



Published in final edited form as:

*Patterns (N Y)*. 2020 May 08; 1(2): .

## Patterns of Reliability: Assessing the Reproducibility and Integrity of DNA Methylation Measurement

Karen Sugden<sup>1,2,9,\*</sup>, Ellis J. Hannon<sup>3</sup>, Louise Arseneault<sup>4</sup>, Daniel W. Belsky<sup>5</sup>, David L. Corcoran<sup>2</sup>, Helen L. Fisher<sup>4</sup>, Renate M. Houts<sup>1</sup>, Radhika Kandaswamy<sup>4</sup>, Terrie E. Moffitt<sup>1,2,4,6</sup>, Richie Poulton<sup>7</sup>, Joseph A. Prinz<sup>2</sup>, Line J.H. Rasmussen<sup>1,8</sup>, Benjamin S. Williams<sup>1,2</sup>, Chloe C.Y. Wong<sup>4</sup>, Jonathan Mill<sup>3</sup>, Avshalom Caspi<sup>1,2,4,6</sup>

<sup>1</sup>Department of Psychology and Neuroscience, Duke University, Grey Building, 2020 West Main Street, Suite 201, Durham, NC 27705, USA <sup>2</sup>Center for Genomic and Computational Biology, Duke University, Durham, NC, USA <sup>3</sup>Complex Disease Epigenetics Group, University of Exeter Medical School, Exeter, UK <sup>4</sup>King's College London, Social, Genetic, and Developmental Psychiatry Research Centre, Institute of Psychiatry, Psychology, and Neuroscience, London, UK <sup>5</sup>Department of Epidemiology & Butler Aging Center, Columbia University Mailman School of Public Health, New York, NY, USA <sup>6</sup>Department of Psychiatry and Behavioral Sciences, Duke University School of Medicine, Durham, NC, USA <sup>7</sup>Dunedin Multidisciplinary Health and Development Research Unit, University of Otago, Dunedin, New Zealand <sup>8</sup>Clinical Research Centre, Copenhagen University Hospital Amager and Hvidovre, Hvidovre, Denmark <sup>9</sup>Lead Contact

### SUMMARY

DNA methylation plays an important role in both normal human development and risk of disease. The most utilized method of assessing DNA methylation uses BeadChips, generating an epigenome-wide “snapshot” of >450,000 observations (probe measurements) per assay. However, the reliability of each of these measurements is not equal, and little consideration is paid to consequences for research. We correlated repeat measurements of the same DNA samples using the Illumina HumanMethylation450K and the Infinium MethylationEPIC BeadChips in 350 blood DNA samples. Probes that were reliably measured were more heritable and showed consistent associations with environmental exposures, gene expression, and greater cross-tissue concordance. Unreliable probes were less replicable and generated an unknown volume of false negatives. This serves as a lesson for working with DNA methylation data, but the lessons are equally applicable

This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

\*Correspondence: karen.sugden@duke.edu.

#### AUTHOR CONTRIBUTIONS

K.S., D.L.C., and A.C. devised the project and the main conceptual ideas. K.S., E.J.H., and D.L.C. processed and generated datasets. K.S. and D.L.C. performed the analyses and designed the figures. L.A., H.L.F., T.E.M., R.P., C.C.Y.W., J.M. and A.C. acquired funding and supervised project activities. K.S. and A.C. drafted the manuscript. All authors provided critical feedback and helped shape the research, analysis, and manuscript.

#### Declaration of Interests

The authors declare no competing interests.

#### Supplemental Information

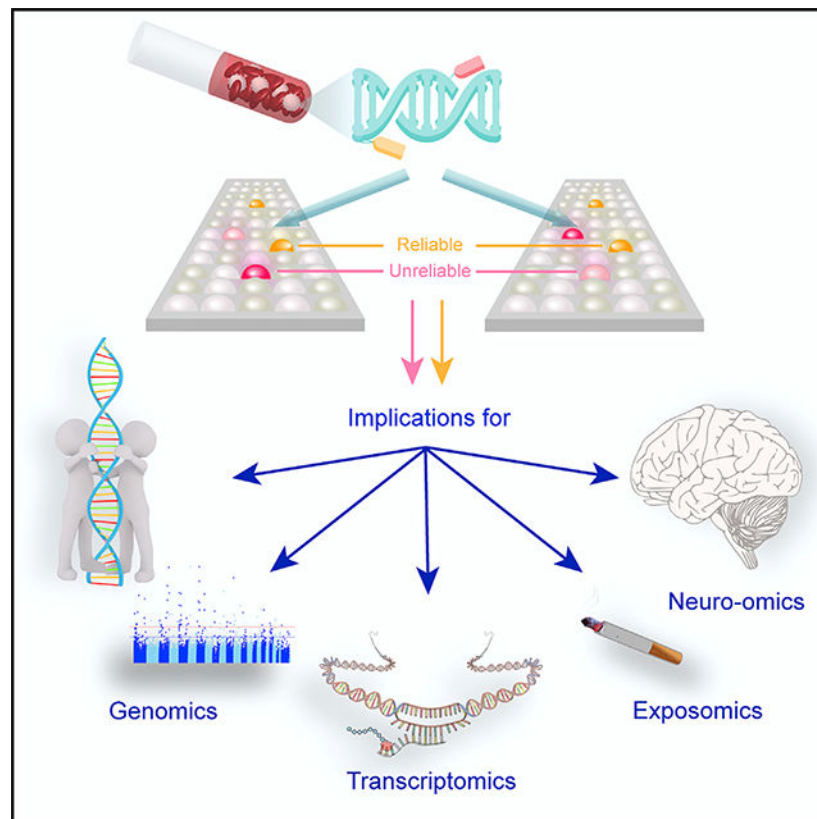
Supplemental Information can be found online at <https://doi.org/10.1016/j.patter.2020.100014>.

to working with other data: as we advance toward generating increasingly greater volumes of data, failure to document reliability risks harming reproducibility.

## In Brief

DNA methylation is an important mechanism of gene regulation. The most popular method to measure methylation is to use BeadChips that contain probes to index hundreds of thousands of methylation sites at once. However, these probes are not equally reliable. In blood DNA, unreliable probes were less heritable and less likely to index gene expression, and associations were less replicable. This has serious downstream consequences for reproducible science and should serve as a caution for all data scientists regardless of discipline.

## Graphical Abstract



## INTRODUCTION

DNA methylation is an epigenetic mechanism that occurs by the addition of a methyl ( $\text{CH}_3$ ) group to DNA, resulting in modification of genetic function without changes to DNA sequence. This mechanism plays an important role in human development and disease, primarily by regulating gene expression.<sup>1</sup> Because of the modifiable nature of epigenetic influence, research into DNA methylation has heralded a new era in the elusive search for the route by which the external world might “get under the skin.”<sup>2</sup> By its very nature, this question spans multiple disciplines; geneticists,<sup>3</sup> biologists,<sup>4</sup> computational scientists,<sup>5</sup>

neuroscientists,<sup>6</sup> social scientists,<sup>7</sup> and philosophers<sup>8</sup> have been drawn to massive new data about the epigenome with an eye toward how it might explain health, disease, and our very nature. The promise of the epigenetics revolution has been sweeping.

In humans, DNA methylation occurs at specific sites across the genome (almost exclusively CpG sites, where a cytosine nucleotide is located next to a guanidine nucleotide), and there exist hundreds of thousands of such sites. Advances in technologies for quantifying site-specific DNA methylation have aided an explosion of research aimed at identifying associations between numerous environmental exposures, disease processes, and methylomic variation.<sup>9–12</sup> One such measurement technology, the Infinium BeadChip produced commercially by Illumina, has fueled much of the research in epigenetic epidemiology. This platform was developed to simultaneously assay thousands of DNA methylation targets in the genome. The relative ease of use, low cost, and modest sample requirements of this technology have enabled a new generation of researchers to add DNA methylation to their research programs, which only a few years ago would have posed an insurmountable challenge. We are among this new generation. This article reports our experience, excitement, and frustration, as a team of multidisciplinary scientists, trying to understand and use these data.

When we began to produce DNA methylation data, we reviewed the literature for best-practice information and guidelines to ensure the highest validity and downstream reproducibility. It was at this point we realized there was no consensus. We had generated data using the Infinium Methylation450 (450K) BeadChip, the gold standard for epigenome-wide DNA methylation data. This provides ~450,000 measurements per individual subject. However, we learned that a significant proportion of the thousands of data points do not yield the equivalent value when quantified twice from the same DNA sample.<sup>13,14</sup> This situation is compounded by the nature of our work, which involves repeated measurement of individuals studied longitudinally. This in itself raises an additional complication: measurement methods become obsolete and are superseded by new, improved products. In this case, the 450K BeadChip was recently replaced by the Infinium MethylationEPIC (EPIC) BeadChip, which contains most of the content (approximately 93%) of the 450K BeadChip augmented with probes covering an additional ~400,000 CpG sites. Published research has suggested that at the array level, DNA methylation values generated using both iterations of Illumina DNA methylation BeadChips are highly correlated, yielding correlations  $>0.9$ ;<sup>15–18</sup> however, the reliability of individual-level probe measurements between the two arrays varies substantially. Using DNA derived from blood collected from 145 adults, one study<sup>17</sup> observed that reliability correlations between probes on the 450K and EPIC BeadChip ranged from  $-0.34$  to  $0.95$  with a median value of  $0.15$ , and only 2.6% of the ~420,000 probes assayed had reliability correlations above 0.8. Using DNA derived from blood collected from 109 newborns and 86 adolescents, a second study<sup>18</sup> observed similarly low correlations (median  $r = 0.23$ , only ~10% of probes with correlations  $>0.8$ ).

These aforementioned reports documented patterns of uneven reliability in the repeated measurement of DNA methylation.<sup>13,14,17,18</sup> However, we were not prepared for the scarcity of information documenting the consequences of these patterns; consequences that, if shown to affect inferences made from DNA methylation data, would have widespread implications

for reproducibility. Most research studies treat the ~450,000 observations as “equals,” each as likely as the next to report true biological differences from a statistical point of view. However, to uncover consistent, replicable signals of DNA methylation dynamics, be it over time, between populations, or between exposures, measurement reliability is crucial. Analysis of probes that cannot be repeatedly measured with precision has the potential to yield irreproducible findings borne from spurious associations, and, just as importantly, may miss discoveries.

Here we share how we went about learning of the cross-disciplinary data challenges of high-throughput DNA methylation data and discuss the implications of these challenges for data processing, analysis, algorithm generation, and interpretation. Our goal is to promote communication about careful practices for working with the new data being generated in this important field.

We first performed test-retest measurement assessments to quantify the reliability of DNA methylation data. We assessed probe reliability between the two types of BeadChips using data on 350 DNA samples measured twice; once using the 450K BeadChip and again using the EPIC BeadChip. The individuals are participants in the E-Risk Study, a birth cohort of 2,232 twins born in 1994–1995 in the United Kingdom. DNA methylation was measured at age 18 years, when participants contributed whole blood for DNA analysis. Probe reliability was defined as the intraclass correlation (ICC) between repeat measures of individual probe  $\beta$  values measured on the two BeadChips. We then assessed the impact of differential reliability on numerous lines of enquiry of interest to many researchers, ourselves included. First, we tested how reliability influenced the ability to detect genetic and environmental effects on the epigenome through (1) analysis of heritability in the E-Risk twin sample and (2) analysis of methylation quantitative trait loci (mQTLs) identified in genome-wide association studies (GWAS) of DNA methylation. Second, we tested the implications of differential reliability for association testing by analyzing results of epigenome-wide association studies of tobacco smoking, one of the most harmful health risks in the modern world.<sup>19</sup> Third, we tested the implications of differential reliability for epigenetic biomarker development by analyzing multi-probe-algorithm-based measurements that are intended to capture information about aging (i.e., “DNA methylation clocks”). Finally, we tested the implications of differential reliability in ascribing biological function to DNA methylation by assessing the impact of reliability on (1) correlations between DNA methylation and gene expression and on (2) correlations between levels of DNA methylation measured in blood tissue and brain tissue.

## RESULTS

### Reliability of CpG Probes Is Low and Highly Variable

We use “reliability” to refer to the reproducibility of methylation probes’ values. We measured probe values twice from the same DNA source (DNA was sourced from a single blood draw via a single extraction). One set of measures was made using the 450K BeadChip, the other set using the EPIC BeadChip. Our analysis was restricted to probes found on both platforms (438,593 probes).

Probe reliabilities were computed using ICCs calculated for each of the 438,593 autosomal probes present on both the EPIC and 450K BeadChips that passed quality control. ICCs are an oft-used metric to assess reliability in test-retest situations,<sup>20</sup> and many different models exist depending on the way in which the test-retest data are generated. Here, we calculated ICCs based on a mean-rating ( $k = 2$ ), absolute-agreement, 2-way random-effects model. We chose this model using the guidelines outlined by Koo and Li,<sup>20</sup> where mean-rating ( $k = 2$ ) relates to the number of repeated measures (i.e., BeadChips per sample); absolute agreement requires that not only do the values across BeadChips correlate but that values are in agreement; and 2-way random effects relates to the generalizability of the ICCs to any subsequent similarly characterized rater (where rater = BeadChip probe).

ICCs between probes ranged from  $-0.28$  to  $1.00$  (Supplemental Information, Section 1.1; Figure S1; Data S1). Probe reliabilities were skewed toward zero, with a mean of  $0.21$  (median =  $0.09$ ). This is low reliability considering that, in the context of establishing reliable measurement, ICCs below  $0.4$  are considered “poor,” those between  $0.4$  and  $0.6$  are considered “fair,” between  $0.6$  and  $0.75$  “good,” and above  $0.75$  “excellent.”<sup>21</sup>

The reliabilities that we observed in our data were highly correlated with the reliabilities observed by Logue et al.,<sup>17</sup> who also compared probes across 450K and EPIC BeadChips ( $r = 0.86$ ,  $p < 0.01$ , Supplemental Information, Section 1.1; Figure S2). This suggests that the low reliabilities that we observed across the arrays are reproducible in other datasets. Importantly, the low reliabilities that we observed were unlikely to be solely due to differences between 450K and EPIC BeadChip probes. First, previous studies have documented similar low reliabilities in 450K-450K probe comparisons<sup>13,14</sup> and EPIC-EPIC probe comparisons.<sup>17</sup> Second, we conducted EPIC-EPIC array comparisons for a subset of Dunedin Study samples ( $n = 28$ ) (for comparison purposes, we restricted analysis to the  $\sim 440,000$  probes overlapping with the 450K array as described throughout this paper). Several noteworthy details emerged. (1) The median reliability in our EPIC-EPIC comparison was  $0.26$ . This is higher than the median reliability ( $0.09$ ) observed in our 450K-EPIC comparisons, but still falls squarely in what is considered to be “poor” reliability.<sup>21</sup> (2) It is not clear what accounts for the higher EPIC-EPIC reliability; it could be due to consistency of the platform or it could be due the fact that, unlike probes for the 450KEPIC comparisons, probes for the EPIC-EPIC were assayed at the same time, using the same reagents, equipment, and so forth. (3) The correlations between the EPIC-EPIC reliabilities estimated by us in the Dunedin Study with the 450K-EPIC reliabilities (estimated by us in the E-Risk Study) was  $0.77$  (Figure S3). (4) When performing the analyses set forth in this manuscript using EPIC-EPIC ICCs rather than 450K-EPIC ICCs, we arrive at the same conclusions: we found that, like between-array reliability, within-array reliability is low, skewed toward zero, and has detrimental effects on research findings, and that differences in 450K and EPIC BeadChip probes are unlikely to be the sole cause of between-array unreliability (Supplemental Information, Section 1.1).

As a sanity check, we also sought to replicate previously observed associations between reliability and (1) the mean and standard deviation (SD) of methylation levels ( $\beta$  values)<sup>13,14,17</sup> (Supplemental Information, Section 1.2) and (2) the genomic annotation (location) of probes<sup>13,18</sup> (Supplemental Information, Section 1.3). We observed the same

associations as previously reported. Taken together, this suggests BeadChip-wide differential reliabilities are reproducible and systematic in pattern.

Previous methodological studies have drawn attention to three factors that might compromise the quality of methylation BeadChip data: probe invariance,<sup>22–24</sup> potential probe hybridization problems,<sup>25</sup> and skewness.<sup>26</sup> We tested whether these features are sufficient to capture unreliability. They are not. Probe unreliability exists in probes that are variable or do not have potential probe hybridization problems, and probe reliabilities calculated on  $\beta$  values resemble the reliabilities of M values, a method for transforming skewed probe distributions<sup>26</sup> (Supplemental Information, Section 1.4).

In summary, we replicated previous reports of low reliability across probes common to the 450K and EPIC BeadChips, demonstrating that, paradoxically, poor reliability is reproducible. Moreover, factors commonly thought to account for unreliability (such as invariance) do not provide a satisfactory account of its ubiquity.

### Evaluating the Consequences of Unreliable Probe Measurements

Our data suggest that the majority of probes we tested have low test-retest reliability. We now examine the practical implications of this observation for epigenetic research by applying our 450K-EPIC reliabilities to the results of previously published epigenetic studies. In all cases, these previously published studies were based on data derived using 450K BeadChips because (1) the EPIC BeadChip is relatively new, and most published research is based on the 450K BeadChip, (2) the probes common to the EPIC and 450K BeadChips reflect almost all ( $\sim 93\%$ <sup>16</sup>) of the probes unchanged from the 450K BeadChip, and (3) earlier 450K-450K comparisons showed patterns of reliabilities similar to those of the 450K-EPIC comparison.<sup>13,14</sup>

### Estimates of Genetic and Environmental Effects on DNA Methylation Are Affected by Unreliable Measurement

Genetic and environmental effects on a phenotype can be estimated by comparing the relative phenotypic differences between monozygotic (MZ) and dizygotic (DZ) twins. The assumptions behind this model are that additive genetic factors are perfectly correlated between MZ twins (i.e., genetic correlation = 1) but are only 50% correlated between DZ twins (i.e., genetic correlation = 0.5) and that shared non-heritable influences are equally similar between MZ and DZ twin pairs. We previously reported the probe-specific genetic and environmental architecture of DNA methylation.<sup>24</sup> Using our twin design, we decomposed variation in each probe into three variance components: additive genetic effects (labeled “A”), shared environmental effects (“C”; environmental effects that each twin in a twin pair share, making twins more similar to each other), and non-shared (or unique) environmental effects (“E”; environmental effects that are specific to each twin within a pair, making twins less similar to each other). Figure 1 shows the association between probe reliability and estimates of A (Figure 1A), C (Figure 1B), and E (Figure 1C). Reliability was significantly correlated with higher heritability ( $r = 0.70$ ,  $p < 0.01$ , Figure 1A). In contrast, low-reliability probes tended to be suffused with more non-shared environmental variance ( $r = -0.58$ ,  $p = 1.00$ , Figure 1C). Given that the non-shared environmental variance component



in biometric models also includes measurement error, these probes are possibly less likely to reflect true environmental effects than they are to reflect unreliable measurement. (The correlation between reliability and estimates of shared environmental variance [C] was low,  $r = -0.07$ , possibly reflecting the fact that the classical twin design has limited power to identify precise estimates of shared environmental influence.<sup>27</sup>)

We further examined how unreliability affects discovery research about the genetic etiology of DNA methylation. A recent GWAS of DNA methylation identified ~55,000 methylation mQTLs, DNA sequence variants that are associated with differential DNA methylation.<sup>28</sup> Figure 2 shows that the reliability of probes indexed by mQTLs in our data ( $N = 50,900$ ) is higher than the reliability of probes that are not ( $N = 387,693$ ).

In summary, given that a significant proportion of probes are suffused with unreliability (as indicated by poor test-retest reliability and as further indexed by high E-components in biometric models), the ability to detect associations between DNA methylation levels and genetic influences will be compromised.

### Probe Reliability Affects Association Testing

We hypothesized that reliability is related to the likelihood that associations between environmental exposures and specific probes would replicate across independent studies. To test this, we focused on one of the most robust findings in epigenetic epidemiology: the effect of tobacco smoking on DNA methylation. We identified 22 studies that reported an epigenomewide analysis of current versus never smoking using the 450K BeadChip platform<sup>12,29–34</sup> (Table S4). For each study, we obtained lists of probe IDs and direction of effect for probes that were significantly associated with current smoking (as determined by the study authors; total number of probes = 3,724; number of probes per study = 84–2,441). We then determined the extent to which individual probes replicated across the 22 studies by summing the number of times each probe was listed with consistent direction of effect (i.e., consistent cross-study increases or decreases in methylation in response to smoking). The number of individual replications across studies was associated with reliability ( $r = 0.52$ ,  $p < 0.001$ , Figure 3). The mean number of replications for low-reliability probes (here defined as reliability  $< 0.4$ ) was 6.84 (median = 1, SD = 6.78,  $n = 1,630$  probes), whereas the mean for high-reliability probes (reliability  $> 0.75$ ) was 13.1 (median = 15, SD = 5.11,  $n = 391$  probes).

In summary, the likelihood of replicating associations between exposures and DNA methylation probes is significantly greater when studying reliable probes. Unreliable probes are likely to generate false positives and to mask true associations and are less likely to be reproducible.

### Publicly Available DNA Methylation Aging Algorithms Contain Unreliable Measurements

There is enormous interest in developing and applying algorithms that use DNA methylation to index biological aging.<sup>35</sup> A critical component of the success of these “DNA methylation clocks” is that probes comprising the algorithms are reliably measured so that they might be applied to any external dataset. We tested the hypothesis that these algorithms are more likely to capture reliable probes than unreliable probes. Figure 4 shows the distribution of

probe reliabilities for three established DNA methylation aging-associated clocks: (1) the “Hannum clock,”<sup>36</sup> (number of probes = 63), (2) the “Horvath clock”<sup>37</sup> (number of probes = 334), and (3) the “Levine clock”<sup>38</sup> (number of probes = 512; the number of probes reflects those available in our data). Each aging algorithm had median probe reliabilities higher than that of the background distribution. However, the distribution for all three algorithms was not solely composed of reliable probes; each algorithm contained many probes whose  $\beta$  values were unreliable.

In summary, externally validated DNA methylation algorithms are generally composed of reliable probes. However, their performance could be improved by utilizing more reliable DNA methylation measurements. This perhaps emphasizes the point that algorithms of this type necessitate careful, extensive external validation; we hypothesize that algorithms over-represented by unreliable probes will, by their very nature, fail to perform well under varied testing situations.

### Reliability Influences the Association between DNA Methylation and Gene Expression

A goal of epigenetic discovery is to assign biological meaning to the observed patterns of DNA methylation (e.g., Schubeler<sup>2</sup> and Teschendorff and Relton<sup>39</sup>). To this end, we tested the hypothesis that DNA methylation probes with higher reliability were more likely to index variation in gene expression, the process by which the information encoded in a gene is used to direct the assembly of a protein molecule. We used two approaches.

First, we used the results of global DNA methylation-gene expression correlation patterns described by Kennedy et al.,<sup>40</sup> wherein 36,485 and 114,536 unique DNA methylation probes were associated with gene expression across two cohorts (GTP and MESA, respectively;  $p < 1 \times 10^{-5}$ ). Figure 5A shows that these significantly correlated methylation probes were more likely to be reliable (median reliability in GTP = 0.21, proportion of these probes with reliability  $>0.75 = 11.2\%$ ; median reliability in MESA = 0.20, proportion of these probes with reliability  $>0.75 = 10.1\%$ ; gene set enrichment analysis [GSEA] enrichment  $p < 1 \times 10^{-4}$  in each) than methylation probes that were not discovered to be related to gene expression. Furthermore, probes that were significantly correlated with gene expression in both datasets had higher reliabilities than those identified in only one dataset (median reliability = 0.36 versus 0.17, proportion of probes with reliability  $>0.75 = 14.7\%$  versus 9.4% for both datasets versus one dataset, respectively; GSEA enrichment  $p < 1 \times 10^{-4}$ ). This suggests that reliability of DNA methylation probes influences the ability to detect correlates of biological function in a reproducible manner.

Second, using gene expression data available in the Dunedin Study, we calculated the correlation between gene expression probeset values with DNA methylation  $\beta$  values for every CpG probe localized to the transcription start site (TSS) of that gene. We restricted analysis to probes within the TSS, as these are hypothesized to have direct effects upon expression of the localized gene. As shown in Figure 5B, DNA methylation probes that significantly correlated with expression probesets ( $\alpha = 1 \times 10^{-7}$ ,  $n = 278$ ) had significantly higher reliabilities than DNA methylation probes that did not ( $n = 23,261$ ; median reliability of correlated probes = 0.64, proportion of these probes with reliability  $>0.75 = 36.0\%$ ;



median reliability of non-correlated probes = 0.04, proportion of these probes with reliability  $>0.75 = 3.4\%$ ; Figure 5C).

In summary, DNA methylation probes were more likely to correlate with transcriptional variation if they were reliably measured. Reliable probes are more likely to index reproducible biological correlates, whereas unreliable probes may mislead about biological function.

### Reliability Influences the Concordance of Blood and Brain Methylation Levels

Most epidemiological investigations into exposure-related differential DNA methylation are undertaken using DNA derived from whole blood. This is an expedient choice due to the relative ease of collecting blood in population-based studies. However, many exposures in which epidemiologists are interested are hypothesized to have their effects (or consequences) in other tissues, such as the brain, raising the question of whether peripheral blood is a problematic surrogate tissue. Previously, we evaluated the similarity of methylation levels between blood DNA and DNA from four brain regions (prefrontal cortex, entorhinal cortex, superior temporal gyrus, and cerebellum) using the 450K BeadChip, and showed that only a small proportion of probes measured in blood correlate with methylation levels in the brain.<sup>41</sup>

We hypothesized that these small numbers of probes that register similar levels of DNA methylation in blood and brain tissue would be over-represented by high-reliability probes. To test this, we cross-referenced the correlations between DNA methylation levels in blood and each of four brain regions (“blood-brain concordance”) with our 450K-EPIC probe reliabilities. Blood-brain concordance was related to reliability ( $\rho = 0.22-0.38$ ,  $p < 0.01$  across the four brain regions). Figure 6 shows the distribution of reliability across low- ( $<0.4$ ), mid- ( $0.4-0.75$ ), and high-concordance ( $>0.75$ ) probes in four brain regions. Median reliabilities for probes with low blood-brain concordance were 0.08 regardless of brain region, while median reliabilities for probes with high blood-brain concordance were 0.90 across the four brain regions. Moreover, probes that showed high blood-brain concordance in all four brain tissues were the most reliable (median reliability = 0.92, number of probes = 6,774, proportion of these probes with reliability  $>0.75 = 78.7\%$ ) while probes that had low blood-brain concordance in each of the four brain tissues were the least reliable (median reliability = 0.08, number of probes = 397,091, proportion of these probes with reliability  $>0.75 = 3.1\%$ ).

In summary, reliable probes are more likely to exhibit cross-tissue concordance in DNA methylation. Unreliable probes may be less likely to prove useful in developing blood-based biomarkers of brain dysregulation.

## DISCUSSION

The reliability of probe-level DNA methylation measurement is highly variable across the ~440,000 sites indexed on the 450K and EPIC BeadChips. This differential reliability has detrimental downstream implications: it undermines published research and masks potential new discoveries.

First, we demonstrated that detection of both environmental and genetic effects on DNA methylation is related to differential probe reliability. The extent to which DNA methylation responds to environmental influences is under intense investigation and is thought to be one route via which environmental exposures “get under the skin.”<sup>2</sup> There is also much interest in the relationship between DNA sequence variation and DNA methylation.<sup>23,24</sup> Here, we showed that the most reliable probes tend to be under significant genetic influence, whereas the least reliable probes are suffused with non-shared environmental variation (which also includes variation arising due to measurement error). These findings suggest that for a proportion of sites that indicate high sensitivity to environmental input, identification of true signal might be hindered by the relatively higher probability of imprecise measurement and that insights into the genetic basis of methylation may be missed due to the poor reliability of DNA methylation.

Second, we demonstrated the implications of differential reliability for epigenome-wide association testing. To achieve this we focused on tobacco smoking, one of the most replicable findings in epigenetic epidemiology. Here we showed that the likelihood of replication across studies increases with probe reliability. We also showed how unreliable probes may slow biomarker discovery. Arguably, “DNA methylation clocks” have been one of the major success stories of epigenetic epidemiology.<sup>36–38</sup> We found that these clocks are enriched for reliable probes but that the algorithms also contain noisy measurements, and it is possible that applying machine learning to uniformly reliable data will improve precision in this and other areas.

Third, we demonstrated the implications of differential reliability for integrating DNA methylation data with sequence and transcriptomic data. Here we showed that probe reliability is necessary to accurately estimate genetic contributions to DNA methylation, to identify gene expression correlates, and to detect correlated DNA methylation signatures across tissues. If the goal is robust and replicable biological inference from site-specific DNA methylation, it is necessary to restrict analysis to those probes that can be reliably measured.

There are some caveats to this study. First, these findings are restricted to DNA derived from blood. However, findings described here will be of use to the majority of researchers in epigenetic epidemiology and to researchers looking for clinical application of epigenetic findings, since blood is the most common substrate from which DNA is derived and biomarkers are developed. Second, our study comprises young adults; it is possible that age-related change in DNA methylation at certain sites in the genome influences the pattern of reliability. That said, the pattern of reliability coefficients observed in our study is consistent with that seen in newborns,<sup>18</sup> 14-year-olds,<sup>18</sup> and ~30-year-olds.<sup>17</sup> Third, findings are restricted to the ~440,000 probes common to both the 450K and EPIC BeadChips. However, Logue et al.<sup>17</sup> reported similar reliability distributions for EPIC-EPIC comparisons in 11 individuals and we found better, but overall poor reliability for EPIC-EPIC comparisons in our data as well. Moreover, for the probes overlapping the two arrays, the EPIC-EPIC reliabilities were highly correlated with the 450K-EPIC reliabilities. The reason we emphasize between-array probe comparisons is that the goal of many researchers’ work is to both make discoveries and replicate discoveries made by others. Given rapid advances in

technologies and the proliferation of available data, it is increasingly the case that researchers need to integrate data that have been created using different arrays; indeed, although the 450K chip is no longer available, the vast bulk of DNA methylation research to date has used this array. As such, an important challenge for data scientists is how to integrate data from different arrays, whether this is in the service of evaluating targets for further scientific interrogation or in meta-analysis (e.g., one needs to know whether published results generated by 450K data are generalizable to new EPIC data and, ultimately, whether EPIC data will be generalizable to new technologies in future). In this case, between-array reliability is the relevant metric.

Taken together, at the very least unreliable probes are uninformative. At worst, they hinder scientific progress. In the GWAS world, much has been done to improve replicability of research, from increasing sample sizes to standardizing data pipelines (e.g., Visscher et al.<sup>42</sup>). In the epigenetic world, researchers have adopted many similar considerations (e.g., Lehne et al.<sup>43</sup> and Yan et al.<sup>44</sup>), but unreliability in the quantitative measurement of DNA methylation is a unique challenge. We list possible responses below.

First, to approximate reliable measurements, it is possible to filter data based on intrinsic properties of probes, such as invariance or hybridization properties. However, restricting analysis to variant probes or to probes without sequence-related performance issues is not sufficient to guarantee reliability; we found that these probes were not uniformly reliable (Supplemental Information, Section S1.4). Furthermore, restricting analysis to only variant probes conveys no enhancement of power to detect associations between reliability and (1) estimates of genetic and environmental influences on DNA methylation, (2) mQTL probes, and (3) concordance in DNA methylation levels between blood and brain tissue (Supplemental Information, Section S2). Second, it is possible to return to the practice, once routine, of using alternative technologies such as pyrosequencing to perform post-analysis validation of positive DNA methylation findings. This approach comes with two caveats, one of which is that it can only detect false positives; false negatives would remain undetected. A second caveat is that as science is shifting toward a culture of open-access and publicly available data, more and more researchers are becoming endpoint data users and as such are not involved in experimental data production. In this scenario, the task of experimental validation of individual findings, potentially in the thousands, is resource heavy, logistically impractical, and financially prohibitive. A third response is to generate pre-analysis reliability metrics, as we did in this report. Indeed, for publicly available data, this is currently the only feasible method of providing individual probe reliability metrics to end-users. To aid standardization, we have made available our reliability metrics for all measured probes (Data S1). Going forward, we suggest that researchers make the assessment of reliability standard practice when designing and measuring DNA methylation. This is because, despite evidence that our individual probe reliabilities correlate highly with those reliabilities reported previously,<sup>17</sup> we do not yet know the full extent to which demographic (e.g., age), measurement (e.g., batch), and source (e.g., tissue) factors may influence measurement reliability. Additionally, specific experimental designs (e.g., longitudinal designs and meta-analyses requiring incorporation of data from different sources, array types and batches, or cross-sectional single time-point designs) would determine which type of reliability metric to employ (e.g., within-array versus between-

array); the reliability metrics reported here might not be the most suitable. By subsetting our repeated samples and calculating reliability, we determined that running just 25 replicates will identify 80% of the reliable probes (reliability >0.75) identified when using 350 replicates (Supplemental Information, Section S3.2). Fulfilling this recommendation would require additional investments during project planning along with commitment of support from funding agencies. The effort associated with incorporating reliability assessment into routine quality control, as we propose, is far outweighed by the benefits to science from improved replicability. The goal would be, at the very least, to report the reliability associated with any probe for which conclusions are drawn; this will allow readers to make independent assessments of the confidence in the probe measurements. Even better would be to filter data, before analysis, on the basis of reliability metrics. Subsetting data in this way should help reduce false positives (by reducing the probability of spurious associations) and possibly false negatives. Although familywise error-rate corrections would not be greatly affected (e.g., within the data we present here, Bonferroni correction would reduce the testing threshold from  $\alpha = 1.14 \times 10^{-7}$  for ~440,000 probes to  $\alpha = 1.77 \times 10^{-6}$  for ~28,000 probes with reliability >0.75), false-discovery-rate corrections may be affected.

Researchers from diverse disciplines have been drawn by the promise of DNA methylation as a convenient vector by which the social environment might exert its effects on an organism's biology. They are also drawn to the relative simplicity of Illumina BeadChip data in both content and comprehensiveness. Anecdotally, we have encountered two reactions to the phenomenon of differential reliability. First, some researchers have expressed little surprise at its existence, coupled with a belief in the self-correcting power of the field to purge false negatives and positives over time. Our response to this is that better use of intellectual and financial resources might be made in analysis of data that are pre-validated, rather than cycling through replication attempts using unreliable measures that are bound to fail. Second, others have expressed shock and alarm that this phenomenon exists at all; these researchers are often new to the field and are not intimately familiar with the nuances of how data are produced or their biological meaning. Our response here is that DNA methylation data are not universally unusable—their suitability for analysis is contextual. Determination of reliability gives researchers confidence in the data they are using, be they new adopters, end-users, or seasoned experts.

Open-access availability of data is accelerating with encouragement from journal publishers and funding agencies. More and more researchers are using these big data; DNA methylation data are only one example of such. End-users rely on providers to verify the integrity of data, but just because data are massive in scale does not preclude the need for careful evaluation of their precision. The reproducibility crisis in science has drawn attention to two Rs: reproducibility (the extent to which consistent results are obtained when an experiment is repeated with the exact same inputs) and replicability (the extent to which consistent results are obtained when an experiment is repeated with the same design but with inputs from other sources). Here, we highlight a potential third “R,” reliability. Reliability is a fundamental aspect of replicability. If desired inputs do not yield the same value when the source differs, replication is impossible. In this sense, test-retest reliability is a tool that has widespread applicability to the entire data-science community, especially where big data are used. The National Academies of Sciences, Engineering, and Medicine recently published a

report<sup>45</sup> on the state of reproducibility and replicability in science, along with suggestions for improvement: "...[r]esearchers should, as applicable to the specific study, provide an accurate and appropriate characterization of relevant uncertainties when they report or publish their research ...". Unreliable probe measurement is one such uncertainty. We hope that our findings will improve the integrity of DNA methylation studies and serve as a cautionary reminder for those generating and implementing big data of any type.

## EXPERIMENTAL PROCEDURES

Full details are provided in Supplemental Information, Section S3.

### Samples

We report data from two samples. The Environmental Risk (E-Risk) Longitudinal Twin Study tracks the development of a 1994–1995 birth cohort of 2,232 British children followed to age 18 years.<sup>46</sup> The Dunedin Study tracks the development of a 1972–1973 birth cohort of 1,037 New Zealand children followed to age 45 years.<sup>47</sup>

### DNA Methylation

In E-Risk, DNA was derived from peripheral blood drawn at age 18 years. In Dunedin, peripheral whole blood was drawn at 38 and 45 years. In E-Risk, DNA from 350 study members was selected for analysis using both Infinium MethylationEPIC (EPIC; Illumina, CA, USA) and Illumina Infinium HumanMethylation450 BeadChip (450K BeadChip; Illumina). The remainder of the cohort ( $n = 1,308$ ) was assayed using the 450K BeadChip only, as previously described.<sup>48</sup> In Dunedin at age 38, DNA from 819 study members was assayed using the 450K BeadChip, as previously described.<sup>48</sup> In Dunedin at age 45, DNA from 28 study members was assayed twice using the EPIC BeadChip. E-Risk DNA methylation assays were run by the Complex Disease Epigenetics Group at the University of Exeter Medical School (UK) ([www.epigenomicslab.com](http://www.epigenomicslab.com)), and Dunedin assays were run by the Molecular Genomics Shared Resource at the Duke Molecular Physiology Institute, Duke University (USA).

### Gene Expression

RNA was derived from peripheral blood drawn into PAXGene RNA tubes at age 38 years in Dunedin. Expression data were generated from RNA using the Affymetrix PrimeView Human Gene Chip (Affymetrix, CA, USA) by the Duke University Microarray Core Facility. Data quality control and RMA (robust multichip average) normalization were carried out using the *affy* Bioconductor package in the R statistical programming environment. After quality control, expression data were available for 836 individuals.

### Probe Reliabilities

Probe reliabilities are computed using intraclass correlation (ICC) estimates, calculated for each autosomal probe present on both the EPIC and 450K BeadChip ( $N = 438,593$ ). ICCs are an oft-used metric to assess reliability in test-retest situations,<sup>20</sup> and many different models exist depending on the way in which the test-retest data are generated. Here, we calculated ICCs based on a mean-rating ( $k = 2$ ), absolute-agreement, 2-way random-effects

model. To compare whether test-retest model choice had an effect on reliability estimates, we also computed Pearson product-moment correlation coefficients. Pearson correlation coefficients and ICC estimates of reliability were highly similar ( $r = 1.00$ ,  $p < 1 \times 10^{-4}$ ; Figure S9).

### Statistical Analysis

All analyses were performed in the R statistical programming environment, often using Bioconductor packages. Unless otherwise noted, correlations are reported as Pearson correlation coefficients. Summary statistics, such as probe mean and SD, were based on the 350 samples processed on the 450K array. GSEA was performed using the *fgsea* Bioconductor package<sup>49</sup> with 10,000 permutations. The proportion of variance in DNA methylation explained by heritable (A), shared environmental (C), and unshared or unique environmental (E) factors was estimated using structural equation modeling implemented with functions from the *OpenMx* R package.<sup>50</sup>

### DATA AND CODE AVAILABILITY

E-Risk 450K DNA methylation data are accessible from the Gene Expression Omnibus (accession code GEO: GSE105018). Data from the Dunedin Study are not publicly available due to lack of informed consent and ethical approval for open access, but are available on request by qualified scientists. Requests require a concept paper describing the purpose of data access, ethical approval at the applicant's institution, and provision for secure data access. We offer secure access on the Duke University, Otago University, and King's College London campuses. All data on probe reliability and characteristics for the 450K-EPIC comparison (Data S1) are available at <https://osf.io/83ucs/>. The data underlying analysis of consequences of unreliability on heritability and blood-brain concordance are available from <https://www.epigenomicslab.com/online-data-resources/>. Code is available on request from the corresponding author.

### Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

### ACKNOWLEDGMENTS

We thank the E-Risk Study and Dunedin Study members, research staff, and Dunedin Study founder Phil Silva. The E-Risk Study is funded by the Medical Research Council (G1002190 to T.E.M., L.A., and A.C.), the National Institute of Child Health and Human Development (HD077482 to A.C.), a joint UK Economic and Social Research Council (ESRC) and Biotechnology and Biological Sciences Research Council (BBSRC) grant (ES/N000277/1 to C.C.Y.W.), and an MQ Fellows Award (MQ14F40 to H.L.F.). The Dunedin Longitudinal Study is funded by the New Zealand Health Research Council, the New Zealand Ministry of Business, Innovation, and Employment (to R.P.), the National Institute on Aging (AG032282 to T.E.M. and A.C. and AG049789 to T.E.M.), and the Medical Research Council (MR/P005918/1 to T.E.M.). Additional support was provided by the Jacobs Foundation (to T.E.M. and A.C.), Charles Lafitte Foundation Duke Faculty Seed Grant (to A.C.), and by a Distinguished Investigator Award from the American Asthma Foundation to J.M. L.A. is the Mental Health Leadership Fellow for the UK Economic and Social Research Council, D.W.B. is a Jacobs Foundation Fellow, and H.L.F. is supported by a British Academy Mid-Career Fellowship (MD\170005). This work used a high-performance computing facility partially supported by grant 2016-IDG-1013 ("HARDAC+: Reproducible HPC for Next-generation Genomics") from the North Carolina Biotechnology Center. We would like to acknowledge the assistance of the Duke Molecular Physiology Institute Molecular Genomics core for the generation of data for this paper.



## REFERENCES

1. Robertson KD (2005). DNA methylation and human disease. *Nat. Rev. Genet* 6, 597–610. [PubMed: 16136652]
2. Schubeler D (2015). Function and information content of DNA methylation. *Nature* 517, 321–326. [PubMed: 25592537]
3. Velasco G, and Francastel C (2019). Genetics meets DNA methylation in rare diseases. *Clin. Genet* 95, 210–220. [PubMed: 30456829]
4. Burggren WW, and Crews D (2014). Epigenetics in comparative biology: why we should pay attention. *Integr. Comp. Biol* 54, 7–20. [PubMed: 24722321]
5. Ruskin HJ, and Barat A (2018). Recent advances in computational epigenetics. *Adv. Genomics Genet.* 8, 12.
6. Tognini P, Napoli D, and Pizzorusso T (2015). Dynamic DNA methylation in the brain: a new epigenetic mark for experience-dependent plasticity. *Front. Cell. Neurosci* 9, 331. [PubMed: 26379502]
7. Chung E, Cromby J, Papadopoulos D, and Tufarelli C (2017). Social epigenetics: a science of social science? *Socio. Rev* 64, 168–185.
8. Van Speybroeck L (2002). Philosophers and biologists exploring epigenetics. *Biol. Philos* 17, 743–746.
9. Klutstein M, Nejman D, Greenfield R, and Cedar H (2016). DNA methylation in cancer and aging. *Cancer Res.* 76, 3446–3450. [PubMed: 27256564]
10. Verma M, Rogers S, Divi RL, Schully SD, Nelson S, Joseph Su L, Ross SA, Pilch S, Winn DM, and Khoury MJ (2014). Epigenetic research in cancer epidemiology: trends, opportunities, and challenges. *Cancer Epidemiol. Biomarkers Prev* 23, 223–233. [PubMed: 24326628]
11. Jones MJ, Goodman SJ, and Kobor MS (2015). DNA methylation and healthy human aging. *Aging Cell* 14, 924–932. [PubMed: 25913071]
12. Joehanes R, Just AC, Marioni RE, Pilling LC, Reynolds LM, Mandaviya PR, Guan W, Xu T, Elks CE, Aslibekyan S, et al. (2016). Epigenetic signatures of cigarette smoking. *Circ. Cardiovasc. Genet* 9, 436–447. [PubMed: 27651444]
13. Bose M, Wu C, Pankow JS, Demerath EW, Bressler J, Fornage M, Grove ML, Mosley TH, Hicks C, North K, et al. (2014). Evaluation of microarray-based DNA methylation measurement using technical replicates: the Atherosclerosis Risk in Communities (ARIC) Study. *BMC Bioinformatics* 15, 312. [PubMed: 25239148]
14. Dugue PA, English DR, MacInnis RJ, Jung CH, Bassett JK, FitzGerald LM, Wong EM, Joo JE, Hopper JL, Southey MC, et al. (2016). Reliability of DNA methylation measures from dried blood spots and mononuclear cells using the HumanMethylation450k BeadArray. *Sci. Rep* 6, 30317. [PubMed: 27457678]
15. Moran S, Arribas C, and Esteller M (2016). Validation of a DNA methylation microarray for 850,000 CpG sites of the human genome enriched in enhancer sequences. *Epigenomics* 8, 389–399. [PubMed: 26673039]
16. Pidsley R, Zotenko E, Peters TJ, Lawrence MG, Risbridger GP, Molloy P, Van Dijk S, Muhlhauser B, Stirzaker C, and Clark SJ (2016). Critical evaluation of the Illumina MethylationEPIC BeadChip microarray for whole-genome DNA methylation profiling. *Genome Biol.* 17, 208. [PubMed: 27717381]
17. Logue MW, Smith AK, Wolf EJ, Maniates H, Stone A, Schichman SA, McGlinchey RE, Milberg W, and Miller MW (2017). The correlation of methylation levels measured using Illumina 450K and EPIC BeadChips in blood samples. *Epigenomics* 9, 1363–1371. [PubMed: 28809127]
18. Solomon O, MacIsaac J, Quach H, Tindula G, Kobor MS, Huen K, Meaney MJ, Eskenazi B, Barcellos LF, and Holland N (2018). Comparison of DNA methylation measured by Illumina 450K and EPIC BeadChips in blood of newborns and 14-year-old children. *Epigenetics* 13, 655–664. [PubMed: 30044683]
19. Ioannidis JPA (2019). Neglecting major health problems and broadcasting minor, uncertain issues in lifestyle science. *JAMA.* 10.1001/jama.2019.17576.

20. Koo TK, and Li MY (2016). A guideline of selecting and reporting Intraclass correlation coefficients for reliability research. *J. Chiropr. Med* 15, 155–163. [PubMed: 27330520]
21. Cicchetti DV, and Sparrow SA (1981). Developing criteria for establishing interrater reliability of specific items: applications to assessment of adaptive behavior. *Am. J. Ment. Defic* 86, 127–137. [PubMed: 7315877]
22. Edgar RD, Jones MJ, Robinson WP, and Kobor MS (2017). An empirically driven data reduction method on the human 450K methylation array to remove tissue specific non-variable CpGs. *Clin. Epigenetics* 9, 11. [PubMed: 28184257]
23. van Dongen J, Ehli EA, Slieker RC, Bartels M, Weber ZM, Davies GE Slagboom PE, Heijmans BT, and Boomsma DI (2014). Epigenetic variation in monozygotic twins: a genome-wide analysis of DNA methylation in buccal cells. *Genes (Basel)* 5, 347–365. [PubMed: 24802513]
24. Hannon E, Knox O, Sugden K, Burrage J, Wong CCY, Belsky DW, Corcoran DL, Arseneault L, Moffitt TE, Caspi A, et al. (2018). Characterizing genetic and environmental influences on variable DNA methylation using monozygotic and dizygotic twins. *PLoS Genet.* 14, e1007544. [PubMed: 30091980]
25. Naeem H, Wong NC, Chatterton Z, Hong MK, Pedersen JS, Corcoran NM, Hovens CM, and Macintyre G (2014). Reducing the risk of false discovery enabling identification of biologically significant genome-wide methylation status using the HumanMethylation450 array. *BMC Genomics* 15, 51. [PubMed: 24447442]
26. Du P, Zhang X, Huang CC, Jafari N, Kibbe WA, Hou L, and Lin SM (2010). Comparison of Beta-value and M-value methods for quantifying methylation levels by microarray analysis. *BMC Bioinformatics* 11, 587. [PubMed: 21118553]
27. Burt SA (2009). Rethinking environmental contributions to child and adolescent psychopathology: a meta-analysis of shared environmental influences. *Psychol. Bull* 135, 608–637. [PubMed: 19586164]
28. McRae AF, Marioni RE, Shah S, Yang J, Powell JE, Harris SE, Gibson J, Henders AK, Bowdler L, Painter JN, et al. (2018). Identification of 55,000 replicated DNA methylation QTL. *Sci. Rep* 8, 17605. [PubMed: 30514905]
29. Zeilinger S, Kuhnel B, Klopp N, Baurecht H, Kleinschmidt A, Gieger C, Weidinger S, Lattka E, Adamski J, Peters A, et al. (2013). Tobacco smoking leads to extensive genome-wide changes in DNA methylation. *PLoS One* 8, e63812. [PubMed: 23691101]
30. Besingi W, and Johansson A (2014). Smoke-related DNA methylation changes in the etiology of human disease. *Hum. Mol. Genet* 23, 2290–2297. [PubMed: 24334605]
31. Guida F, Sandanger TM, Castagne R, Campanella G, Polidoro S, Palli D, Krogh V, Tumino R, Sacerdote C, Panico S, et al. (2015). Dynamics of smoking-induced genome-wide methylation changes with time since smoking cessation. *Hum. Mol. Genet.* 24, 2349–2359. [PubMed: 25556184]
32. Wilson R, Wahl S, Pfeiffer L, Ward-Caviness CK, Kunze S, Kretschmer A, Reischl E, Peters A, Gieger C, and Waldenberger M (2017). The dynamics of smoking-related disturbed methylation: a two time-point study of methylation change in smokers, non-smokers and former smokers. *BMC Genomics* 18, 805. [PubMed: 29047347]
33. Dogan MV, Shields B, Cutrona C, Gao L, Gibbons FX, Simons R, Monick M, Brody GH, Tan K, Beach SR, et al. (2014). The effect of smoking on DNA methylation of peripheral blood mononuclear cells from African American women. *BMC Genomics* 15, 151. [PubMed: 24559495]
34. Dogan MV, Beach SRH, and Philibert RA (2017). Genetically contextual effects of smoking on genome wide DNA methylation. *Am. J. Med. Genet. B Neuropsychiatr. Genet* 174, 595–607. [PubMed: 28686328]
35. Horvath S, and Raj K (2018). DNA methylation-based biomarkers and the epigenetic clock theory of ageing. *Nat. Rev. Genet* 19, 371–384. [PubMed: 29643443]
36. Hannum G, Guinney J, Zhao L, Zhang L, Hughes G, Sada S, Klotzle B, Bibikova M, Fan JB, Gao Y, et al. (2013). Genome-wide methylation profiles reveal quantitative views of human aging rates. *Mol. Cell* 49, 359–367. [PubMed: 23177740]
37. Horvath S (2013). DNA methylation age of human tissues and cell types. *Genome Biol.* 14, R115. [PubMed: 24138928]

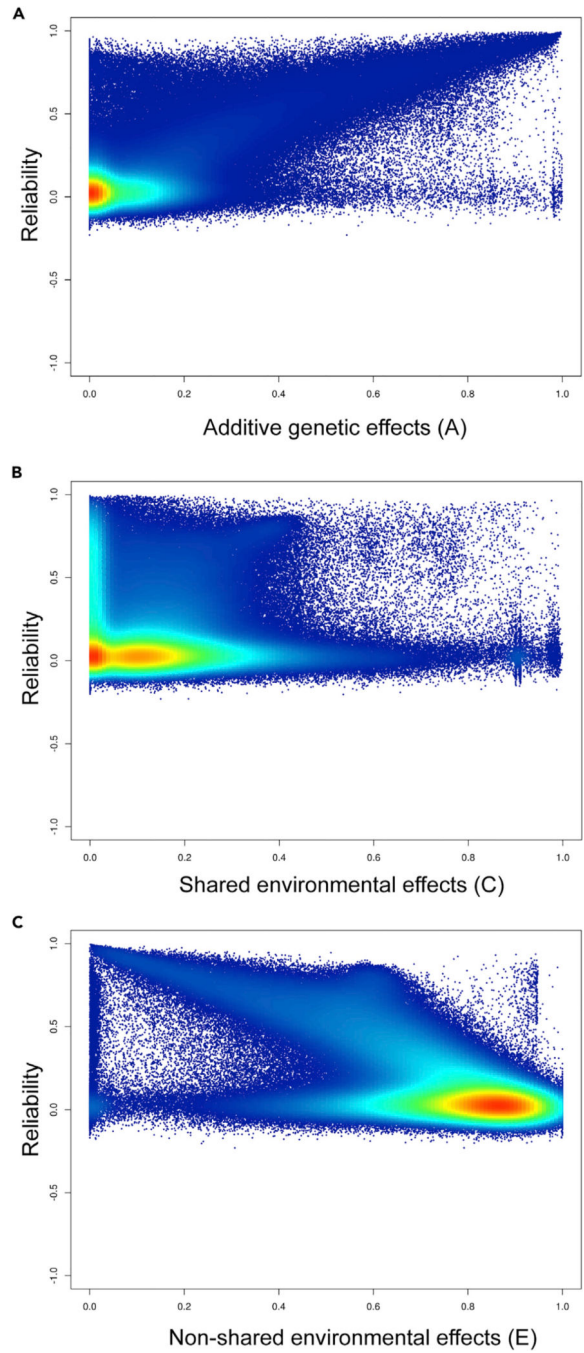
38. Levine ME, Lu AT, Quach A, Chen BH, Assimes TL, Bandinelli S, Hou L, Baccarelli AA, Stewart JD, Li Y, et al. (2018). An epigenetic biomarker of aging for lifespan and healthspan. *Aging (Albany NY)* 10, 573–591. [PubMed: 29676998]
39. Teschendorff AE, and Relton CL (2018). Statistical and integrative system-level analysis of DNA methylation data *Nat. Rev. Genet* 19, 129–147. [PubMed: 29129922]
40. Kennedy EM, Goehring GN, Nichols MH, Robins C, Mehta D, Klengel T, Eskin E, Smith AK, and Conneely KN (2018). An integrated -omics analysis of the epigenetic landscape of gene expression in human blood cells *BMC Genomics* 19, 476. [PubMed: 29914364]
41. Hannon E, Lunnon K, Schalkwyk L, and Mill J (2015). Interindividual methylomic variation across blood, cortex, and cerebellum: implications for epigenetic studies of neurological and neuropsychiatric phenotypes. *Epigenetics* 10, 1024–1032. [PubMed: 26457534]
42. Visscher PM, Wray NR, Zhang Q, Sklar P, McCarthy MI, Brown MA, and Yang J (2017). 10 Years of GWAS discovery: biology, function, and translation. *Am. J. Hum. Genet* 101, 5–22. [PubMed: 28686856]
43. Lehne B, Drong AW, Loh M, Zhang W, Scott WR, Tan ST, Afzal U, Scott J, Jarvelin MR, Elliott P, et al. (2015). A coherent approach for analysis of the Illumina HumanMethylation450 BeadChip improves data quality and performance in epigenome-wide association studies. *Genome Biol.* 16, 37. [PubMed: 25853392]
44. Yan H, Tian S, Slager SL, Sun Z, and Ordog T (2016). Genome-wide epigenetic studies in human disease: a primer on -omic technologies. *Am. J. Epidemiol* 183, 96–109. [PubMed: 26721890]
45. The National Academies of Sciences, Engineering, and Medicine (2019). *Reproducibility and Replicability in Science* (National Academies Press).
46. Moffitt TE; E-Risk Study Team (2002). Teen-aged mothers in contemporary Britain. *J. Child Psychol. Psychiatry* 43, 727–742. [PubMed: 12236608]
47. Poulton R, Moffitt TE, and Silva PA (2015). The Dunedin multidisciplinary health and development study: overview of the first 40 years, with an eye to the future. *Soc. Psychiatry Psychiatr. Epidemiol* 50, 679–693. [PubMed: 25835958]
48. Marzi SJ, Sugden K, Arseneault L, Belsky DW, Burrage J, Corcoran DL, Danese A, Fisher HL, Hannon E, Moffitt TE, et al. (2018). Analysis of DNA methylation in young people: limited evidence for an association between victimization stress and epigenetic variation in blood. *Am. J. Psychiatry* 175, 517–529. [PubMed: 29325449]
49. Sergushichev AA (2016). An algorithm for fast preranked gene set enrichment analysis using cumulative statistic calculation. *bioRxiv*. 10.1101/060012.
50. Neale MC, Hunter MD, Pritikin JN, Zahery M, Brick TR, Kirkpatrick RM, Estabrook R, Bates TC, Maes HH, and Boker SM (2016). OpenMx 2.0: extended structural equation and statistical modeling. *Psychometrika* 81, 535–549. [PubMed: 25622929]

### THE BIGGER PICTURE

Although DNA methylation data are used widely by researchers in many fields, the reliability of these data are surprisingly variable. Our findings remind us that, in an age of increasingly big data, research is only as robust as its foundations. We hope that our findings will improve the integrity of DNA methylation studies. We also hope that our findings serve as a cautionary reminder for those generating and implementing big data of any type: reliability is a fundamental aspect of replicability. Conducting analysis with reliable data will improve chances of replicable findings, which might lead to more actionable targets for further research. To the extent that reliable data improve replicability, the knock-on effect will be more public confidence in research and less effort spent trying to replicate findings that are bound to fail.

### Highlights

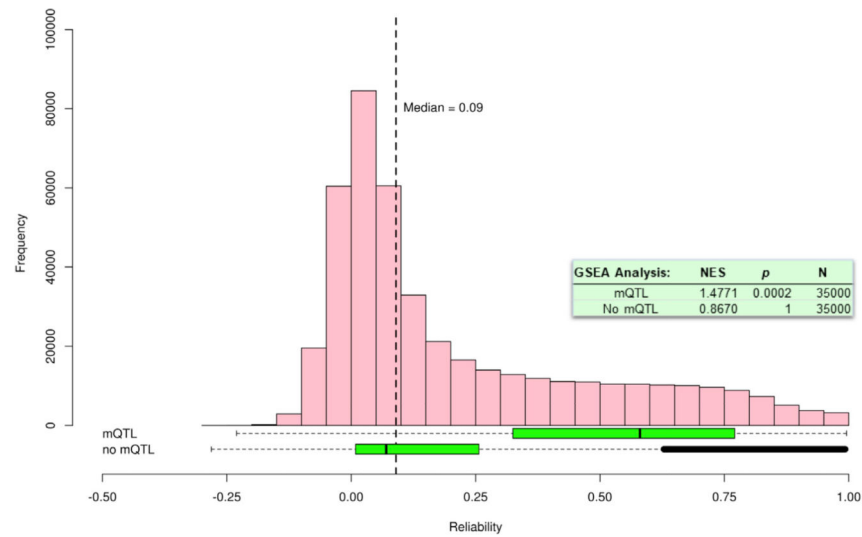
- Measurements of DNA methylation made using BeadChip probes are differentially reliable
- Unreliable probes were less heritable, less replicable, and less functionally relevant
- This has serious implications for reporting and evaluating DNA methylation findings
- Reliability joins replicability and reproducibility to make three fundamental Rs of STEM



**Figure 1. Density Heatmap of Probe Reliability Plotted against Estimates of Genetic and Environmental Effects on DNA Methylation**

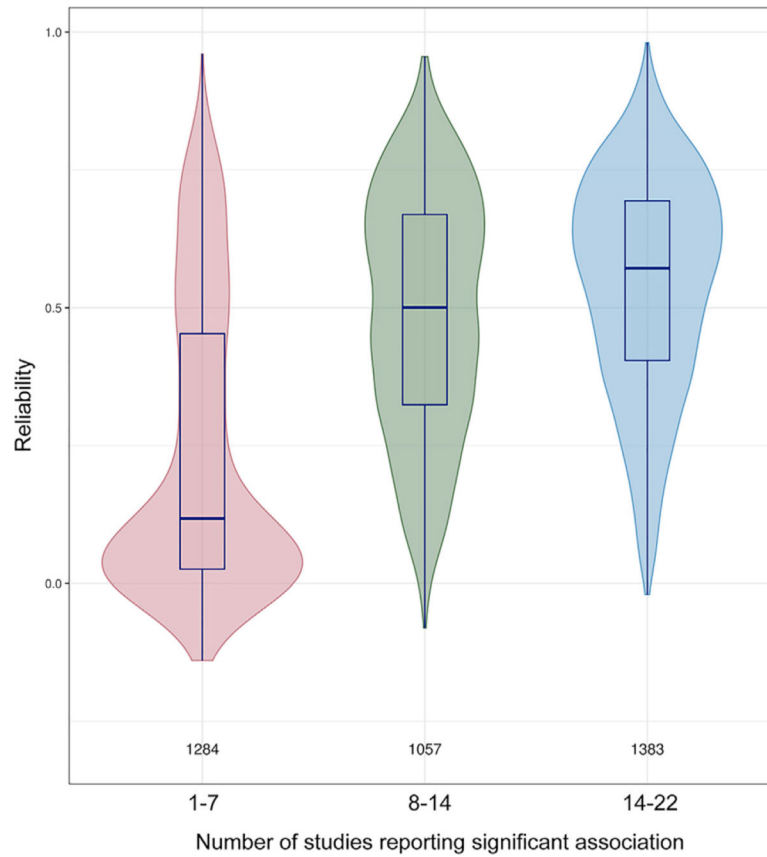
(A) Additive genetic effects (denoted as “A”), (B) shared environmental effects (denoted as “C”), and (C) non-shared (or unique) environmental effects (denoted as “E”). The variance component is plotted on the x axis and the reliability is plotted on the y axis. Probes with the highest reliability have the highest value of A and lowest value of E. Density is depicted on a spectral scale from low (dark blue) to high (red).





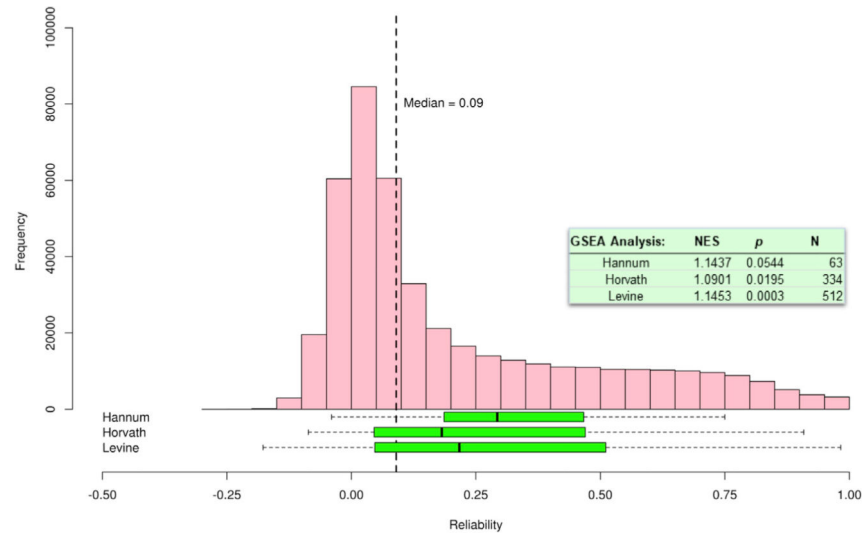
**Figure 2. The Distribution of Reliabilities of Probes Identified in a Large-Scale mQTL Analysis Compared with Non-mQTL Probes**

Distributions are depicted as box-and-whisker plots of the reliability coefficients of the probes identified as having mQTLs (“mQTL”) and the remainder not included in the mQTL list (“no mQTL”). Boxes correspond to interquartile range (IQR), and whiskers extend to  $1.5 \times$  IQR. Observations beyond the whiskers (outliers) are represented by individual points. As a reference, the distribution (pink bars) and median (vertical dashed line) of all ~440,000 probe reliabilities in the E-Risk dataset is shown above the box-and-whisker plots. The text box shows the results of gene set enrichment analysis (GSEA; NES, normalized enrichment score; N, number of probes); probes associated with mQTLs are enriched for reliable probes, suggesting that reliable probe measurement is important for uncovering genetic effects on methylation.



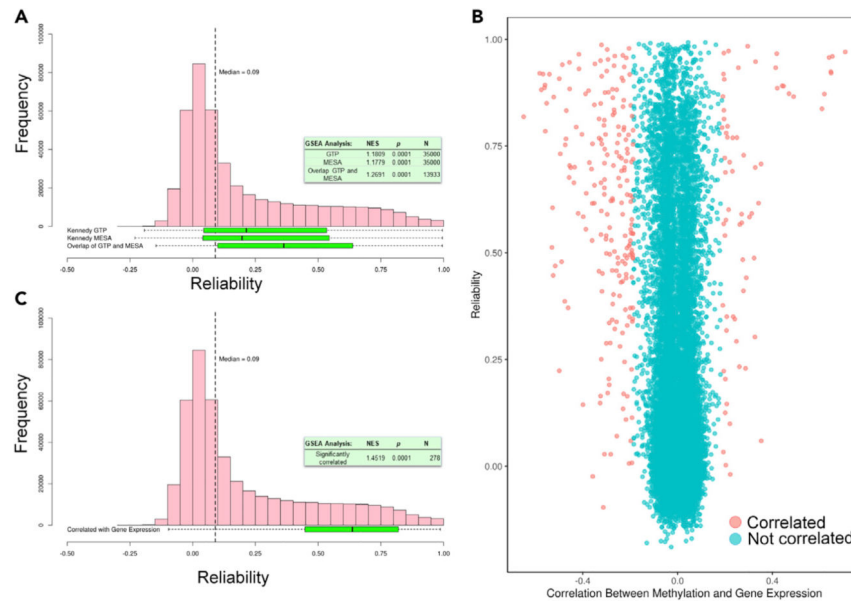
**Figure 3. Probes Consistently Associated with Smoking across Studies Have Higher Reliabilities Than Probes that Are Not**

We identified 22 epigenome-wide association studies of smoking and DNA methylation. For ease of visualization, probes have been binned into three groups representing 1–7 replications (pink), 8–14 replications (green), and 15–22 replications (blue). The values above the x axis represent the number of probes per group. In the 1–7 replication bin, the highest density of probes was at the low-reliability end of the distribution, and the median reliability (as depicted by the median line of the box plot within the violin) was the lowest of the three groups. Boxes correspond to IQR and whiskers extend to  $1.5 \times$  IQR.



**Figure 4. Reliabilities of Probes Included in Established, Publicly Available DNA Methylation Algorithms (“Clocks”)**

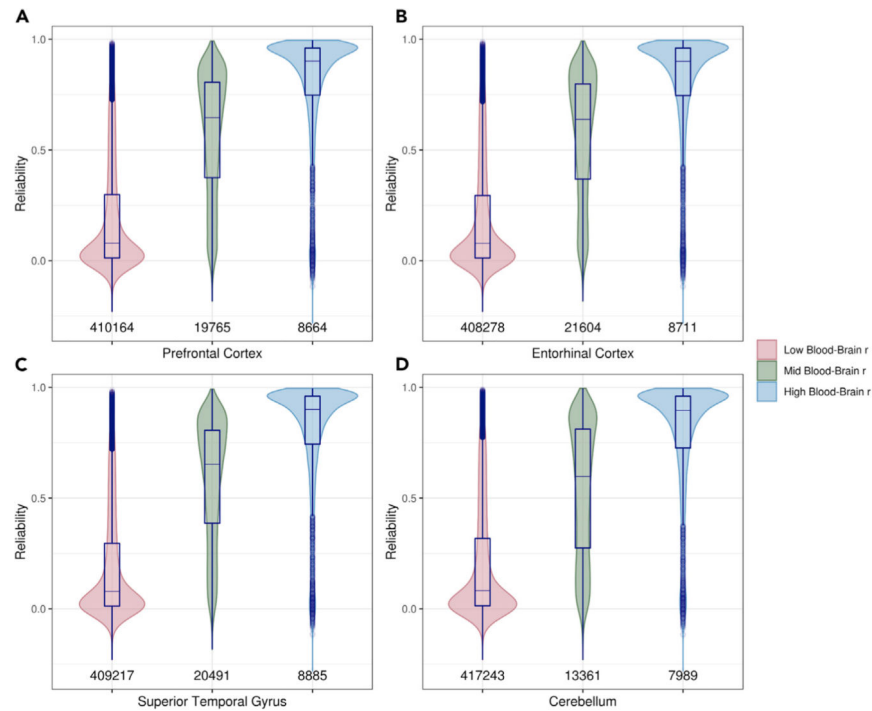
Distributions are depicted as box-and-whisker plots of the reliability coefficients of the probe constituents of the Hannum et al.<sup>36</sup> aging clock (63 probes), Horvath<sup>37</sup> DNAMAge clock (334 probes), and Levine et al.<sup>38</sup> biological aging clock (512 probes). Boxes correspond to IQR and whiskers extend to  $1.5 \times$  IQR. Observations beyond the whiskers (outliers) are represented by individual points. As a reference, the distribution (pink bars) and median (vertical dashed line) of all  $\sim 440,000$  probe reliabilities in the E-Risk dataset is shown above the box-and-whisker plots. The aging clocks are enriched for reliable probes (values to the right of the figure; NES, normalized enrichment score; N, number of probes). Median reliabilities of probes included in aging clocks are higher than those of the general distribution; however, each algorithm contained many unreliable probes.



**Figure 5. Reliabilities of Probes Significantly Correlated with Gene Expression Have Higher Reliabilities Than Non-correlated Probes**

(A) Distributions of the reliability coefficients of the probes identified as correlated with gene expression by Kennedy et al.<sup>40</sup> in the GTP and MESA cohorts (N probes = 36,485 and 114,536, respectively). Probes that are correlated with gene expression in both cohorts are shown in the bottommost box-and-whisker plot. Boxes correspond to IQR and whiskers extend to  $1.5 \times$  IQR. As a reference, the distribution (pink bars) and median (vertical dashed line) of all ~440,000 probe reliabilities in the E-Risk dataset is shown above the box-and-whisker plots. The text box shows the results of GSEA for the GTP cohort, MESA cohort, and the intersection of both cohorts (NES, normalized enrichment score; N, number of probes). Each cohort's set of significantly correlated DNA methylation probe-gene expression pairs is enriched for reliable probes; pairs that are significantly correlated in both datasets are further enriched.

(B) TSS-localized DNA methylation probe-gene expression probeset correlation (x axis) plotted against DNA methylation probe reliability (y axis) in the Dunedin Study dataset. Probes that were significantly correlated with gene expression are shown in pink (n = 278) and those not correlated are shown in blue. (C) Distribution of reliabilities of these significantly correlated DNA methylation probes as a box-and-whisker plot. The text box shows the results of GSEA (NES, normalized enrichment score; N, number of probes); DNA methylation probes that were significantly correlated with expression probesets are enriched for reliable probes.



**Figure 6. Violin Plots of the Distribution of Reliability in Probes with Low (<0.4, Pink), Medium (0.4–0.75, Green), and High (>0.75, Blue) Blood-Brain Correlation in DNA Methylation**  
 Distributions are shown across four brain regions: prefrontal cortex (A), entorhinal cortex (B), superior temporal gyrus (C), and cerebellum (D). Number of probes per group is listed above the x axis. Box-and-whisker plots of the distribution are plotted within violin plots. Values below each violin correspond to the number of probes in that group. Probes with high blood-brain concordance are concentrated at the high-reliability end of the distribution. Boxes correspond to IQR and whiskers extend to  $1.5 \times$  IQR.