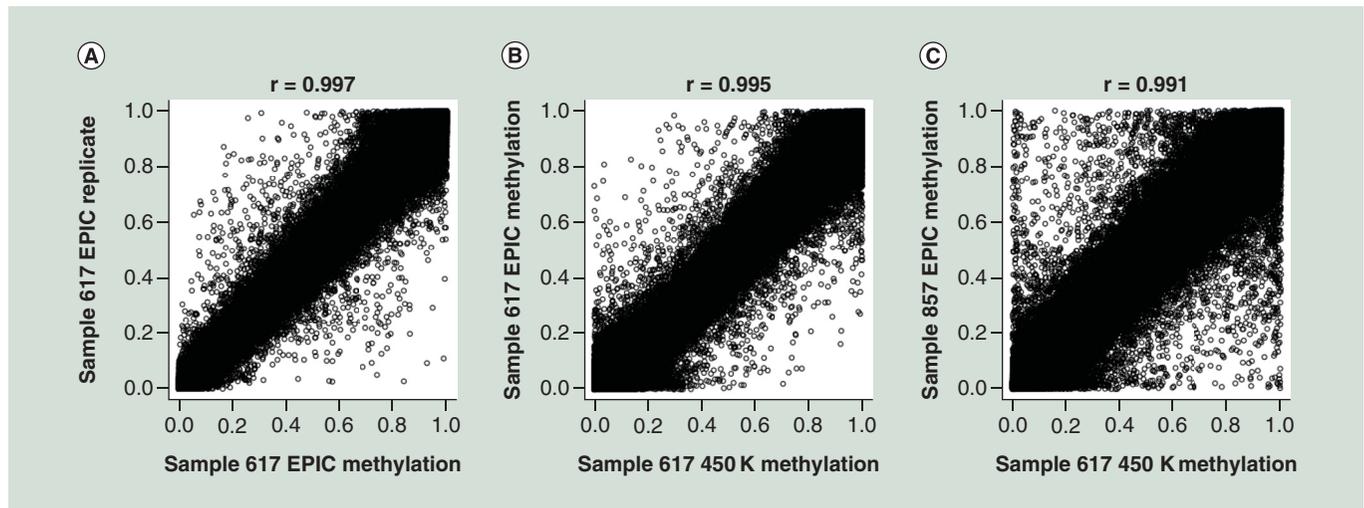
## Methods

### Samples

The samples used in this study were from the Translational Research Center for TBI and Stress Disorders (TRACTS), a VA Rehabilitation Research and Development National Center for TBI Research at VA Boston Healthcare System. We have previously published on a subset of TRACTS with 450k BeadChip methylation and genotype data for 283 subjects [4]. Here, we have obtained methylation data for 145 of these previously assessed DNA samples using EPIC BeadChips. This research was performed under the oversight of appropriate human-subjects review boards and informed consent was obtained from all subjects at the time of inclusion to the study.

### Data cleaning & processing

Both the 450k and EPIC chips were run at the Pharmacogenomics Analysis Laboratory at Central Arkansas Veterans Healthcare System according to manufacturer's protocols. The quality control (QC) process for the 450k data, including missing rates and sample fails, is described in detail elsewhere [4]. It is the same process that we will describe here for the EPIC data. Briefly, we used a pipeline developed by the PGC-PTSD Epigenetics workgroup for 450k data [5]. A total of 524 DNA samples were run on the EPIC chip, 145 of which overlapped with the 450k chip data. Individual-level background-corrected probe data and idat files were output from GenomeStudio. Sample identity across the 450k and EPIC chips were confirmed using the SNP probes included on both chips for this purpose. Subsequent cleaning was performed within the CpGassoc package and the ChAMP package in R [6]. Individual probes failing to meet a detection  $p$ -value threshold of 0.001 were set to missing. One chip had seven out of eight failed samples ( $>10\%$  missing) and was rerun. After replacing this chip, no samples had more than 10% missing data. No samples had intensity of less than 50% of the experiment-wide mean or intensity less than 2000 arbitrary units. Probes with more than 10% missing data were excluded from the analysis ( $n = 2093$ ). We additionally excluded EPIC probes that can cross-hybridize between autosomes and sex chromosomes [7] ( $n = 44,132$ ) and the 977 'underperforming' EPIC probes included in Illumina Product Quality Notification PQN0223 dated 19 April 2017. Normalization was performed using the  $\beta$  mixture quantile dilation method [8] as implemented in the watermelon [9,10] R package. Batch and chip effect adjustments were implemented using an empirical Bayes batch-correction method (ComBat) [11] as implemented in the Bioconductor sva package [12]. As differential white-cell counts were not available for our samples, cell counts were estimated from the methylation data itself using the R minfi [13] package. Post QC, there were 453,747 probes available for analysis in the 450k data and 819,942 probes available for analysis in the EPIC chip data. Of these, there were 422,524 CpG sites that were interrogated by both chips. For the EPIC chip data, technical replicates (the same DNA sample run more than one-time using the EPIC chip) were available for 11 samples. Of these, ten were run twice and one was run three times. Comparisons of the number of beads measuring each CpG site (nbeads) and missing rates observed in the 450k BeadChip and the EPIC BeadChip data are presented in the Supplementary Materials.



**Figure 1.** Example of the overall correlation between chips for representative samples. (A) Methylation values for a sample (ID 617) run twice on the EPIC BeadChip. (B) Methylation values for the same sample run on an EPIC and a 450k BeadChip. (C) The correlation between samples from two different individuals, one run with the 450k BeadChip (ID 617) chip and the other run with the EPIC BeadChip (ID 857).

### Analysis methods

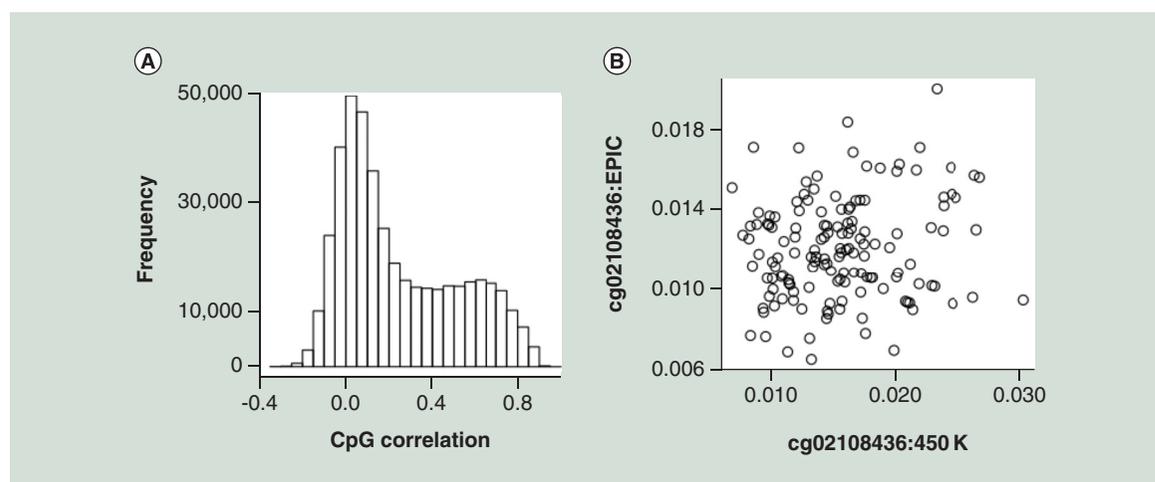
Data analysis, including the calculation of correlations and plotting of figures was performed in R [6]. Pearson correlations were computed using R's `cor` function. Intraclass correlations were computed using the R `ICCest` function from the `ICC` library. Linear models of the correlation for probes were run using R's `lm` function. We assessed the impact of average nbeads across samples, missing rate (number of probes to fail a detection p-value threshold of  $p > 0.001$ ) and the range of  $\beta$  values ( $\max \beta - \min \beta$ ) on the correlation observed at each CpG site. Predictors were standardized (mean = 0, standard deviation = 1) prior to analysis. As they are based on different chemistries, separate models were fit on type I probes, where two probe types are used per site (A and B) and type II probes, which only have a single probe type per site. For simplicity's sake, we will only discuss the results for the linear model of type II probes in the main body of the manuscript. Results for the type I probes are similar and are presented in the Supplementary Materials.

We additionally examined the methylation data and its correspondence to externally measurable sample quantities. First, we looked at cg10636246, which had been previously associated with plasma CRP levels [14]. Association between plasma CRP and cg10636246 was assessed using a linear model that included covariates for age, estimated cell counts, sex and principal components for ancestry. Additionally, we examined methylation age estimates as computed by the Horvath [15] and Hannum [16] methods and their correlation with chronological age in both the 450k and EPIC BeadChip data. These analyses are presented in detail in the Supplementary Materials.

## Results

### Assessing the overall correlation between BeadChips

Similar to what was observed in prior reliability studies of oncological tissues [1–3], for technical replicates of the EPIC chip (the same DNA sample run twice), the overall correlation across 819,942 CpG sites was quite high (correlation from  $r = 0.994$  to  $r = 0.998$  across all DNA sample pairs), with a median overall correlation of 0.997. See Figure 1A for a scatter plot of methylation for a technical replicate for a representative DNA sample run on an EPIC chip. The overall correlation for a single sample run on both the 450k and EPIC BeadChips across the 422,524 CpG sites assessed on both BeadChips was also quite high. For the 145 samples, overall correlation values between chips ranged from  $r = 0.991$  to  $r = 0.996$ , with a median of  $r = 0.995$  (Figure 1B). However, it must be noted that these overall correlations were high for any two samples from TRACTS. In fact, in 1000 randomly sampled pairs, overall correlation between any two samples ranged from  $r = 0.985$  to  $r = 0.994$ , with a median of  $r = 0.991$  (Figure 1C). Hence overall correlation across all CpG sites can be expected to be high for any two blood samples, regardless of whether or not they were from the same subject.



**Figure 2.** CpG Correlation between 450K and EPIC BeadChips. **(A)** Histogram of correlation for 422,524 CpG sites assessed by both the 450k and EPIC BeadChip based on 145 samples measured with both chips; and **(B)** a scatterplot of EPIC and 450k measured methylation values of 145 samples for a representative CpG site (cg02108436) with a correlation of  $r = 0.15$ , the median correlation across all probes.

**Table 1.** Model of correlation for CpG sites measured by type II probes.

Model term <sup>†</sup>	Estimate	SE	T	p-value
(Intercept)	0.297	0.0003	954.03	<2e-16
Nbeads 450	0.0119	0.00029	40.45	<2e-16
Nbeads EPIC	0.016	0.00030	52.60	<2e-16
N missing 450	-0.020	0.00038	-52.50	<2e-16
N missing EPIC	-0.00065	0.00038	-1.70	0.09
Range 450	0.25	0.0010	250.96	<2e-16
(Range 450) <sup>2</sup>	-0.034	0.00024	-141.74	<2e-16
Range EPIC	0.052	0.00092	56.44	<2e-16
(Range EPIC) <sup>2</sup>	-0.00031	0.00022	-1.42	0.15

Residual standard error: 0.16 on 311,293 degrees of freedom.

Multiple R-squared: 0.63, adjusted R-squared: 0.63.

F-statistic: 6.68e+04 on 8 and 311,293 DF, p-value: <2.2e-16.

<sup>†</sup>All predictors standardized (mean 0, SD = 1) prior to modeling. Nbeads is average number of beads across samples. N missing represents the proportion of missing values per probe. Range:  $\max(\beta) - \min(\beta)$ .

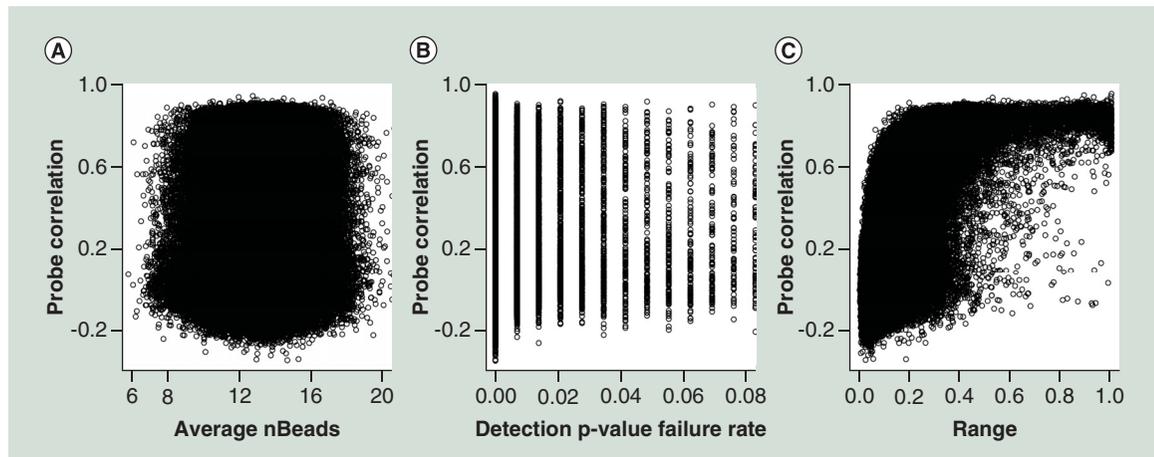
DF: Degrees of freedom; SD: Standard deviation; T: Student's T statistic.

### Assessing CpG site-specific correlations

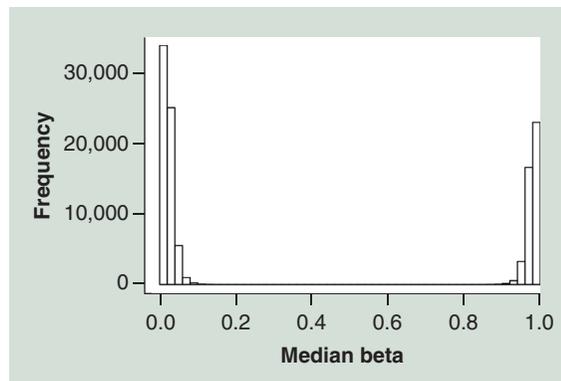
For studies attempting to replicate epidemiological findings from the 450k BeadChip with methylation data generated using the EPIC array, or for studies interested in doing a joint analysis of 450K and EPIC data, the degree to which methylation values for an individual CpG site interrogated with the 450k correlates with methylation values for that same CpG site interrogated with the EPIC array is far more informative. The correlation of CpGs measured on the 450k and EPIC arrays across 145 subjects varied widely. Figure 2A presents a histogram of the correlation values for CpG sites assessed on both chips. All of the correlation values, as well as QC information from each CpG site, are presented in Supplementary Table 1. For the 422,254 probes on both chips, the median and mean correlation were 0.15 and 0.25, respectively ( $r = -0.34$  to 0.95). See Figure 2B for a scatterplot from a representative CpG site. Of the sites measured by both probes, 234,727 (55.6%) had correlation values <0.20, 95,950 (23%) had correlation values >0.50 and 10,987 (2.6%) had correlations >0.80.

### Examining the causes of low CpG correlation

We then investigated the impact of nbeads, missing rate and range of observed  $\beta$  values on the correlations for each CpG site (Figure 3 A–C) to determine what factors influence CpG correlation. Table 1 presents the results of a linear model of these predictors and methylation values for CpG sites assessed with type II probes, with a quadratic



**Figure 3.** The correlation for CpG sites assessed by both 450k and EPIC BeadChips (type II probes only). Presented as a function of (A) Bead Count per CpG for the 450k data; (B) the proportion of missing values for the 450k; and (C) the range of the observed 450k B values (max  $\beta$ –min  $\beta$ ). The range of values is trimmed to enhance detail.



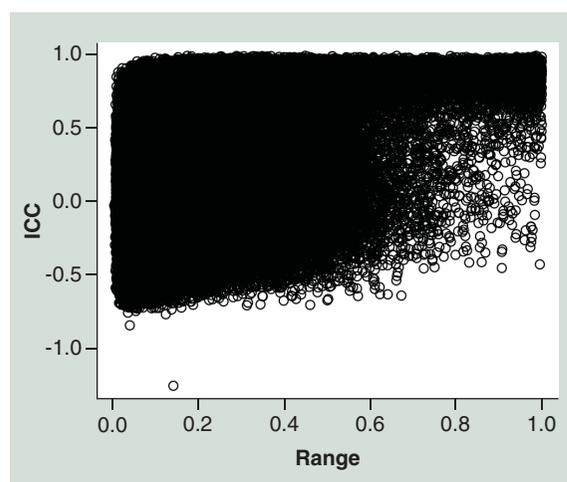
**Figure 4.** A histogram of the median  $\beta$  value for CpG sites with  $\beta$  ranges  $<0.05$  demonstrating that these probes are nearly completely methylated or completely unmethylated.

term for range included due to the curvature apparent in Figure 3C. All predictors except the number of missing probes and the range<sup>2</sup> term for the EPIC data were significant. We note that this may be due to multicollinearity, as the range and the missing rate were correlated across chips ( $r = 0.91$  and  $r = 0.63$ , respectively). There was a strong positive association between the  $\beta$  range variables and CpG correlation with a quadratic term indicating downward curvature. A comparison of the multiple  $r^2 = 0.63$  for the model presented in Table 1 to the multiple  $r^2 = 0.0048$  value for a model that excludes the  $\beta$  range variables (presented in Supplementary Table 2) indicates that the majority of the variance in correlation for the CpG sites is explained by probe variability, such that sites with low variation tend to have lower correlation. The results for the type I probes are similar, and are presented in the Supplementary Tables 3 & 4.

We further examined the low-variation sites. For CpG sites with a 450k chip methylation range of less than 0.05 ( $n = 109,279$ ), the median correlation was  $r = 0.031$  and 94.64% of sites had correlations less than  $r = 0.20$ . The vast majority of these low variance sites were nearly completely methylated or nearly completely unmethylated (Figure 4). That is, 95.81% of these sites had a median methylation value  $<0.05$  or  $>0.95$ . The correlation was much higher for CpGs with a greater range of  $\beta$  values. For the CpGs, which had a  $\beta$  range of  $>0.20$  ( $n = 61,248$ ), the median correlation was  $r = 0.71$  and 93.36% had correlations greater than 0.20.

#### Examining the ICC for EPIC chip data

Finally, we examined the effect of  $\beta$  range on the ICC of CpG sites based on 11 DNA samples measured repeatedly using the EPIC chip. This was done to determine if the reliability of low-variation sites would be similarly low for those analyzing data from the EPIC chip alone. Supplementary Table 5 presents the ICC values for 819,942 CpG sites along with  $\beta$  range, nbeads and missing rate for each CpG site. Figure 5 presents the ICC for 819,942



**Figure 5.** Intraclass correlation as a function of the range of observed  $\beta$  values for EPIC-measured CpGs.

CpG sites as a function of the  $\beta$  range. The pattern is more diffuse, but broadly matches the pattern observed in Figure 3C. In a linear model of EPIC type II probes, both range and range<sup>2</sup> were significantly associated with ICC ( $p < 2 \times 10^{-16}$ ) (Supplementary Table 6). The  $r^2$  of the model of ICC including nBeads, N missing and range and range<sup>2</sup> was 0.24 and all predictors were significant (all  $p < 2 \times 10^{-16}$ ). For the corresponding models excluding the  $\beta$  range variables, the  $r^2$  was 0.022 (Supplementary Table 7). While the  $r^2$  for the models of ICC with the range variables (Supplementary Table 6) is smaller than the  $r^2$  of the models of CpG correlation (as presented in Table 1), this may be due to the higher standard errors of ICC estimates computed from only 11 samples. Results for the type I probes were similar (Supplementary Tables 8 & 9).

#### Two examples of comparability of performance between the 450K & EPIC chip

While there are many CpG sites with low variability and hence low correlation, this should not necessarily be interpreted as unreliability of the chips in general or a general lack of validity of results that have been previously reported based on BeadChip technology. To emphasize this fact and to demonstrate that the low correlation observed at many probes was not an artifact of mishandling of our data, we examined the relationship between EPIC chip measured methylation values and several externally measurable quantities available for our cohort. First, an epigenome-wide association study of CRP that used the 450k BeadChip [14] reported that methylation of cg10636246 was associated with (serum/plasma) CRP levels. Here, we find the cross-chip correlation for cg10636246 was  $r = 0.74$  and correlation between the EPIC chip measurement of cg10636246 and CRP levels was similar to that reported with the 450K chip (see Supplementary Results for details). We also examined measures of DNA methylation age [15,16] and their relationship to chronological age. Again, performance of the EPIC BeadChip was similar to the 450k chip (see Supplementary Results for details). Hence, we stress that even though many individual sites display low correlation between chips, there are also many instances in which the EPIC array performs similarly to the 450k chip.

#### Discussion

We have shown that the concordance of individual level methylation data for CpG sites assessed on both the EPIC BeadChip and 450k chip is quite low for many CpG sites. We note that data for both chips passed QC filters based on a pipeline developed for use by a genome-wide methylation consortium. We showed that the correlation at a given CpG site is largely determined by the amount of true variability in DNA methylation levels across subjects, such that sites with low variation tend to also have low correlation and low ICCs. The power to detect an association is a function of the reliability of measurement [17,18], and hence these low correlation sites will have difficulty replicating. This represents a substantial proportion of sites measured by these BeadChips, for example, 25.8% of overlapping sites have a range of  $\beta$  values less than 0.05, and within these sites, 95.8% have CpG correlations of less than 0.20. When the overall variability of a CpG site in a given tissue is low, this may be a signal that the ratio of the true variability relative to the measurement error is low, and the correlation can suffer as a consequence. This is true whether we are comparing 450k data to EPIC data, or whether we are investigating intraclass correlations within EPIC chip data. However, this should not be interpreted as a blanket excoriation of the

performance of either the EPIC or the 450k chip, or to necessarily impugn the results of any study which used these chips. As our examples of the CRP-associated probe and the methylation age calculation demonstrate, there are clear cases where the data from these chips can be consistently associated with independently measured quantities, and the performance for these purposes can be similar across chips. We showed that the low observed correlation at many sites is a technical artifact of low variability at CpG sites that are either nearly entirely methylated or entirely unmethylated in blood. This is a consequence of chip design, as the 450k and EPIC chips are designed to assess methylation broadly, regardless of tissue type, and not specifically tailored to epidemiological studies of whole blood samples, and hence did not specifically screen for CpG sites, which were variable in blood in the population. This represents an opportunity for manufacturers. A BeadChip focusing on sites that are variable in blood (or saliva) tailored for use in epidemiological studies could increase the number of useful probes while lowering the number of sites assessed overall, potentially producing a 'human whole blood methylation chip' with higher power but a lower price point. The full complement of probes present on the EPIC chip would still be useful for researchers examining different tissues and for use in oncological studies.

There are several limitations that should be kept in mind while interpreting these results. We had a reasonably large sample with which to examine the correlation for CpG sites measured on both chips, but much less data to examine the intraclass correlation of technical replicates performed on the EPIC chip, and no technical replicates for DNA samples assessed with the 450k chip to use as a comparison. Also, patterns of methylation can vary across tissues, and sites which have high variability in one tissue may have little or no variability in the other. As our data is from whole blood, the ranges and correlations observed here should only be used to guide studies of whole blood samples. Also, as we do not have independent measurements of methylation using another methodology (e.g., pyrosequencing), we are not able to determine whether the EPIC or the 450k is more representative of the true underlying methylation values for a CpG site.

In conclusion, we urge caution when interpreting results or attempting to replicate findings from CpG sites with low observed variability. This is not an entirely new concept. Previous studies have suggested filtering out low variation probes (see for example [19]). In oncological research, validation studies have led to a customary minimum detectible mean  $\beta$  difference [20–22]. However, in epidemiological studies, there is currently no corresponding convention and there is not a general practice of filtering sites for variability (see e.g. [23–27]). Researchers conducting a blood-based EWAS study may wish to consider whether or not it is worthwhile to analyze low variability sites, or if it is better to focus on sites with either high variability or high estimated correlation. It might be a good strategy for an epidemiological EWAS or a candidate-gene methylation study to drop all CpG sites with a range of  $\beta$  values  $<0.1$ . Technical replicates could also be used to filter out CpG sites with reliability problems, for example, by excluding all sites with an estimated ICC of  $<0.50$ . However, the relatively high noise pattern observed in Figure 4 suggests that more than 11 technical replicates would be necessary to generate a clean distinction between high and low performing CpG sites. Filtering based on  $\beta$  range (as suggested in [19]), estimated correlation for overlapping probes (based on Supplementary Table 1) or estimated ICC (as suggested in [28]) would allow for greater power through a reduction of the multiple testing burden as well as increase the chance that the probe will replicate in another dataset. It is also worth noting that even in situations where the correlation is high at a particular CpG site, there may be batch effects, so that EPIC and 450k data should not be analyzed together without an appropriate adjustment.

Going forward, it is worth noting that even for epidemiological studies exclusively using EPIC chips, it may be worthwhile to exclude low variance sites. In addition, the results here emphasize the importance of validation of methylation values using other technologies such as bisulfite methylation sequencing. Researchers should strongly consider validation of top reported sites using another independent technology. This may be especially important to validate the performance of sites with a relatively low observed range of values. For CpG sites, which are nearly completely methylated or unmethylated, double sequencing, or assessing each site twice, may be necessary to improve reliability beyond the levels achieved by the methylation chips. The estimated correlation between sites measured on both chips presented here represent important information on the reliability and reproducibility of epidemiological studies of methylation in blood-derived DNA samples conducted using BeadChips in general, and those transitioning from Illumina 450k to EPIC BeadChips in particular.

#### Disclaimer

The views expressed in this article are those of the authors and do not necessarily reflect the position or policy of the Department of Veterans Affairs, the Department of Defense or the US government.

#### Financial & competing interests disclosure

MW Miller owns stock in Illumina, Inc. This work was funded by I01BX003477, a VA BLR&D grant to MW Logue, 1R03AG051877 and 1I01CX001276–01A2 to EJ Wolf, R21MH102834 to MW Miller, 1R01MH108826 to AK Smith/MW Logue/C Nievergelt/M Uddin and the Translational Research Center for TBI and Stress Disorders (TRACTS), a VA Rehabilitation Research and Development (RR&D) Traumatic Brain Injury Center of Excellence (B9254-C) at VA Boston Healthcare System. The authors have no other relevant affiliations or financial involvement with any organization or entity with a financial interest in or financial conflict with the subject matter or materials discussed in the manuscript apart from those disclosed.

No writing assistance was utilized in the production of this manuscript.

#### Ethical conduct of research

The authors state that they have obtained appropriate institutional review board approval or have followed the principles outlined in the Declaration of Helsinki for all human or animal experimental investigations. In addition, for investigations involving human subjects, informed consent has been obtained from the participants involved.

#### Open access

This work is licensed under the Attribution-NonCommercial-NoDerivatives 4.0 Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>

#### Summary points

- We obtained measures of DNA methylation for whole blood-derived DNA samples from 145 subjects using both the Illumina Infinium HumanMethylation450 BeadChip and the Infinium MethylationEPIC BeadChip.
- We examined the correlation of methylation values measured on the same sample using both chips for all overlapping CpG sites (overall correlation; one correlation estimate per sample), as well as the correlation of 450k and EPIC methylation values at each of the overlapping CpG sites individually (CpG correlation; one correlation estimate per CpG site).
- As previously observed in earlier studies of oncological samples, the overall correlation for a single sample measured on both chips was quite high ( $r > 0.99$  for each of the 145 samples).
- When we examined the CpG site correlation, we found that some CpGs had high correlation, with 23% having  $r > 0.50$ , but many had low observed correlation, with 55% having  $r < 0.20$ .
- We found that CpG site correlation is largely a function of the range of observed  $\beta$  values (the proportion of methylated DNA), such that sites with low variability in blood samples tended to have low correlation. This is likely due to the lower S/N ratio for these probes.
- Many high variability probes have high correlation. We present several examples of previously reported associations where EPIC chip performance is similar to that of the 450K chip.
- Researchers performing epigenome-wide association studies should consider analyzing only probes with a higher range of variability, regardless of the chip used to measure methylation.
- This would allow for lower multiple-testing corrected significance thresholds as well as increase the likelihood that observed associations are replicable.

#### References

Papers of special note have been highlighted as: • of interest; •• of considerable interest

- 1 Moran S, Arribas C, Esteller M. Validation of a DNA methylation microarray for 850,000 CpG sites of the human genome enriched in enhancer sequences. *Epigenomics* 8(3), 389–399 (2016).
  - **Validating the EPIC BeadChips in different tissues by assessing the overall correlation.**
- 2 Pidsley R, Zotenko E, Peters TJ *et al.* Critical evaluation of the Illumina MethylationEPIC BeadChip microarray for whole-genome DNA methylation profiling. *Genome Biol.* 17(1), 208 (2016).
  - **Validating the EPIC BeadChips in different tissues by assessing the overall correlation.**
- 3 Kling T, Wenger A, Beck S, Caren H. Validation of the MethylationEPIC BeadChip for fresh-frozen and formalin-fixed paraffin-embedded tumours. *Clin. Epigenetics* 9, 33 (2017).
  - **Validating the EPIC BeadChips in different tissues by assessing the overall correlation.**
- 4 Sadeh N, Spielberg JM, Logue MW *et al.* SKA2 methylation is associated with decreased prefrontal cortical thickness and greater PTSD severity among trauma-exposed veterans. *Mol. Psychiatry* 21(3), 357–363 (2016).
- 5 Ratanatharathorn A, Boks M, Aiello A *et al.* Epigenome-wide association of PTSD from heterogeneous cohorts with a common multi-site analysis pipeline. *Am. J. Med. Genet. B Neuropsychiatr. Genet.* doi:10.1002/ajmg.b.32568 (2017) (Epub ahead of print).

- **Presents the development of a new pipeline for the analysis of BeadChip data. This pipeline was used for both 450K and EPIC BeadChip data.**
- 6 R Development Core Team. R: a language and environment for statistical computing. (2008). [www.r-project.org/](http://www.r-project.org/)
- 7 Chen YA, Lemire M, Choufani S *et al.* Discovery of cross-reactive probes and polymorphic CpGs in the Illumina Infinium HumanMethylation450 microarray. *Epigenetics* 8(2), 203–209 (2013).
- 8 Teschendorff AE, Marabita F, Lechner M *et al.* A beta-mixture quantile normalization method for correcting probe design bias in Illumina Infinium 450k DNA methylation data. *Bioinformatics* 29(2), 189–196 (2013).
- 9 Touleimat N, Tost J. Complete pipeline for Infinium(®) Human Methylation 450K BeadChip data processing using subset quantile normalization for accurate DNA methylation estimation. *Epigenomics* 4(3), 325–341 (2012).
- 10 Pidsley R, Wong CC, Volta M, Lunnon K, Mill J, Schalkwyk LC. A data-driven approach to preprocessing Illumina 450K methylation array data. *BMC Genomics* 14, 293 (2013).
- 11 Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics* 8(1), 118–127 (2007).
- 12 Leek JT, Johnson WE, Parker HS, Jaffe AE, Storey JD. sva: Surrogate Variable Analysis. R package version 3.10.0.
- 13 Jaffe AE, Irizarry RA. Accounting for cellular heterogeneity is critical in epigenome-wide association studies. *Genome Biol.* 15(2), R31 (2014).
- 14 Ligthart S, Marzi C, Aslibekyan S *et al.* DNA methylation signatures of chronic low-grade inflammation are associated with complex diseases. *Genome Biol.* 17(1), 255 (2016).
- 15 Horvath S. DNA methylation age of human tissues and cell types. *Genome Biol.* 14(10), R115 (2013).
- 16 Hannum G, Guinney J, Zhao L *et al.* Genome-wide methylation profiles reveal quantitative views of human aging rates. *Mol. Cell* 49(2), 359–367 (2013).
- 17 Kanyongo GY, Brook GP, Kyei-Blankson L, Gocmen G. Reliability and statistical power: how measurement fallibility affects power and required sample sizes for several parametric and nonparametric statistics. *J. Modern Appl. Stat. Methods* 6(1), 81–90 (2007).
- 18 Cleary TA, Linn RL, Walster GW. Effect of reliability and validity on power of statistical tests. *Sociol. Methodol.* 2, 130–138 (1970).
- 19 Meng H, Joyce AR, Adkins DE *et al.* A statistical method for excluding non-variable CpG sites in high-throughput DNA methylation profiling. *BMC Bioinformatics* 11, 227 (2010).
- **Advocating for the filtering out of low variability sites in BeadChip-based epigenetic studies.**
- 20 Bibikova M, Lin Z, Zhou L *et al.* High-throughput DNA methylation profiling using universal bead arrays. *Genome Res.* 16(3), 383–393 (2006).
- 21 Bjornsson HT, Brown LJ, Fallin MD *et al.* Epigenetic specificity of loss of imprinting of the IGF2 gene in Wilms tumors. *J. Natl Cancer Inst.* 99(16), 1270–1273 (2007).
- 22 Espinal AC, Wang D, Yan L *et al.* A methodological study of genome-wide DNA methylation analyses using matched archival formalin-fixed paraffin embedded and fresh frozen breast tumors. *Oncotarget* 8(9), 14821–14829 (2017).
- 23 Wahl S, Drong A, Lehne B *et al.* Epigenome-wide association study of body mass index, and the adverse outcomes of adiposity. *Nature* 541(7635), 81–86 (2017).
- 24 Lee MK, Hong Y, Kim SY, Kim WJ, London SJ. Epigenome-wide association study of chronic obstructive pulmonary disease and lung function in Koreans. *Epigenomics* 9(7), 971–984 (2017).
- 25 Rask-Andersen M, Martinsson D, Ahsan M *et al.* Epigenome-wide association study reveals differential DNA methylation in individuals with a history of myocardial infarction. *Hum. Mol. Genet.* 25(21), 4739–4748 (2016).
- 26 Montano C, Taub MA, Jaffe A *et al.* Association of DNA methylation differences with schizophrenia in an epigenome-wide association study. *JAMA Psychiatry* 73(5), 506–514 (2016).
- 27 Shimada-Sugimoto M, Otowa T, Miyagawa T *et al.* Epigenome-wide association study of DNA methylation in panic disorder. *Clin. Epigenetics* 9, 6 (2017).
- 28 Chen J, Just AC, Schwartz J *et al.* CpGFilter: model-based CpG probe filtering with replicates for epigenome-wide association studies. *Bioinformatics* 32(3), 469–471 (2016).
- **Argues for the filtering of sites based on ICC values computed from replicates of the data, which would be ideal in the presence of a suitable number of replicates.**

