

# **Title: Tissue- and ethnicity-independent hypervariable DNA methylation states show evidence of establishment in the early human embryo**

**Authors:** Maria Derakhshan<sup>1</sup>, Noah J. Kessler<sup>2</sup>, Miho Ishida<sup>3</sup>, Charalambos Demetriou<sup>3</sup>, Nicolas Brucato<sup>4</sup>, Gudrun E. Moore<sup>3</sup>, Caroline H.D. Fall<sup>5</sup>, Giriraj R. Chandak<sup>6</sup>, Francois-Xavier Ricaut<sup>4</sup>, Andrew M. Prentice<sup>7</sup>, Garrett Hellenthal<sup>8\*</sup> and Matt J. Silver<sup>1,7\*</sup>

## **Affiliations:**

<sup>1</sup>London School of Hygiene and Tropical Medicine, UK.

<sup>2</sup>Department of Genetics, University of Cambridge, Cambridge, CB2 3EH, UK.

<sup>3</sup>UCL Great Ormond Street Institute of Child Health.

<sup>4</sup>Laboratoire Évolution and Diversité Biologique (EDB UMR 5174), Université de Toulouse Midi-Pyrénées, CNRS, IRD, UPS, Toulouse, France.

<sup>5</sup>MRC Lifecourse Epidemiology Unit, University of Southampton, Southampton, United Kingdom

<sup>6</sup>Genomic Research on Complex Diseases (GRC Group), CSIR-Centre for Cellular and Molecular Biology, Hyderabad, India.

<sup>7</sup>Medical Research Council Unit The Gambia at the London School of Hygiene and Tropical Medicine, The Gambia.

<sup>8</sup>UCL Genetics Institute, University College London, Gower Street, London WC1E 6BT, UK.

\*Corresponding authors

## **Abstract**

We analysed DNA methylation data from 30 datasets comprising 3,474 individuals, 19 tissues and 8 ethnicities at CpGs covered by the Illumina450K array. We identified 4,143 hypervariable CpGs (“hvCpGs”) with methylation in the top 5% most variable sites across multiple tissues and ethnicities. hvCpG methylation was influenced but not determined by genetic variation, and was not linked to probe reliability, epigenetic drift, age, sex or cell heterogeneity effects. hvCpG methylation tended to covary across tissues derived from different germ-layers and hvCpGs were enriched for associations with periconceptional environment, proximity to ERV1 and ERVK retrovirus elements and parent-of-origin-specific methylation. They also showed distinctive methylation signatures in monozygotic twins. Together, these properties position hvCpGs as strong candidates for studying how stochastic and/or environmentally influenced DNA methylation states which are established in the early embryo and maintained stably thereafter can influence life-long health and disease.

## Introduction

DNA methylation (DNAm) plays a critical role in mammalian development, underpinning X-chromosome inactivation, genomic imprinting, silencing of repetitive regions and cell differentiation<sup>1</sup>. DNAm states that vary between individuals have been a focus of Epigenome-Wide Association Studies (EWAS) due to their potential to drive phenotypic variation<sup>2,3</sup>. Factors influencing interindividual methylation differences include genetic variation<sup>4,5</sup>, cell heterogeneity effects<sup>6,7</sup>, sex<sup>8,9</sup>, age<sup>10,11</sup>, and pre- and post- natal environment<sup>12–14</sup>. Growing evidence from studies investigating DNAm patterns in multiple tissues suggests that these factors have both shared and tissue-specific influences on DNAm variation<sup>12,15–18</sup>.

In this study, we sought to identify loci with high interindividual methylation variability in multiple tissues and ethnicities, and to gain insights into the biological mechanisms influencing methylation variation. By using a large number of diverse sample types, we reasoned that identified loci would be robust to tissue-specific drivers of methylation variability such as those mentioned above, and to dataset-specific technical artefacts, including batch effects and poorly performing probes<sup>19–22</sup>. We began by characterising hypervariable CpGs ('hvCpGs') covered on the widely used Illumina HumanMethylation450K (hereafter 'Illumina450K') array<sup>23</sup> that showed high interindividual variability across multiple datasets covering 19 different tissue/cell types and 8 ethnicities spanning a wide range of ages. We next investigated the influence of genetic variation, sex, age and probe reliability on methylation variability at hvCpGs. We additionally determined whether methylation states at hvCpGs covary across tissues by exploring their overlap with loci at which methylation varies between individuals but is correlated across tissues within a given individual, termed systemic interindividual variation or 'SIV'. Since loci showing SIV have been linked to variable methylation establishment before germ-layer differentiation<sup>24–27</sup>, we further explored evidence for early embryo methylation at hvCpGs by determining their overlap with loci that show unique methylation patterns in MZ twins<sup>25,28</sup> and with loci that show sensitivity to the periconceptual environment<sup>29</sup>. We assessed the genomic context of hvCpGs by exploring their association with multi-tissue histone marks and their proximity to transposable elements and regions of parent-of-origin-specific methylation. Finally, we probed putative functional roles of hvCpGs by interrogating EWAS trait associations and by performing gene ontology enrichment analysis.

Our curated list of hvCpGs show evidence of establishment in the early embryo and of correlation across tissues. They therefore serve as a useful resource for studying the influence of early environmental and/or stochastic effects on DNAm in diverse tissues and ethnicities, and for studying the impact of DNAm differences on life-long health and disease.

## Materials and Methods

### Methylation data used for identifying hvCpGs

Publicly available methylation data were downloaded from The Cancer Genome Atlas (TCGA) (<https://www.cancer.gov/tcga>) and the Gene Expression Omnibus (GEO)<sup>30</sup> databases as methylation Beta matrices (Supplementary Tables 1 and 2). TCGA methylation data were downloaded using the *TCGAbiolinks* (v2.18.0) R package<sup>31–33</sup>, selecting only samples annotated as ‘Solid Tissue Normal’. Of the 33 TCGA datasets, 10 were selected for our study as these had methylation data in at least 20 samples. GEO methylation Beta matrices were downloaded from 11 unique accessions using the *GEOquery* (v2.58.0) R package<sup>34</sup>. Where available, detection p-values (measuring signal intensity), and metadata on age, sex, and disease status were also downloaded. We split GEO beta matrices into separate groups based on ethnicity and tissue/cell type and refer to the resulting 17 separated groups as ‘datasets’. Non-public datasets internal to this study include IlluminaEPIC<sup>35</sup> array data from whole blood samples from Gambian 8-9-year olds (ISRCTN14266771<sup>36</sup>) and Illumina450K data from Bornean and Kenyan saliva samples<sup>37</sup> (Supplementary Table 3). For IlluminaEPIC datasets we selected probes covered on the Illumina450K array. In total, we analysed 30 datasets (3 internal, 10 TCGA and 17 GEO) that covered 8 ethnicities and 19 different tissue/cell types (Supplementary Table 4).

### Methylation data processing

For each methylation dataset used in our main analysis, we used the *ChAMP* (v2.20.1) R package<sup>38</sup> to remove: i) probes with a detection p-value > 0.01 in > 5% samples (where detection p-values were available), ii) probes mapping to multiple genomic positions<sup>39</sup>, iii) probes mapping to the X and Y chromosomes, and iv) single nucleotide polymorphism (SNP)-related probes identified by Zhou *et al.*<sup>39</sup> that contain SNPs (MAF > 1%) within 5 bp of the CpG interrogation site and/or SNPs effecting probe hybridisation. Where ethnicity information was available, we removed probes with population-specific SNPs identified by Zhou *et al.* using 1000 Genomes populations (MAF > 1%), otherwise we removed the General Recommended Probes<sup>40</sup>. Probes that had a missing value in any of the samples in a specific dataset were removed from that dataset. To reduce technical biases introduced by differing type I and type II probe designs on the Illumina450K and IlluminaEPIC arrays, we applied Beta Mixture Quantile normalisation (BMIQ)<sup>41</sup> to all datasets using the *champ.norm()* function from the *ChAMP* R package. All datasets were adjusted for the first 10 principal components (PCs) of variation to account for methylation variability driven by known and/or unknown technical artefacts (such as plate and array position) and cell heterogeneity. Methylation values were adjusted for these 10 PCs, age (where available) and sex by taking the residuals from a linear regression on methylation M values, where M is defined as  $\log_2(\text{beta}/1-\text{beta})$ . Finally, for each probe, we removed outlier methylation values, defined according to Tukey’s outer fences ( $Q1 - 3 \cdot \text{IQR}$  and  $Q3 + 3 \cdot \text{IQR}$ ). The hg19 reference genome was used throughout all relevant analyses as the Illumina450K array metadata manifest uses this version.

## Identification of hvCpGs

For each dataset, we identified CpGs within the top  $i\%$  of CpGs by methylation Beta variance in  $\geq j\%$  of datasets in which the CpG was covered, for increasing values of  $i$  and  $j$ . We additionally required that selected CpGs were covered in a minimum of 15 datasets (Supplementary Fig. 1A). To define the ethnicity- and tissue-independent hypervariable CpGs (hvCpGs) explored in this paper, we set the threshold at  $i, j = 5,65$  (Supplementary Fig. 1B).

## Probe reliability

Technically unreliable probes were identified by examining intra-class correlation coefficients (ICCs) from two studies. The first study compared methylation consistency between the Illumina450K and IlluminaEPIC platforms using 365 blood DNA samples, defining poor quality probes as those with  $ICC \leq 0.4^{22}$ . The second study examined methylation reliability between technical replicates from 265 African American peripheral blood leukocyte samples on the Illumina450K platform, defining poor quality probes as those with  $ICC \leq 0.37^{42}$ . We defined technically unreliable probes as those reported as being poor quality in at least one of these two studies (Supplementary Table 5).

## Methylation quantitative trait locus (mQTL) analysis

mQTL summary statistics from the Genetics of DNA Methylation Consortium (GoDMC), a meta-GWAS of 36 European blood cohorts ( $N = 27,750$ ) generated using imputed genotype data ( $\sim 10$  million SNPs) and  $\sim 420,000$  CpGs<sup>43</sup> were used for this analysis. Significance thresholds of  $p < 1 \times 10^{-8}$  and  $p < 1 \times 10^{-14}$  were applied for *cis* and *trans* mQTLs respectively<sup>43</sup>, giving 271,724 significant SNP-CpG associations comprising 190,102 CpGs and 224,648 SNPs. The variance in DNA methylation explained by a given mQTL was estimated as  $2 * \beta * MAF(1-MAF)$ , where  $\beta$  is the effect size and MAF is the minor allele frequency<sup>44</sup>.

## Monozygotic twin discordance

We analysed CpGs identified as being ‘equivalently variable’ between MZ co-twins and between unrelated individuals (‘evCpGs (blood)’ by Planterose Jiménez *et al.*<sup>45</sup> using Illumina450K data in whole blood. 154 of these evCpGs replicated in adipose tissue from 97 MZ twin pairs (‘evCpGs (blood & adipose)’). evCpGs are candidates for methylation states that are established stochastically after MZ twin splitting.

## Control CpG sets

### *Distribution-matched controls*

hvCpGs are enriched for intermediate methylation states (Supplementary Fig. 2). This property of hvCpGs has the potential to bias several downstream analyses, for example because this can affect power to find association with phenotypes in EWAS. We therefore constructed a set of CpGs with similar distribution of methylation Beta values to hvCpGs in the Caucasian blood dataset



(‘Blood\_Cauc’, Supplementary Table 1). This dataset was chosen as it has the highest number of post-natal samples and because several downstream analyses leverage published studies that used blood methylation data. For each of the 4,108 hvCpGs covered in the ‘Blood\_Cauc’ dataset, a two-sided Kolmogorov-Smirnov (KS) test (*ks.test()* in R) was used to test for the divergence in methylation Beta distributions between the hvCpG and technically reliable (see ‘Probe reliability’, Methods) background probes, selecting the CpG with the greatest p-value (requiring a p-value > 0.1). In total, 3,566 hvCpGs were each matched to a control CpG (‘distribution-matched controls’, Table 1, Supplementary Fig. 3).

### ***mQTL-matched controls***

To determine the degree to which hypervariability at hvCpGs is explained by mQTL effects, each hvCpG was matched to a CpG amongst those reported in the GoDMC meta-analysis<sup>43</sup>. Controls were selected to have i) the same number of mQTL associations, ii) a similar mean % variance explained by mQTL (across all significant mQTL) and iii) presence in at least as many datasets as the hvCpG (Table 1, Supplementary Fig. 4).

## **Identification of hvCpG clusters**

hvcpg clusters were identified by considering the decay of methylation correlation with distance at hvCpGs. To do this, we calculated the average pairwise Spearman correlation ( $\rho$ ) across hvCpG pairs with inter-CpG distance falling within 100 bp bins, for datasets with at least 100 samples (Supplementary Fig. 5B). The distance threshold for defining hvCpG clusters was chosen to be 4,000 bp as this is approximately the point at which pairwise correlations levelled out (Supplementary Fig. 5B). In total, 2,219 (54%) hvCpGs fell into 716 clusters comprising at least 2 CpGs, with the remaining 1,924 (46%) hvCpGs falling outside of these clusters (Supplementary Fig. 5C). In 563 (79%) of these clusters, the average Spearman correlation ( $\rho$ ) across hvCpG pairs was > 0.5 (Supplementary Fig. 5D).

### **‘De-clustering’ of hvCpGs**

To account for the possibility that our analyses may be biased by the non-random distribution and inter-dependence of hvCpGs in CpG clusters, we generated a de-clustered set of hvCpGs in which no CpG was within 4 kb of another CpG. 2,640 de-clustered hvCpGs were generated by randomly selecting one CpG from each of the clusters and then including all ‘singleton’ CpGs falling outside of clusters.

## **Age stability**

To examine temporal stability of hvCpGs we used published intra-class correlation coefficients (ICCs) for probes on the Illumina450K array determined using white blood cell samples taken ~6 years apart<sup>46</sup>. The ICC scores compare within-sample variability (across the two time-points) to between-sample variability, with

ICC  $\geq 0.5$  defined as temporally stable by Flanagan *et al.*<sup>46</sup>. Because methylation data from the Flanagan *et al.* dataset were publicly available (GSE61151), we compared ICC scores at hvCpGs to those at CpGs with similar methylation Beta distributions to hvCpGs at the first time point (Supplementary Fig. 6A) to ensure that high hvCpG ICC scores were not biased by the high variability of hvCpGs. These CpGs were matched to each hvCpG using the same Kolmogorov-Smirnov method detailed in ‘Distribution-matched controls’ but using the Flanagan *et al.* methylation data instead the ‘Blood\_Cauc’ dataset<sup>44</sup>. Longer-term susceptibility to epigenetic drift was examined by determining the proportion of hvCpGs that overlap a published set of 6,108 CpGs identified using whole blood Illumina450K data from 3,295 individuals aged 18 to 88 years that show an increased methylation variability with age of more than 5% every 10 years<sup>11</sup> (Supplementary Fig. 6B).

## Published CpG sets used to investigate early embryo establishment

We used the following publicly available data to examine evidence that methylation states at hvCpGs are established in the early embryo. See Table 2 for a summary of these datasets.

### **Systemic Interindividual Variation (‘SIV’) CpGs**

SIV-CpGs were collated from four published datasets that used either whole genome bisulfite sequencing (WGBS) or Illumina450K data from multiple tissues derived from different germ layers to identify CpGs displaying high interindividual variation and low intra-individual (cross-tissue) variation. These properties are suggestive of variable methylation establishment before germ layer differentiation<sup>24–27</sup>. Further details on the four SIV screens used in this study are given in Supplementary Table 6.

### **Epigenetic Supersimilarity (‘ESS’) CpGs**

Epigenetic supersimilarity (ESS) loci were identified by van Baak *et al.*<sup>25</sup> using Illumina450K data from adipose tissue from 97 MZ and 162 dizygotic (DZ) twin pairs<sup>47</sup>. In that study, 1,580 ESS sites were identified within the top decile of methylation variance, with an interindividual methylation range  $> 0.4$  and greater-than-expected concordance in MZ twins vs DZ twins. This supersimilarity is attributed to methylation establishment before MZ twin splitting.

### **MZ twinning CpGs**

Van Dongen *et al.*<sup>48</sup> performed an epigenome-wide association analysis on each of 6 cohorts with methylation data from both MZ and DZ twins (5 blood and 1 buccal) to identify probes differentially methylated between MZ twins and DZ (dizygotic) twins. A meta-analysis was then performed using the blood datasets to identify 834 Bonferroni-significant differentially methylated CpGs, which we refer to as ‘MZ twinning CpGs’.

### **Season of conception ('SoC') CpGs**

Silver *et al.*<sup>29</sup> used Illumina450K data to identify 259 CpGs associated with season-of-conception ('SoC') in Gambian 2-year olds with a minimum methylation difference of 4% between the peaks of the Gambian rainy and seasons.

### **Transposable elements and telomeres**

Locations of ERV1 and ERVK transposable elements determined by RepeatMasker were downloaded from the UCSC annotations repository as previously described<sup>26</sup>. Telomere coordinates were downloaded from the UCSC hg19 annotations repository. (<http://genome.ucsc.edu>).

### **Imprinted genes, parent-of-origin-specific methylation (PofOm)**

Imprinted genes classified as 'predicted' or 'known' were downloaded from <https://www.geneimprint.com>. Parent-of-origin-specific CpGs were identified by Zink *et al.*<sup>49</sup> using WGBS data from peripheral blood from Icelandic individuals.

## **SIV power calculation**

To assess power to detect SIV in previous screens with small numbers of samples, we analysed the 4-individual multi-tissue dataset used by van Baak *et al.*<sup>25,50</sup>. We downloaded this dataset from GEO (GSE50192), selecting the same tissues (gall bladder, abdominal aorta sciatic nerve) used by van Baak *et al.*<sup>25</sup>. For each of the 1,042 SIV-CpGs reported by van Baak *et al.*<sup>25</sup>, we generated methylation values for three tissues for each simulated individual by randomly sampling from a 3-dimensional multivariate normal distribution, with mean equal to the mean of each tissue's sampled methylation values at the CpG, and standard deviation specified by a 3x3 cross-tissue co-variance matrix of the sampled methylation values at the CpG. For each SIV-CpG, we sampled four simulated individuals and determined if this random sample met the SIV definition specified by van Baak *et al.*<sup>25</sup>, repeating this process 1000 times to give a power estimate (Supplementary Fig. 7).

## **Processing and analysis of fetal multi-tissue dataset**

The fetal multi-tissue dataset comprised 60 samples, corresponding to 30 individuals each with methylation data from two tissues derived from different germ layers (ectoderm: brain, spinal cord, skin; mesoderm: kidney, rib, heart, tongue; endoderm: intestine, gut, lung, liver). These fetal tissues were obtained from the 'Moore Fetal Cohort' from the termination of pregnancies at Queen Charlotte's and Chelsea Hospital (London, UK). Ethical approval for obtaining fetal tissues was granted by the Research Ethics Committee of the Hammersmith, Queen Charlotte's and Chelsea and Acton Hospitals (2001/6028). DNA was extracted from fetal tissues using the AllPrep DNA/RNA/Protein Mini Kit (Qiagen) and bisulfite conversion was carried out using EZ DNA Methylation Kits (Zymo Research). Samples were then processed using the Illumina

InfiniumEPIC array. Derived methylation data were imported as *.idat* files into R and analysed using the *meffil* R package (v 1.1.2)<sup>51</sup> with default parameters. Briefly, methylation predicted sex was used to remove 2 sex outliers (samples with methylation > 5 SDs from mean). Next, 1 sample was removed for which the predicted median methylation signal was more than 3 SDs from the expected signal, leaving 57 samples. 515 probes with detection-p-value value > 0.1 and 307 probes with bead number < 3 in more than 20% of samples respectively were removed. Array data were then corrected for dye-bias and background effects and functional normalisation was applied, specifying the number of PCs to be 7 (the PC at which the variance explained at control probes levelled out). Next, the *ChAMP* (v2.20.1) R package<sup>38</sup> was used to remove cross-hybridising and multi-mapping probes, probes on XY chromosomes, and SNP-related probes, leaving 746,492 CpGs. We selected the 452,016 probes that overlapped the Illumina450K array and the 27 individuals for which both tissue samples passed quality control: 9 individuals with methylation data from endoderm and mesoderm, 10 individuals with methylation data from endoderm and ectoderm and 8 individuals with methylation data from mesoderm and ectoderm (see Supplementary Table 7). Methylation was then adjusted for predicted sex and batch using a linear model. For the 9 individuals with available endoderm-mesoderm samples we calculated the Pearson r between germ layer methylation values for each hvCpG, and repeated this for individuals with endoderm-ectoderm and mesoderm-ectoderm samples. The inter-germ layer correlation was then defined as the average Pearson r across these three comparisons. Following van Baak *et al.*<sup>25</sup>, interindividual variation was determined by calculating the mean methylation value across both tissues within each of the 27 individuals, before taking the range of these means for every CpG.

### Chromatin states at hvCpGs

Chromatin states were predicted by a ChromHMM 15-state model<sup>52</sup> using Chromatin Immunoprecipitation Sequencing (ChIP-Seq) data generated by the Roadmap Epigenomics Consortium<sup>53</sup>. These data were downloaded for H1 ESCs (E003), fetal brain (E071), fetal muscle (E090), fetal small intestine (E085), foreskin fibroblasts (E055), adipose (E063) and primary mononuclear cells (E062) from the Washington University Roadmap repository. Chromatin states were collapsed into 8 states for clarity (Supplementary Table 8).

### EWAS trait associations at hvCpGs

hvCpG trait associations were determined using the EWAS catalogue (<http://ewascatalog.org/>), which details significant results ( $p\text{-value} < 1 \times 10^{-4}$ ) from published EWAS studies. Considering only those traits for which at least 1% of hvCpGs overlapped associated CpGs (highlighted in green in Supplementary Table 9), we first extracted the array background CpGs overlapping the 'Blood\_Cauc' dataset that were associated with each trait. We then calculated the proportion of these CpGs that comprised hvCpGs and blood distribution-matched controls (Table 1). Traits that were significantly enriched or depleted for hvCpGs relative to controls were those for which bootstrapped 95% confidence intervals did not overlap.

### **Gene ontology term enrichment analysis**

Gene Ontology (GO) term enrichment analysis was performed using the *missMethyl* R package (v1.24.0)<sup>54</sup> using the *gometh()* function, setting arguments sig.cpg = hvCpGs, all.cpg = array.background, sig.genes = T, collection = "GO", array.type = "450K" and prior.prob = T to adjust for variation in the number of 450K probes mapping to each gene.

### **Bootstrapped confidence intervals**

All bootstrapped 95% confidence intervals were calculated over 1,000 bootstrap samples.

# Results

## Identification of hypervariable CpGs

We analysed methylation data from 3,474 individuals across 30 datasets (28 Illumina450K and 2 EPIC array) comprising 19 unique tissue/cell types and 8 ethnicities covering a range of ages (Supplementary Tables 1-4). We focussed on CpGs covered by the Illumina450K array and began by excluding probes with poor detection p-values, cross-hybridising probes, probes on the X and Y chromosomes and probes associated with known SNPs (see Methods for details).

We aimed to identify CpGs with consistently high interindividual variation in methylation across diverse datasets, so minimising the effects of dataset-specific drivers of variability. We reasoned that removal of unmeasured technical, batch and cell heterogeneity effects within each dataset would maximise power to detect true biological variation across datasets and therefore adjusted all methylation values for the first ten principal components (PCs) of methylation variation, as well as sex (in datasets with both sexes) and age (where available).

Our strategy for identifying tissue- and ethnicity- independent hypervariable CpGs ('hvCpGs') is summarised in Fig. 1 and detailed in 'Methods'. We first selected CpGs within the top x% of each dataset by methylation Beta variance, and then took the intersection of these CpGs across an increasing proportion of covered datasets, ensuring that each CpG was present in at least 15 datasets (Supplementary Fig. 1A). Using this approach, we identified 4,330 hvCpGs, defined as CpGs with methylation Beta variance in the top 5% of all CpGs in at least 65% of datasets for which that CpG passed QC criteria (Table 1). This definition met our required criteria of selecting CpGs that are highly variable in a large number of tissues and ethnicities (median [IQR] for each hvCpG = 13[10,15] and 7[6,7] respectively; Supplementary Fig. 1B). Note that no CpGs are expected to meet these criteria if the top 5% most variable CpGs in each dataset are entirely independent of those in the others. While re-defining these thresholds will change the set of hvCpGs, we noted that ~80% of identified hvCpGs overlapped with an alternative set obtained when selecting CpGs with methylation Beta variance in the top 20% in at least 90% of datasets (Supplementary Fig. 1C), meaning that the majority of hvCpGs are within the top 20% of variable loci in almost all covered datasets.

We next compared the 4,330 hvCpGs with an alternative set obtained using the same method but without prior adjustment of each dataset for the first ten PCs. This alternative set contained only 1,302 CpGs, which confirmed our intuition that PC adjustment maximises power to identify true dataset-independent hypervariability by removing unwanted technical variation (Supplementary Fig. 8). Finally, we used reported measures of methylation variability among technical replicates<sup>22,42</sup> to remove 187 technically unreliable probes (see Methods), leaving a final set of 4,143 hvCpGs (Table 1; Supplementary Table 5).

hvcpgs are enriched for intermediate methylation values in all datasets compared to the array background (Supplementary Fig. 2; see Table 1 for definition of array background) and are distributed throughout the genome (Supplementary Fig. 5A), with 2,219 (54%) falling within 716 ‘clusters’ containing two or more hvcpgs separated by < 4 kb (Supplementary Fig. 5C). To account for the possibility that our downstream analyses may be biased by these distributional properties, we generated a set of controls that were distribution-matched in a whole blood dataset (Supplementary Fig. 3) and a set of ‘de-clustered hvcpgs’ (Table 1, ‘Methods’).

### **hvcpg variability is not driven by age, sex, or cell heterogeneity**

Evidence from multiple studies suggests that methylation variability can increase with age (termed epigenetic drift)<sup>11,55</sup>, raising the possibility that cross-dataset hypervariability of hvcpgs is driven in part by a large proportion of adult/elderly samples. However, 3,815 (92%) out of 4,122 hvcpgs with methylation measured in cord blood and/or buccal samples from infants showed methylation variance within the top 5% of CpGs in those datasets (Supplementary Table 10), suggesting that high variability at hvcpgs arises in early life. We further probed age stability of hvcpgs by leveraging two studies of age effects in blood. The first study reported methylation consistency in individuals sampled at two time points six years apart<sup>46</sup>. The temporal stability of hvcpgs was significantly greater than that of controls with similar methylation Beta distributions to hvcpgs at the first time point (Wilcoxon paired signed-rank test p-value < 5.7 x10<sup>-81</sup>), with 95% of hvcpgs considered temporally stable versus 89% of controls (Supplementary Fig. 6A). The second measured epigenetic drift in a cross-sectional study of 3,295 whole blood samples from individuals aged 18 to 88<sup>11</sup>. Only 7% of hvcpgs overlapped CpGs that show increased methylation variability with age, compared to 16.5% of blood distribution-matched controls (Supplementary Fig. 6B). This suggests that the majority of hvcpgs are stable over a broad time period in whole blood and further supports the notion that hypervariability of hvcpgs in multiple datasets is not an artefact of epigenetic drift effects.

While methylation values were pre-adjusted for sex in all datasets where sex was available as a covariate (24 out of 30 datasets), we further investigated the potential for sex effects to drive methylation variance by considering the four female-only datasets. 4,102 (99%) of the 4,136 hvcpgs covered in any of these datasets had methylation variance among the highest 5% in at least one (Supplementary Table 10). Furthermore, we found no significant difference in mean methylation at hvcpgs between male and females in a diverse set of tissues (Supplementary Fig. 6C). Finally, 3,548 (96%) of the 3,678 hvcpgs covered in purified CD4+ and CD8+ datasets had methylation variance among the top 5% in at least one dataset (Supplementary Table 10). Together, these data strongly suggest that variability at hvcpgs is not driven by sex, age, or cell heterogeneity effects.



## Hypervariability is not driven by genetic variants

Genetic variation is an important driver of interindividual methylation differences<sup>4,5</sup>. There is evidence that mQTLs can be shared across different tissues<sup>15,16,56,57</sup> and ethnic groups<sup>5</sup>, raising the possibility that ‘universal’ (multi-tissue and multi-ethnic) mQTLs might drive cross-dataset variability at hvCpGs. We therefore investigated the potential influence of methylation quantitative trait loci (mQTL) on methylation variability at hvCpGs by leveraging a recent large meta-GWAS (36 cohorts,  $n = 27,750$  individuals) that identified common genetic variants associated with methylation in blood from Europeans<sup>43</sup>, reasoning that by definition ‘universal’ mQTLs would be included in this meta-analysis.

We considered multiple methylation variance thresholds (5%, 10% and 20%) and observed a positive relationship between hypervariability and both the probability of a significant mQTL association and the mean mQTL effect size (Fig. 2A). Amongst hvCpGs, there were 6,985 *cis* mQTL (covering 3,635 hvCpGs and 6,417 SNPs) and 971 *trans* mQTL (covering 713 hvCpGs and 753 SNPs). Overall, 3,722 (90%) of hvCpGs were reported to have at least one (*cis* or *trans*) mQTL association that were estimated to explain, on average, 4% of methylation variance (Fig. 2B). This suggests that additive genetic effects explain a small to moderate proportion of methylation variability at these hypervariable loci in blood. Noting that the statistical power to detect mQTL associations will be greater at loci that are inherently variable, we matched hvCpGs to CpGs with the same number of mQTL associations and similar average % variance explained by mQTL (‘mQTL-matched controls’, Table 1, Supplementary Fig. 4). hvCpGs showed an average 5-fold increase in methylation variance compared to mQTL-matched controls across datasets (Fig. 2C), further supporting the notion that methylation variation at hvCpGs is not principally driven by universal genetic effects.

To further probe the influence of genetic effects on hvCpG methylation we examined the overlap between hvCpGs and CpGs that show DNAm variation between monozygotic (MZ) co-twins that is *equivalently variable* (*ev*) to that between unrelated individuals, suggestive of genetically independent variable methylation establishment after MZ twin splitting<sup>45</sup>. In total, hvCpGs comprise 122 (42%) of the 317 *ev*CpGs identified in blood (1.9-fold enrichment relative to distribution-matched controls) and 62% of those that were replicated as *ev*CpGs in adipose tissue (2.8-fold enrichment relative to controls) (Supplementary Table 11), supporting the notion that hvCpGs are likely influenced but not determined by genetic variation in multiple tissues.

## hvcpgs show covariation across tissues derived from different germ layers

Variable DNAm states that covary across tissues derived from different germ layers and that are influenced but not determined by genotype may have been established before germ layer separation in early embryonic development<sup>26</sup>. None of the datasets considered here had multi-tissue data from the same individuals. We

therefore examined the overlap between hvCpGs and 3,089 CpGs that show systemic (cross-tissue) interindividual variation (SIV), collated from four published sources<sup>24–27</sup> (Supplementary Table 6). 24% of hvCpGs overlap a known SIV-CpG as do 21% of hvCpGs with a blood distribution-matched control, the latter representing a ~5-fold enrichment (Fig. 3A, Supplementary Table 12, Supplementary Fig. 10A). We note that a further 540 (13%) hvCpGs are within 1 kb of a SIV-CpG, ~5-fold greater than array background CpGs. This suggests that many hvCpGs directly overlap or co-localise with a known SIV-CpG.

The set of all hvCpGs comprises 32.1% of the 3,089 CpGs reported as SIV in any of the four independent studies analysed despite comprising <1% of the 450K array. When considering ‘high-confidence’ SIV-CpGs reported in at least two or three of the four screens, the proportion identified rises to 76.5% and 95.1% respectively (Fig.3B). This suggests that our approach of identifying hypervariable loci across multiple datasets may be a more powerful method for identifying putative SIV loci, compared to existing SIV screens that necessarily rely on rare datasets with multi-tissue, multi-germ layer methylation data from small numbers of individuals. To confirm this, we estimated the power to detect SIV using the multi-tissue data from four individuals analysed by van Baak *et al.*<sup>25</sup>. Using a permutation framework (‘Methods’), we estimated the mean power to detect SIV as 56% (median [IQR] = 0.58 [0.44, 0.72]; Supplementary Fig. 7). As expected, given the small sample size of this multi-tissue dataset, a large proportion of hvCpGs (75%) did not meet the minimum interindividual variation threshold of 0.2 used by van Baak *et al.* to define SIV. On the assumption that hvCpGs are highly enriched for true SIV, this could explain why hvCpGs constitute 61.7% of the van Baak *et al.* SIV-CpGs, while just 13.5% of hvCpGs are identified as SIV-CpGs in the van Baak *et al.* analysis.

To directly test our hypothesis that hvCpGs comprise previously unidentified SIV loci, we analysed a dataset of fetal tissues from 27 individuals, each with methylation data from two tissues derived from different germ layers (see Supplementary Table 7). Inter-germ layer correlations at hvCpGs had a median average Pearson *r* of 0.42, compared to array background CpGs which had a median average Pearson *r* of 0.05 (Fig. 3C left). Of the 3,878 hvCpGs covered in this fetal multi-tissue dataset, 1,653 (42%) had an average inter-germ layer Pearson *r*  $\geq 0.5$ . Of these, 58% did not overlap previously identified SIV loci, suggesting that hvCpGs comprise novel SIV loci. A comparison of the average inter-germ layer correlation at hvCpGs and at previously identified SIV-CpGs showed that hvCpGs and SIV-CpGs had similar inter-germ layer correlations (Fig. 3C right).

### hvCpGs are enriched for loci with distinctive methylation patterns in MZ twins

We further investigated evidence for establishment of hvCpG methylation states in the early embryo by testing the overlap between hvCpGs and 1,217 “epigenetic supersimilarity” (ESS) CpGs overlapping array background. ESS CpGs show high interindividual variation with greater-than-expected methylation

concordance between monozygotic co-twins in adipose tissue, suggestive of methylation establishment in the early zygote before MZ cleavage<sup>25</sup>. 13% of hvCpGs overlap an ESS CpG, showing a ~9.5-fold enrichment for ESS CpGs relative to distribution-matched controls (Fig. 3A, Supplementary Table 12, Supplementary Fig. 10B).

We next examined the overlap between hvCpGs and a set of CpGs showing a unique methylation signature in adult tissues from MZ vs DZ twins ('MZ twinning CpGs', Table 2), implicating these CpGs in MZ twin splitting events in early development<sup>28</sup>. 7% of hvCpGs overlap an MZ twinning CpG, showing a 3.7-fold enrichment for MZ twinning CpGs compared to distribution-matched controls (Fig. 3A, Supplementary Table 12).

Notably, 54% of ESS and 37% of MZ twinning CpGs overlapping array background are hvCpGs (Fig. 3B).

## Reconciling the timing of variable methylation establishment at hvCpGs

The enrichments that we observe for SIV, ESS, evCpGs and MZ twinning CpGs offer a potential insight into the timing of methylation establishment at hvCpGs. 38% of hvCpGs overlap at least one of these CpG sets and enrichment is stronger amongst CpGs that show at least two of these properties (Supplementary Fig. 11). In particular, hvCpGs comprise 78% of SIV-ESS loci and 65% of SIV-MZ twinning loci, suggesting that SIV loci with evidence of establishment in the pre-gastrulation embryo are enriched for hvCpGs<sup>26</sup>.

Variable methylation states identified at evCpGs are thought to originate in embryonic development and/or early post-natal life<sup>43</sup>. We note that 41 out of 317 evCpGs overlap SIV and/or MZ twinning CpGs, suggesting that at least a subset may be established in the pre-gastrulation embryo. hvCpGs comprise 67% of evCpGs that overlap SIV-CpGs, and 76% of that overlap MZ twinning CpGs (Supplementary Fig. 11).

## hvCpGs are enriched for parent-of-origin methylation and proximal TEs

In mice, variable methylation states have been associated with the Intracisternal A Particle (IAP) class of endogenous retrovirus<sup>58,59</sup>, with growing evidence that methylation variability may in part be driven by incomplete silencing of IAPs in early development<sup>60–62</sup>. In humans, SIV-CpGs are enriched for proximal endogenous retrovirus elements (ERVs)<sup>63</sup>, including the subclasses ERV1 and ERVK<sup>26</sup>. This is also the case with hvCpGs, which show a ~1.3-fold and ~1.7-fold enrichment for proximal (within 10 kb) ERV1 and ERVK elements respectively, relative to both array background and blood distribution-matched controls (Fig. 3D, Supplementary Fig. 10C, Supplementary Table 12). Approximately 4.7% of hvCpGs are also located within

1Mb of telomeric regions, showing a 1.8-fold enrichment relative to distribution-matched controls and array background CpGs (Supplementary Table 12).

Maintenance of parent of origin-specific methylation (PofOm) in the pre-implantation embryo is critical for genomic imprinting<sup>64</sup>, and several previously identified SIV loci have been found to be associated with imprinted genes and/or PofOm<sup>25,63</sup>. 58 hvCpGs (1.4%) were annotated to 32 imprinted genes (Supplementary Table 13), no more than expected by chance since 1.9% of array background CpGs are annotated to imprinted genes. 10 hvCpGs were annotated to the polymorphically imprinted non-coding RNA *VTRNA2-1*, a well-established SIV locus that is associated with periconceptional environmental exposures<sup>25,63,65–67</sup>. Although only a small proportion (2.2%) of hvCpGs overlap regions of PofOm identified in peripheral blood<sup>68</sup>, this overlap represents a 3.5-fold and 11-fold enrichment relative to distribution-matched controls and array background respectively that is maintained after de-clustering (Fig. 3A, Supplementary Fig. 10D, Supplementary Table 12). This overlap constitutes 13% of all PofOm CpGs overlapping array background (Fig 3B).

### hvCpGs show sensitivity to pre-natal environment

Variable methylation states established in early development that are sensitive to environmental perturbation are promising candidates for exploring the developmental origins of health and disease<sup>69–71</sup>. We explored whether hvCpGs show sensitivity to pre-natal environment by examining their overlap with loci associated with season of conception ('SoC') in a rural Gambian population exposed to seasonal fluctuations in diet and other factors<sup>72–74</sup>. hvCpGs comprise 70 (29%) out of 242 previously identified SoC-CpGs<sup>29</sup> overlapping array background, an approximately 3-fold enrichment relative to distribution-matched controls (Supplementary Table 11).

We next leveraged a recent meta-analysis of 2,365 cord blood samples that modelled genetic (G), genetic by environment (GxE) and additive genetic and environment (G+E) effects at variably methylated probes, where E represents a range of prenatal exposures including pre-pregnancy BMI, maternal smoking, gestational age, hypertension, anxiety and depression<sup>14</sup>. Of the 703 hvCpGs overlapping the neonatal blood variably methylated regions explored in that study, G, GxE, and G+E effects were the 'winning' models for 30%, 30% and 40% of probes respectively, representing an increase in G+E effects compared to array background (Supplementary Fig. 12). This analysis supports our intuition that hvCpGs are influenced but not determined by genetic variation, with pre-natal environment as an additional influencing factor.

### Chromatin states at hvCpGs

Compared to array background, hvCpGs are enriched within intergenic regions and CpG island 'shores' but are depleted within gene bodies and regions directly upstream of transcription start sites (Supplementary

Fig. 9). We predicted chromatin states at hvCpGs by examining the overlaps of hvCpGs with histone modifications using the chromHMM 15-state model<sup>52</sup> for seven tissues including embryonic stem cells (H1 ESCs), and fetal and adult tissues<sup>53</sup>. Although many hvCpGs were associated with regulatory elements in all tissues, hvCpGs were generally depleted in these regions compared to array background, except within predicted enhancers in H1 ESCs (Supplementary Fig. 13).

### Association with the clustered protocadherin gene locus on chromosome 5

Gene ontology enrichment analysis revealed that hvCpGs were significantly enriched for terms associated with cell-cell adhesion (Fig. 4A), which is largely driven by the colocalization of 3.3% of hvCpGs to clustered protocadherin (*cPCDH*) genes on chromosome 5. This region comprises three clusters of protocadherin genes (*cPCDHα*, *cPCDHβ*, *cPCDHγ*), each containing many variable exons whose promoter choice is determined stochastically via differential methylation by DNA-methyltransferase 3 beta (DNMT3B) in early embryonic development<sup>75,76</sup>, resulting in the expression of distinct *cPCDH* isoforms of cell-surface proteins that are critical for establishing neuronal circuits<sup>77</sup>. The *cPCDH* gene locus has also been found to be influenced by age<sup>11,78–80</sup>. Accordingly, although a minority (5%) of hvCpGs showed evidence of epigenetic drift in blood<sup>11</sup>, these are enriched within the *cPCDH*\_locus relative to those that did not show evidence of epigenetic drift (Fisher's Exact Test (FET) p-value =  $9.4 \times 10^{-9}$ , OR = 4.02). Hypervariable methylation states at the *cPCDH* gene locus may therefore be driven by early developmental and/or aging effects. Noting that evCpGs and MZ twinning CpGs (Table 2) have also been reported to colocalise with this locus<sup>45,81</sup>, hvCpGs annotated to *cPCDH* genes were ~8.5-fold enriched for MZ twinning CpGs (FET p-value =  $1.04 \times 10^{-22}$ ) and ~3-fold enriched for evCpGs (FET p-value =  $1.6 \times 10^{-3}$ ) relative to hvCpGs that were not.

### Association of hvCpGs with reported EWAS trait associations

To probe the potential functional role of hvCpGs, we analysed their overlap with traits reported in the epigenome-wide association studies (EWAS) catalogue (<http://ewascatalog.org/>). 86% of hvCpGs show significant associations (reported p-value <  $1 \times 10^{-4}$ ) with one or more of 231 unique traits covered in the catalogue (Supplementary Table 9). However, compared to blood distribution-matched controls, a suitable comparator given that the majority of EWAS have been carried out in blood, we found that hvCpGs were enriched amongst CpGs associated with sex and Alzheimer's disease only (Fig. 4B, left panel).

Noting that sex-associated hvCpGs are not influenced by sex differences in the datasets that we analysed (Supplementary Fig. 6C, bottom) and that a similar proportion of SIV-CpGs are also associated with sex (23% of hvCpGs and 20% of the 3,089 SIV-CpGs considered in our study), we speculate that the association with sex may be a feature of variable methylation states established in early development. Amongst the 64 hvCpGs

associated with Alzheimer's disease, 23 overlap previously identified SIV and/or ESS loci, 9 of which annotated to *CYP2E1*, a gene that has also been associated with Parkinson's disease and rheumatoid arthritis<sup>82,83</sup>.

hvcpgs were notably depleted amongst age-related traits relative to distribution-matched controls (Fig. 4B, right panel), in agreement with our earlier findings that hvcpgs are largely stable with age (Supplementary Fig. 6). hvcpgs are also depleted amongst CpGs that are differentially methylated between buccal cells and peripheral blood mononucleocytes ('Tissue' in Fig. 4B), supporting the notion that hvcpgs may be established before cell differentiation and that the method used to identify the hvcpgs is robust to tissue-specific methylation variation.

## Discussion

We have identified and characterised tissue- and ethnicity- independent hypervariable methylation states at CpGs covered by the 450k array. Our methodological approach was designed to be robust to dataset-specific drivers of methylation variability, including sex, age, cell type heterogeneity and technical artefacts. We identified 4,143 hvCpGs and found strong evidence that methylation states at many hvCpGs are likely to be established in the early embryo and are stable postnatally. Our analysis positions hvCpGs as tissue- and ethnicity- independent age-stable biomarkers of early stochastic and/or environmental effects on DNA methylation.

hvCpGs cover ~1% of the 450K array and were in the top 5% variable methylation states in an average of 13 distinct tissues and 7 ethnicities. Our study is not the first to investigate DNAm patterns in multiple tissues. Previous studies have identified CpGs that are differentially methylated between tissues<sup>85–87</sup>; determined the extent to which variable methylation states in accessible tissues (such as blood) reflect those in inaccessible tissues such as brain<sup>56,86–90</sup>; compared methylation patterns between peripheral tissues<sup>57,91,92</sup>; directly identified SIV loci using tissues derived from different germ layers<sup>24–27,63,93</sup>; functionally characterised tissue-specific variably methylated regions<sup>94</sup>; and examined the extent to which common drivers of methylation variation, such as genetics, age, sex and environment, are tissue-specific<sup>8,12,15–18,95,96</sup>. The majority of these studies used a comparatively small number of tissues or cell-types, and few have used multi-tissue datasets from different ethnicities<sup>15</sup>. To our knowledge, ours is the first study to explore the extent to which variably methylated CpGs are shared across diverse tissues and ethnicities in the human genome.

The majority of hvCpGs were associated with at least one mQTL suggesting that additive genetic effects influence methylation variation at these loci. However, a comparison with mQTL-matched controls, together with evidence of enrichment for sensitivity to periconceptional environment and methylation discordance between MZ twins, suggests that stochastic and/or environmental effects have a relatively large influence on methylation variability at hvCpGs. In line with this, a large proportion of hvCpGs show evidence of systemic interindividual variation (SIV), that is, intra-individual correlation in methylation across tissues derived from different germ layers. Whilst loci that covary across different tissue types are enriched for mQTL effects<sup>16,56,57,91</sup>, it has been suggested that SIV loci are putative human metastable epialleles with variable methylation states established before gastrulation that are influenced but not determined by genetic variation<sup>26</sup>.

Our fetal multi-tissue analysis supports the notion that SIV at hvCpGs arises during development and is not, for example, driven by post-natal environmental influences that act across many tissues. hvCpGs were also highly enriched for epigenetic supersimilarity loci and MZ twinning-associated CpGs, both of which have been



linked to establishment of methylation in the cleavage stage pre-implantation embryo<sup>25,28</sup>. The degree of overlap between variably methylated regions in different cell types has also been linked to their common developmental origin<sup>94</sup>. If this pattern holds true, it follows that stochastic and/or environmentally influenced variably methylated loci that are shared across a large number of diverse tissues are likely to have originated before germ-layer differentiation. We note that it is possible that variable methylation variation at some hvCpGs is influenced by later gestational or post-natal environmental effects, acting in addition to or independently of early environmental effects across multiple tissues, as has been suggested at the *VTRNA2-1* locus in the context of folate supplementation in pregnancy<sup>97</sup>, maternal age at delivery<sup>67</sup>, and smoking<sup>98</sup>.

The association of hvCpGs with parent-of-origin-specific methylation and proximal ERV1 and ERVK elements is notable because these features have been linked to SIV-CpGs<sup>26</sup>. This suggests that genomic regions targeted by epigenetic silencing or maintenance mechanisms during early embryonic reprogramming may be enriched for stochastic and/or environmentally influenced methylation variation. For example, it has been suggested that regions of PofOm may be vulnerable to stochastic or environmentally-sensitive loss of methylation on the usually-methylated allele or gain of methylation on the usually-unmethylated allele at a later time-point, leading to interindividual methylation variation<sup>29,64,99</sup>. Similarly, certain IAP elements (a class of ERVK LTR retrotransposon) show methylation variation between isogenic mice<sup>58,59</sup> that in several cases can be influenced by pre-natal environment<sup>100–103</sup>. Whilst transposable elements are usually silenced to prevent insertion events from damaging the genome, recent evidence suggests that methylation variability at IAP elements is partly driven by low-affinity binding of *trans*-acting Krüppel-associated box (KRAB)-containing zinc finger proteins (KZFPs)<sup>60</sup> and by sequence variation in KZFP-binding sites<sup>60,104</sup>. Whilst KZFPs are known to target TEs in humans<sup>105,106</sup>, the extent of their role in driving methylation variation is an ongoing area of research.

The large overlap between hvCpGs and ‘high confidence’ SIV-CpGs identified in at least two independent screens suggests that the identification of hvCpGs might constitute a high-powered method for detecting novel SIV loci. Supporting this, the largest SIV screen to date with 10 individuals was reported to be underpowered to detect the well-established SIV locus at the non-coding RNA gene *VTRNA2-1*<sup>27</sup> (represented by 10 hvCpGs), and we found that a 4-individual multi-tissue dataset analysed by van Baak *et al.*<sup>25</sup> had limited power to detect SIV loci. Another consideration is that SIV screens to date have used different sets of tissues. Since loci that covary between one pair of tissues do not necessarily covary between another pair<sup>56</sup>, the enrichment for high confidence SIV loci might reflect the fact that methylation states at hvCpGs covary across a large number of tissues. Importantly, our analysis of a fetal multi-tissue dataset offers a strong validation of previously unreported SIV at hvCpGs.

Our analysis of EWAS trait associations revealed a moderate enrichment for hvCpGs amongst CpGs associated with Alzheimer’s disease and SIV loci have been linked to this and other disease outcomes

including autism, cancer and obesity<sup>26,63,93</sup>. For example, 10 hvCpGs overlap the *PAX8* gene which is a known SIV locus. *PAX8* methylation measured in peripheral blood of Gambian 2-year olds was recently shown to be correlated with thyroid volume and hormone levels in the same children in mid-childhood, and the latter was associated with changes in body fat and bone mineral density<sup>107</sup>. This suggests that hvCpGs are interesting candidates for exploring how stochastic and/or environmentally influenced DNAm states established in early development might influence life-long health.

hvCpGs are variable in diverse ethnicities, raising the possibility that regions of hypervariable methylation may be a conserved feature in the human genome. Stochastic methylation patterns established in the early embryo that are sensitive to early environment and that are able to influence gene expression might mediate a *predictive-adaptive-response* mechanism that senses the pre-natal environment in order to prime the developing embryo to its post-natal environment<sup>70,71</sup>. This would require environmentally responsive variable methylation states to be genetically hardwired into the human genome, providing a means of rapid adaptation to changing local environments on a scale much faster than is attainable through Darwinian evolution<sup>108</sup>. Associations between genotype and methylation variance have been previously reported, for example at the putative metastable epiallele *PAX8*<sup>107</sup>, at the master regulator of genomic imprinting *ZFP57*<sup>25</sup> and at several probes in the major histocompatibility complex (MHC) region associated with rheumatoid arthritis<sup>109</sup>. Interestingly, 4% of hvCpGs are located within the MHC, representing an enrichment relative to the array background (FET p-value =  $2.7 \times 10^{-10}$ , OR = 1.7). Further analysis of genotype-methylation variance effects is required to determine if this region, which contains a large amount of sequence variation and is implicated in many immune-mediated diseases<sup>110</sup>, might contain other examples of genetically-driven phenotypic plasticity that is mediated by DNA methylation.

Whilst our method of adjusting for the first 10 PCs of variation may not have controlled for all technical artifacts within each dataset, if technical issues were to cause a random CpG to be in the top 5% of variance in one dataset, this CpG would be unlikely to be in the top 5% of variance across a majority of datasets. The consequence would therefore be a loss of power to identify hvCpGs rather than the identification of spurious hypervariability. This is supported by our sensitivity analysis with unadjusted methylation data (Supplementary Fig. 8). An exception would be if the probe were consistently unreliable. We tested this using reliability metrics derived from analysis of technical replicates and found no evidence that hvCpGs are driven by technically unreliable probes. However, we note that better adjustment for technical artifacts within datasets<sup>111</sup> and the addition of further datasets would likely lead to the identification of more hvCpGs.

The vast majority of publicly available methylation datasets use the Illumina 450k array. Therefore, a major limitation of this study is that we were only able to analyse the small proportion of the methylome covered by this array, which has been found to miss a disproportionate amount of variable CpGs<sup>27</sup>. However, we note

that our method for identifying hypervariable CpGs can easily be applied to whole methylome sequencing data which is becoming increasingly available.

Through the joint analysis of methylation data from multiple tissues, we have identified a large set of hypervariable loci on the 450K array that are present across multiple tissues and ethnicities. Comparisons with a diverse range of data sources reveal that stochastic and/or environmentally-responsive methylation states at these loci are likely to have been established in the early embryo and appear to be stable with age, making them interesting candidates for studying the developmental origins of life-long health and disease.

## Data availability

The large majority of datasets analysed in this paper are in the public domain. GEO accession numbers and/or further details are provided in Supplementary Tables. A small number of analysed datasets have restricted access. Requests to access these should be submitted to the corresponding authors in the first instance with researcher access requiring an application to the relevant institutional review boards.

## Acknowledgements

GoDMC meta-GWAS summary statistics were kindly provided by Jordana Bell. This data has now been published (Min et al, 2021).

## Funding

This work was supported by the Biotechnology and Biological Sciences Research council (BB/M009513/1); the UK Medical Research Council (MR/M01424X/1, MC\_PC\_MR/RO20183/1 and MR/N006208/1); the French National Research Agency (ANR OCEOADAPTO and ANR PAPUADEVOL) and the Department of Biotechnology, Ministry of Science and Technology, India (BT/IN/DBT-MRC/DFID/24/GRC/2015–16). The fetal tissues supplied are part of the UCL-Imperial Baby Bio Bank collection. <https://directory.biobankinguk.org>

## *Conflict of Interest Disclosure*

None declared.

## References

1. Smith ZD, Meissner A. DNA methylation: roles in mammalian development. *Nature Reviews Genetics*. 2013;14(3):204-220. doi:10.1038/nrg3354
2. Rakyan VK, Chong S, Champ ME, et al. Transgenerational inheritance of epigenetic states at the murine Axin(Fu) allele occurs after maternal and paternal transmission. *Proceedings of the National Academy of Sciences of the United States of America*. 2003;100(5):2538-2543. doi:10.1073/pnas.0436776100
3. Lappalainen T, Greally JM. Associating cellular epigenetic models with human phenotypes. *Nature Reviews Genetics*. 2017;18(7):441-451. doi:10.1038/nrg.2017.32
4. Gaunt TR, Shihab HA, Hemani G, et al. Systematic identification of genetic influences on methylation across the human life course. *Genome Biology*. 2016;17(1):61. doi:10.1186/s13059-016-0926-z
5. Bell JT, Pai AA, Pickrell JK, et al. DNA methylation patterns associate with genetic and gene expression variation in HapMap cell lines. *Genome Biology*. 2011;12(1):R10-R10. doi:10.1186/gb-2011-12-1-r10
6. Jaffe AE, Irizarry RA. Accounting for cellular heterogeneity is critical in epigenome-wide association studies. *Genome Biology*. 2014;15(2):R31. doi:10.1186/gb-2014-15-2-r31
7. Houseman EA, Kile ML, Christiani DC, Ince TA, Kelsey KT, Marsit CJ. Reference-free deconvolution of DNA methylation data and mediation by cell composition effects. *BMC Bioinformatics*. 2016;17(1):259. doi:10.1186/s12859-016-1140-4
8. Singmann P, Shem-Tov D, Wahl S, et al. Characterization of whole-genome autosomal differences of DNA methylation between men and women. *Epigenetics & chromatin*. 2015;8:43. doi:10.1186/s13072-015-0035-3
9. Yousefi P, Huen K, Davé V, Barcellos L, Eskenazi B, Holland N. Sex differences in DNA methylation assessed by 450 K BeadChip in newborns. *BMC genomics*. 2015;16:911. doi:10.1186/s12864-015-2034-y
10. Horvath S. DNA methylation age of human tissues and cell types. *Genome Biology*. 2013;14(10):R115-R115. doi:10.1186/gb-2013-14-10-r115
11. Slieker RC, van Iterson M, Luijk R, et al. Age-related accrual of methylomic variability is linked to fundamental ageing mechanisms. *Genome Biology*. Published online 2016. doi:10.1186/s13059-016-1053-6
12. Busche S, Shao X, Caron M, et al. Population whole-genome bisulfite sequencing across two tissues highlights the environment as the principal source of human methylome variation The Multiple Tissue Human Expression Resource. Published online 2015. doi:10.1186/s13059-015-0856-1
13. Hannon E, Knox O, Sugden K, et al. Characterizing genetic and environmental influences on variable DNA methylation using monozygotic and dizygotic twins. Greally JM, ed. *PLOS Genetics*. 2018;14(8):e1007544-e1007544. doi:10.1371/journal.pgen.1007544
14. Czamara D, Eraslan G, Page CM, et al. Integrated analysis of environmental and genetic influences on cord blood DNA methylation in new-borns. *Nature Communications*. 2019;10(1). doi:10.1038/s41467-019-10461-0
15. Smith AK, Kilaru V, Kocak M, et al. Methylation quantitative trait loci (meQTLs) are consistently detected across ancestry, developmental stage, and tissue type. *BMC Genomics*. 2014;15(1):145. doi:10.1186/1471-2164-15-145
16. Lin D, Chen J, Perrone-Bizzozero N, et al. Characterization of cross-tissue genetic-epigenetic effects and their patterns in schizophrenia. *Genome Medicine*. 2018;10(1):13. doi:10.1186/s13073-018-0519-4
17. Day K, Waite LL, Thalacker-Mercer A, et al. Differential DNA methylation with age displays both common and dynamic features across human tissues that are influenced by CpG landscape. *Genome Biology*. 2013;14(9). doi:10.1186/gb-2013-14-9-r102
18. Slieker RC, Relton CL, Gaunt TR, Slagboom PE, Heijmans BT. Age-related DNA methylation changes are tissue-specific with ELOVL2 promoter methylation as exception. *Epigenetics and Chromatin*. 2018;11(1). doi:10.1186/s13072-018-0191-3
19. Price EM, Robinson WP. Adjusting for Batch Effects in DNA Methylation Microarray Data, a Lesson Learned. *Frontiers in Genetics*. 2018;0(MAR):83. doi:10.3389/FGENE.2018.00083
20. Leek JT, Scharpf RB, Bravo HC, et al. Tackling the widespread and critical impact of batch effects in high-throughput data. *Nature Reviews Genetics*. 2010;11(10):733-739. doi:10.1038/nrg2825
21. Harper KN, Peters BA, Gamble M v. Batch effects and pathway analysis: Two potential perils in cancer studies involving DNA methylation array analysis. *Cancer Epidemiology Biomarkers and Prevention*. 2013;22(6):1052-1060. doi:10.1158/1055-9965.EPI-13-0114

22. Sugden K, Hannon EJ, Arseneault L, et al. Patterns of Reliability: Assessing the Reproducibility and Integrity of DNA Methylation Measurement. *Patterns*. 2020;1(2). doi:10.1016/j.patter.2020.100014
23. Bibikova M, Barnes B, Tsan C, et al. High density DNA methylation array with single CpG site resolution. *Genomics*. 2011;98(4):288-295. doi:10.1016/J.YGENO.2011.07.007
24. Harris RA, Nagy-Szakal D, Kellermayer R. Human metastable epiallele candidates link to common disorders. *Epigenetics*. 2013;8(2):157-163. doi:10.4161/epi.23438
25. van Baak TE, Coarfa C, Dugué PA, et al. Epigenetic supersimilarity of monozygotic twin pairs. *Genome Biology*. 2018;19(1):2. doi:10.1186/s13059-017-1374-0
26. Kessler NJ, Waterland RA, Prentice AM, Silver MJ. Establishment of environmentally sensitive DNA methylation states in the very early human embryo. *Science Advances*. Published online 2018. doi:10.1126/sciadv.aat2624
27. Gunasekara CJ, Scott CA, Laritsky E, et al. A genomic atlas of systemic interindividual epigenetic variation in humans. *Genome Biology*. 2019;20(1):105. doi:10.1186/s13059-019-1708-1
28. van Dongen J, Gordon SD, McRae AF, et al. Identical twins carry a persistent epigenetic signature of early genome programming. *Nature Communications*. 2021;12(1):5618. doi:10.1038/s41467-021-25583-7
29. Silver MJ, Saffari A, Kessler NJ, et al. Environmentally sensitive hotspots in the methylome of the early human embryo. *bioRxiv*. Published online January 1, 2021:777508. doi:10.1101/777508
30. Barrett T, Wilhite SE, Ledoux P, et al. NCBI GEO: archive for functional genomics data sets—update. *Nucleic Acids Research*. 2013;41(D1):D991-D995. doi:10.1093/nar/gks1193
31. Silva TC, Colaprico A, Olsen C, et al. TCGA Workflow: Analyze cancer genomics and epigenomics data using Bioconductor packages. *F1000Research*. 2016;5:1542. doi:10.12688/f1000research.8923.1
32. Mounir M, Lucchetta M, Silva TC, et al. New functionalities in the TCGAbiolinks package for the study and integration of cancer data from GDC and GTEx. *PLoS computational biology*. 2019;15(3):e1006701-e1006701. doi:10.1371/journal.pcbi.1006701
33. Colaprico A, Silva TC, Olsen C, et al. TCGAbiolinks: an R/Bioconductor package for integrative analysis of TCGA data. *Nucleic Acids Research*. 2016;44(8):e71-e71. doi:10.1093/nar/gkv1507
34. Davis S, Meltzer PS. GEOquery: a bridge between the Gene Expression Omnibus (GEO) and BioConductor. *Bioinformatics*. 2007;23(14):1846-1847. doi:10.1093/bioinformatics/btm254
35. 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*. 2016;8(3):389-399.
36. Chandak GR, Silver MJ, Saffari A, et al. Protocol for the EMPHASIS study; epigenetic mechanisms linking maternal pre-conceptional nutrition and children's health in India and Sub-Saharan Africa. *BMC Nutrition*. 2017;3(1):81. doi:10.1186/s40795-017-0200-0
37. Brucato N, Fernandes V, Mazières S, et al. The Comoros Show the Earliest Austronesian Gene Flow into the Swahili Corridor. *American journal of human genetics*. 2018;102(1):58-68. doi:10.1016/j.ajhg.2017.11.011
38. Morris TJ, Butcher LM, Feber A, et al. ChAMP: 450k Chip Analysis Methylation Pipeline. *Bioinformatics*. Published online 2014. doi:10.1093/bioinformatics/btt684
39. Nordlund J, Bäcklin CL, Wahlberg P, et al. Genome-wide signatures of differential DNA methylation in pediatric acute lymphoblastic leukemia. *Genome Biology*. 2013;14(9):r105. doi:10.1186/gb-2013-14-9-r105
40. Zhou W, Laird PW, Shen H. Comprehensive characterization, annotation and innovative use of Infinium DNA methylation BeadChip probes. *Nucleic Acids Research*. 2016;45(4):gkw967. doi:10.1093/nar/gkw967
41. Teschendorff AE, Marabita F, Lechner M, et al. A beta-mixture quantile normalization method for correcting probe design bias in Illumina Infinium 450 k DNA methylation data. *Bioinformatics*. 2013;29(2):189-196. doi:10.1093/bioinformatics/bts680
42. Bose M, Wu C, Pankow JS, et al. Evaluation of microarray-based DNA methylation measurement using technical replicates: The atherosclerosis risk in communities (ARIC) study. *BMC Bioinformatics*. 2014;15(1):312. doi:10.1186/1471-2105-15-312
43. Min JL, Hemani G, Hannon E, et al. Genomic and phenotypic insights from an atlas of genetic effects on DNA methylation. *Nature genetics*. 2021;53(9):1311-1321.
44. Bonilla C, Bertoni B, Min JL, Hemani G, Consortium G of DNAM, Elliott HR. Investigating DNA methylation as a potential mediator between pigmentation genes, pigmentary traits and skin cancer. *Pigment Cell & Melanoma Research*. 2020;n/a(n/a). doi:https://doi.org/10.1111/pcmr.12948



45. Planterose Jiménez B, Liu F, Caliebe A, et al. Equivalent DNA methylation variation between monozygotic co-twins and unrelated individuals reveals universal epigenetic inter-individual dissimilarity. *Genome Biology*. 2021;22(1):18. doi:10.1186/s13059-020-02223-9
46. Flanagan JM, Brook MN, Orr N, et al. Temporal stability and determinants of white blood cell DNA methylation in the breakthrough generations study. *Cancer Epidemiology Biomarkers and Prevention*. 2015;24(1):221-229. doi:10.1158/1055-9965.EPI-14-0767
47. Grundberg E, Meduri E, Sandling JK, et al. Global analysis of DNA methylation variation in adipose tissue from twins reveals links to disease-associated variants in distal regulatory elements. *American journal of human genetics*. 2013;93(5):876-890. doi:10.1016/j.ajhg.2013.10.004
48. van Dongen J, Gordon SD, McRae AF, et al. Identical twins carry a persistent epigenetic signature of early genome programming. *Nature Communications*. 2021;12(1):5618. doi:10.1038/s41467-021-25583-7
49. Zink F, Magnusdottir DN, Magnusson OT, et al. Insights into imprinting from parent-of-origin phased methylomes and transcriptomes. *Nature Genetics*. 2018;50(11):1542-1552. doi:10.1038/s41588-018-0232-7
50. Loh K, Modhukur V, Rajashekar B, et al. DNA methylome profiling of human tissues identifies global and tissue-specific methylation patterns. *Genome biology*. Published online 2014. doi:10.1186/gb-2014-15-4-r54
51. Min JL, Hemani G, Smith GD, Relton C, Suderman M. Meffil: efficient normalization and analysis of very large DNA methylation datasets. *Bioinformatics*. 2018;34(23):3983. doi:10.1093/BIOINFORMATICS/BTY476
52. Ernst J, Kellis M. ChromHMM: Automating chromatin-state discovery and characterization. *Nature Methods*. 2012;9(3):215-216. doi:10.1038/nmeth.1906
53. Roadmap Epigenomics Consortium, Kundaje A, Meuleman W, et al. Integrative analysis of 111 reference human epigenomes. *Nature*. Published online 2015. doi:10.1038/nature14248
54. Phipson B, Maksimovic J, Oshlack A. missMethyl: an R package for analyzing data from Illumina's HumanMethylation450 platform. *Bioinformatics*. 2016;32(2):286-288. doi:10.1093/bioinformatics/btv560
55. Fraga MF, Ballestar E, Paz MF, et al. *Epigenetic Differences Arise during the Lifetime of Monozygotic Twins*; 2005. <http://www.pnas.org/cgi/doi/10.1073/pnas.0500398102>
56. Hannon E, Lunnon K, Schalkwyk L, Mill J. Interindividual methylomic variation across blood, cortex, and cerebellum: Implications for epigenetic studies of neurological and neuropsychiatric phenotypes. *Epigenetics*. 2015;10(11):1024-1032. doi:10.1080/15592294.2015.1100786
57. Islam SA, Goodman SJ, MacIsaac JL, et al. Integration of DNA methylation patterns and genetic variation in human pediatric tissues help inform EWAS design and interpretation 06 Biological Sciences 0604 Genetics. *Epigenetics and Chromatin*. 2019;12(1):1. doi:10.1186/s13072-018-0245-6
58. Rakyan VK, Blewitt ME, Druker R, Preis JJ, Whitelaw E. Metastable epialleles in mammals. *Trends in Genetics*. 2002;18(7):348-351. doi:10.1016/S0168-9525(02)02709-9
59. Kazachenka A, Bertozzi TM, Sjöberg-Herrera MK, et al. Identification, Characterization, and Heritability of Murine Metastable Epialleles: Implications for Non-genetic Inheritance. *Cell*. 2018;175(5):1259-1271.e13. doi:10.1016/j.CELL.2018.09.043
60. Bertozzi TM, Elmer JL, Macfarlan TS, Ferguson-Smith AC. KRAB zinc finger protein diversification drives mammalian interindividual methylation variability. *Proceedings of the National Academy of Sciences*. Published online 2020:202017053. doi:10.1073/pnas.2017053117
61. Elmer JL, Hay AD, Kessler NJ, Bertozzi TM, Ainscough EAC, Ferguson-Smith AC. Genomic properties of variably methylated retrotransposons in mouse. *Mobile DNA*. 2021;12(1). doi:10.1186/s13100-021-00235-1
62. Costello KR, Leung A, Trac C, et al. Mechanisms regulating interindividual epigenetic variability at transposable elements. doi:10.1101/2021.06.01.446659
63. Silver MJ, Kessler NJ, Hennig BJ, et al. Independent genomewide screens identify the tumor suppressor VTRNA2-1 as a human epiallele responsive to periconceptional environment. *Genome Biology*. Published online 2015. doi:10.1186/s13059-015-0660-y
64. Monk D, Mackay DJG, Eggermann T, Maher ER, Riccio A. Genomic imprinting disorders: lessons on how genome, epigenome and environment interact. *Nature Reviews Genetics*. 2019;20(4):235-248. doi:10.1038/s41576-018-0092-0
65. Finer S, Iqbal MS, Lowe R, et al. Is famine exposure during developmental life in rural Bangladesh associated with a metabolic and epigenetic signature in young adulthood? A historical cohort study. Published online 2016. doi:10.1136/bmjopen-2016

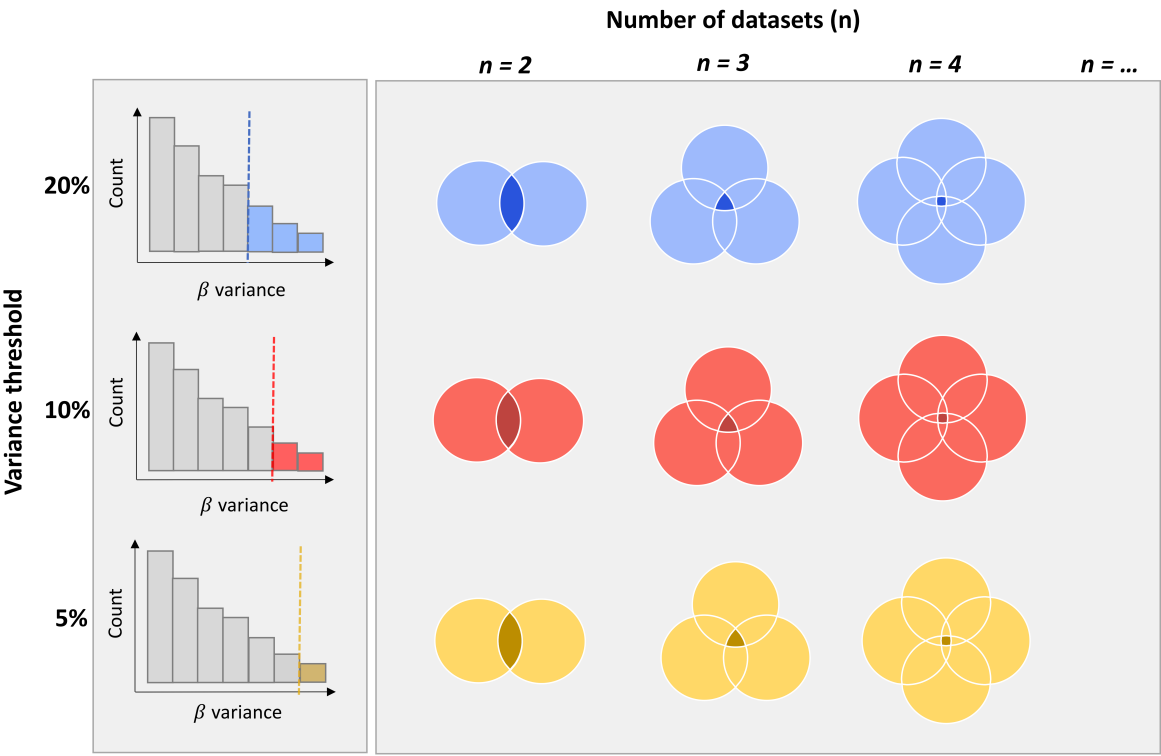


66. Carpenter BL, Remba TK, Thomas SL, et al. Oocyte age and preconceptional alcohol use are highly correlated with epigenetic imprinting of a noncoding RNA (nc886). *Proceedings of the National Academy of Sciences*. 2021;118(12). doi:10.1073/PNAS.2026580118
67. Markunas CA, Wilcox AJ, Xu Z, et al. Maternal age at delivery is associated with an epigenetic signature in both newborns and adults. *PLoS ONE*. 2016;11(7). doi:10.1371/journal.pone.0156361
68. Zink F, Magnusdottir DN, Magnusson OT, et al. Insights into imprinting from parent-of-origin phased methylomes and transcriptomes. *Nature Genetics*. 2018;50(11):1542-1552. doi:10.1038/s41588-018-0232-7
69. Gluckman PD, Hanson MA, Low FM. The role of developmental plasticity and epigenetics in human health. *Birth Defects Research Part C - Embryo Today: Reviews*. 2011;93(1):12-18. doi:10.1002/bdrc.20198
70. Low FM, Gluckman PD, Hanson MA. Developmental Plasticity, Epigenetics and Human Health. *Evolutionary Biology*. Published online 2012. doi:10.1007/s11692-011-9157-0
71. Fleming TP, Watkins AJ, Velazquez MA, et al. Origins of lifetime health around the time of conception: causes and consequences. *The Lancet*. 2018;391(10132):1842-1852. doi:10.1016/S0140-6736(18)30312-X
72. Moore SE, Cole TJ, Collinson AC, Poskitt EME, McGregor IA, Prentice AM. Prenatal or early postnatal events predict infectious deaths in young adulthood in rural Africa. *International Journal of Epidemiology*. Published online 1999. doi:10.1093/ije/28.6.1088
73. Dominguez-Salas P, Moore SE, Cole D, et al. DNA methylation potential: Dietary intake and blood concentrations of one-carbon metabolites and cofactors in rural African women. *American Journal of Clinical Nutrition*. 2013;97(6):1217-1227. doi:10.3945/ajcn.112.048462
74. James PT, Dominguez-Salas P, Hennig BJ, Moore SE, Prentice AM, Silver MJ. Maternal One-Carbon Metabolism and Infant DNA Methylation between Contrasting Seasonal Environments: A Case Study from The Gambia. *Current Developments in Nutrition*. 2019;3(1). doi:10.1093/cdn/nzy082
75. el Hajj N, Dittrich M, Haaf T. Epigenetic dysregulation of protocadherins in human disease. *Seminars in Cell and Developmental Biology*. 2017;69:172-182. doi:10.1016/j.semcdb.2017.07.007
76. Toyoda S, Kawaguchi M, Kobayashi T, et al. Developmental epigenetic modification regulates stochastic expression of clustered Protocadherin genes, generating single neuron diversity. *Neuron*. 2014;82(1):94-108. doi:10.1016/j.neuron.2014.02.005
77. Flaherty E, Maniatis T. The role of clustered protocadherins in neurodevelopment and neuropsychiatric diseases. *Current Opinion in Genetics and Development*. 2020;65:144-150. doi:10.1016/j.gde.2020.05.041
78. Salpea P, Russanova VR, Hirai TH, et al. Postnatal development- and age-related changes in DNA-methylation patterns in the human genome. *Nucleic Acids Research*. 2012;40(14):6477-6494. doi:10.1093/NAR/GKS312
79. McClay JL, Aberg KA, Clark SL, et al. A methylome-wide study of aging using massively parallel sequencing of the methyl-CpG-enriched genomic fraction from blood in over 700 subjects. *Human Molecular Genetics*. 2014;23(5):1175-1185. doi:10.1093/HMG/DDT511
80. Kim S, Wyckoff J, Morris AT, et al. DNA methylation associated with healthy aging of elderly twins. *GeroScience* 2018 40:5. 2018;40(5):469-484. doi:10.1007/S11357-018-0040-0
81. van Dongen J, Nivard MG, Willemsen G, et al. Genetic and environmental influences interact with age and sex in shaping the human methylome. *Nature Communications*. 2016;7. doi:10.1038/ncomms11115
82. Mok A, Rhead B, Holingue C, et al. Hypomethylation of CYP2E1 and DUSP22 Promoters Associated With Disease Activity and Erosive Disease Among Rheumatoid Arthritis Patients. *Arthritis and Rheumatology*. 2018;70(4):528-536. doi:10.1002/art.40408
83. Kaut O, Schmitt I, Wüllner U. Genome-scale methylation analysis of Parkinson's disease patients' brains reveals DNA hypomethylation and increased mRNA expression of cytochrome P450 2E1. *Neurogenetics*. 2012;13(1):87-91. doi:10.1007/S10048-011-0308-3
84. Gunasekara CJ, Waterland RA. A new era for epigenetic epidemiology. Published online 2019.
85. Byun HM, Siegmund KD, Pan F, et al. Epigenetic profiling of somatic tissues from human autopsy specimens identifies tissue- and individual-specific DNA methylation patterns. *Human Molecular Genetics*. 2009;18(24):4808-4817. doi:10.1093/hmg/ddp445
86. Sliker RC, Bos SD, Goeman JJ, et al. Identification and systematic annotation of tissue-specific differentially methylated regions using the Illumina 450k array. *Epigenetics and Chromatin*. 2013;6(1). doi:10.1186/1756-8935-6-26

87. Davies MN, Volta M, Pidsley R, et al. Functional annotation of the human brain methylome identifies tissue-specific epigenetic variation across brain and blood. *Genome biology*. 2012;13(6):R43. doi:10.1186/gb-2012-13-6-r43
88. Walton E, Hass J, Liu J, et al. Correspondence of DNA methylation between blood and brain tissue and its application to schizophrenia research. *Schizophrenia Bulletin*. 2016;42(2):406-414. doi:10.1093/schbul/sbv074
89. Edgar RD, Jones MJ, Meaney MJ, Turecki G, Kobor MS. BECon: A tool for interpreting DNA methylation findings from blood in the context of brain. *Translational Psychiatry*. 2017;7(8). doi:10.1038/tp.2017.171
90. Braun PR, Han S, Hing B, et al. Genome-wide DNA methylation comparison between live human brain and peripheral tissues within individuals. *Translational Psychiatry*. 2019;9(1):1-10. doi:10.1038/s41398-019-0376-y
91. Hannon E, Mansell G, Walker E, et al. Assessing the co-variability of DNA methylation across peripheral cells and tissues: Implications for the interpretation of findings in epigenetic epidemiology. *PLOS Genetics*. 2021;17(3):e1009443. doi:10.1371/JOURNAL.PGEN.1009443
92. Jiang R, Jones MJ, Chen E, et al. Discordance of DNA methylation variance between two accessible human tissues. *Scientific Reports*. 2015;5:8257. doi:10.1038/srep08257
93. Waterland RA, Kellermayer R, Laritsky E, et al. Season of conception in rural gambia affects DNA methylation at putative human metastable epialleles. *PLoS Genetics*. Published online 2010. doi:10.1371/journal.pgen.1001252
94. Garg P, Joshi RS, Watson C, Sharp AJ. A survey of inter-individual variation in DNA methylation identifies environmentally responsive co-regulated networks of epigenetic variation in the human genome. Relton C, ed. *PLOS Genetics*. 2018;14(10):e1007707-e1007707. doi:10.1371/journal.pgen.1007707
95. Gordon L, Joo JE, Powell JE, et al. Neonatal DNA methylation profile in human twins is specified by a complex interplay between intrauterine environmental and genetic factors, subject to tissue-specific influence. *Genome Research*. 2012;22(8):1395-1406. doi:10.1101/gr.136598.111
96. Do C, Lang CF, Lin J, et al. Mechanisms and Disease Associations of Haplotype-Dependent Allele-Specific DNA Methylation. *American Journal of Human Genetics*. 2016;98(5):934-955. doi:10.1016/j.ajhg.2016.03.027
97. Richmond RC, Sharp GC, Herbert G, et al. The long-term impact of folic acid in pregnancy on offspring DNA methylation: Follow-up of the Aberdeen folic acid supplementation trial (AFAST). *International Journal of Epidemiology*. 2018;47(3):928-937. doi:10.1093/ije/dyy032
98. Ambatipudi S, Cuenin C, Hernandez-Vargas H, et al. Tobacco smoking-associated genome-wide DNA methylation changes in the EPIC study. *Epigenomics*. 2016;8(5):599-618.
99. James P, Sajjadi S, Tomar AS, et al. Candidate genes linking maternal nutrient exposure to offspring health via DNA methylation: A review of existing evidence in humans with specific focus on one-carbon metabolism. *International Journal of Epidemiology*. Published online 2018. doi:10.1093/ije/dyy153
100. Wolff GL, Kodell RL, Moore SR, Cooney CA. *Maternal Epigenetics and Methyl Supplements Affect Agouti Gene Expression in a vy/a Mice.*; 1998. <http://www.fasebj.org>
101. Cooney CA, Dave AA, Wolff GL. Maternal methyl supplements in mice affect epigenetic variation and DNA methylation of offspring. In: *Journal of Nutrition*. ; 2002. doi:10.1093/jn/132.8.2393s
102. Dolinoy DC, Weidman JR, Waterland RA, Jirtle RL. Maternal genistein alters coat color and protects Avy mouse offspring from obesity by modifying the fetal epigenome. *Environmental Health Perspectives*. Published online 2006. doi:10.1289/ehp.8700
103. Dolinoy DC, Huang D, Jirtle RL. *Maternal Nutrient Supplementation Counteracts Bisphenol A-Induced DNA Hypomethylation in Early Development.*; 2007. <http://www.pnas.org/cgi/doi/10.1073/pnas.0703739104>
104. Costello KR, Leung A, Trac C, et al. Mechanisms regulating interindividual epigenetic variability at transposable elements. doi:10.1101/2021.06.01.446659
105. Imbeault M, Helleboid PY, Trono D. KRAB zinc-finger proteins contribute to the evolution of gene regulatory networks. *Nature*. 2017;543(7646):550-554. doi:10.1038/nature21683
106. Shen P, Xu A, Hou Y, et al. Conserved paradoxical relationships among the evolutionary, structural and expressional features of KRAB zinc-finger proteins reveal their special functional characteristics. *BMC Molecular and Cell Biology* 2021 22:1. 2021;22(1):1-15. doi:10.1186/S12860-021-00346-W

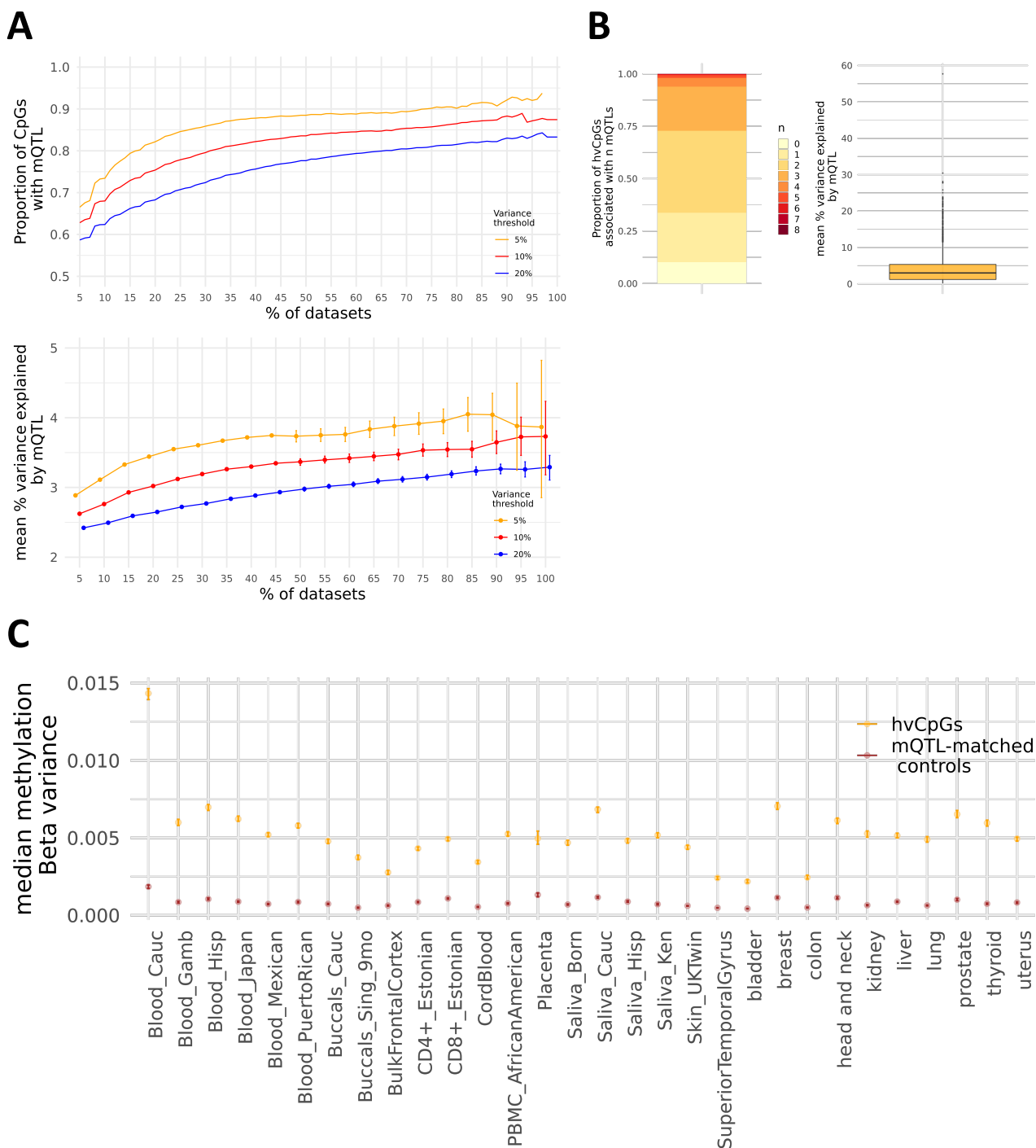
107. Candler T, Kessler N, Gunasekara C, et al. DNA methylation at a nutritionally sensitive region of the PAX8 gene is associated with thyroid volume and function in Gambian children. *Science Advances*. 2021;7(45):eabj1561.
108. Feinberg AP, Irizarry RA. Stochastic epigenetic variation as a driving force of development, evolutionary adaptation, and disease. *Proceedings of the National Academy of Sciences of the United States of America*. 2010;107 Suppl(suppl 1):1757-1764. doi:10.1073/pnas.0906183107
109. Liu Y, Aryee MJ, Padyukov L, et al. Epigenome-wide association data implicate DNA methylation as an intermediary of genetic risk in rheumatoid arthritis. *Nature biotechnology*. 2013;31(2):142-147. doi:10.1038/nbt.2487
110. Matzaraki V, Kumar V, Wijmenga C, Zhernakova A. The MHC locus and genetic susceptibility to autoimmune and infectious diseases. *Genome Biology*. 2017;18(1):76. doi:10.1186/s13059-017-1207-1
111. Sala C, Lena P di, Durso DF, Prodi A, Castellani G, Nardini C. Evaluation of pre-processing on the meta-analysis of DNA methylation data from the Illumina HumanMethylation450 BeadChip platform. *PLOS ONE*. 2020;15(3):e0229763. doi:10.1371/JOURNAL.PONE.0229763

**FIGURE 1**



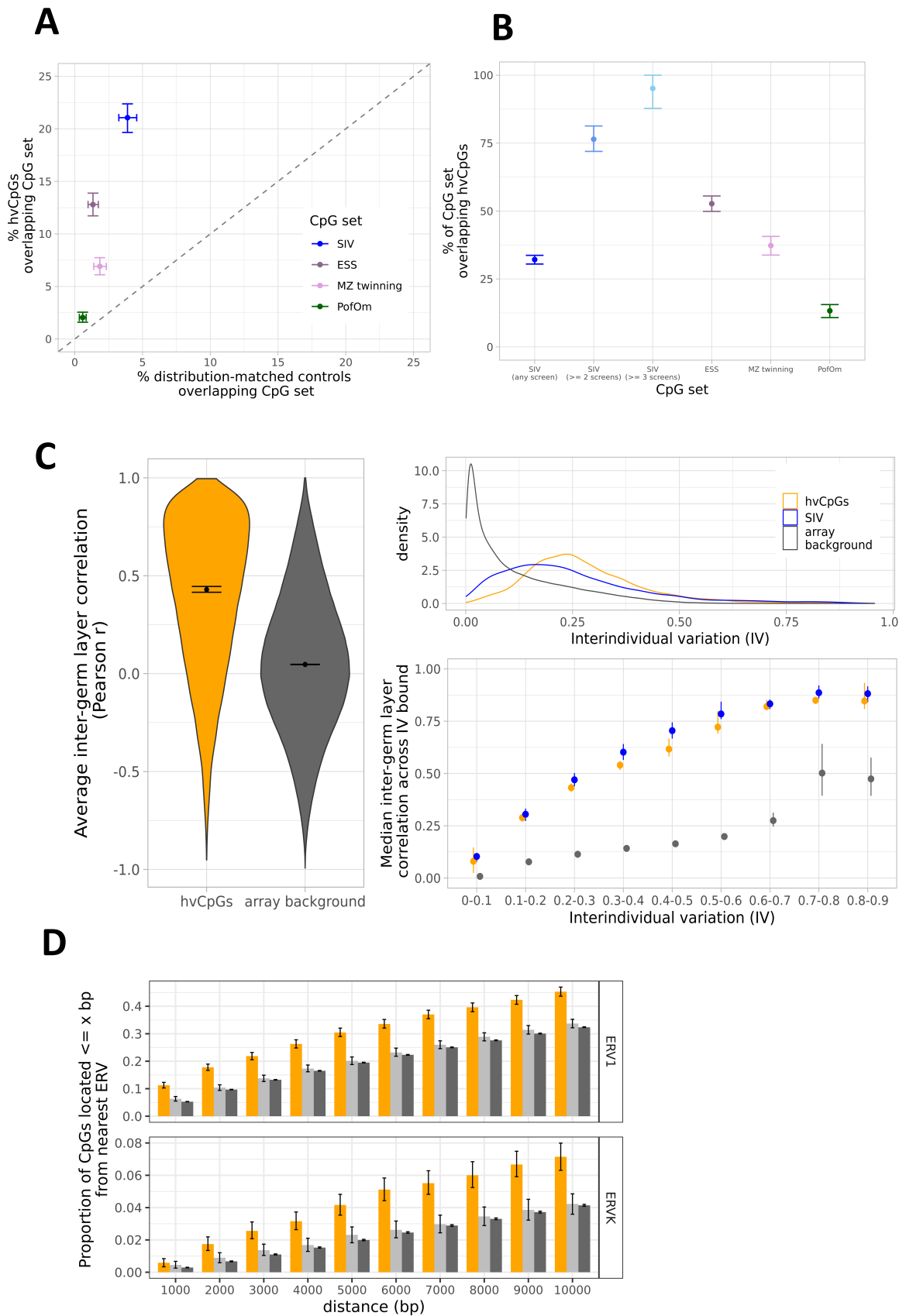
**Figure 1. Schematic of the method for identifying tissue- and ethnicity- independent hypervariable CpGs (hvCpGs).** The top 20%, 19%, 18% ... 1% of variable CpGs by methylation Beta variance were first extracted from each of the 30 methylation datasets used in this study. The intersection of these CpGs was then taken over an increasing number of datasets ( $n$ ), requiring each CpG to be present in a minimum of 15 out of the 30 datasets analysed (Supplementary Fig. 1).

FIGURE 2



**Figure 2. Genetic effects at hvCpGs using mQTL data from a large meta GWAS in blood (Min et al. 2021).** **A)** The relationship between hypervariability and the proportion of CpGs with at least one mQTL association (top) and the mean mQTL effect size (bottom). Coloured curves represent CpGs with top 5% (orange), 10% (red) and 20% (blue) methylation Beta variance in at least x% of datasets. **B)** The distributions of the number of mQTL associations (left) and mean % variance explained by mQTL (right) at hvCpGs. **C)** Median methylation Beta variance at 3,722 hvCpGs overlapping the 'Blood\_Cauc' dataset (orange) and corresponding controls matched on number of mQTL associations and mean % explained by mQTL ('mQTL-matched controls', Table1; Supplementary Fig. 4), in each dataset. Error bars in A and C are bootstrapped 95% confidence intervals. Note, error bars in C are very small.

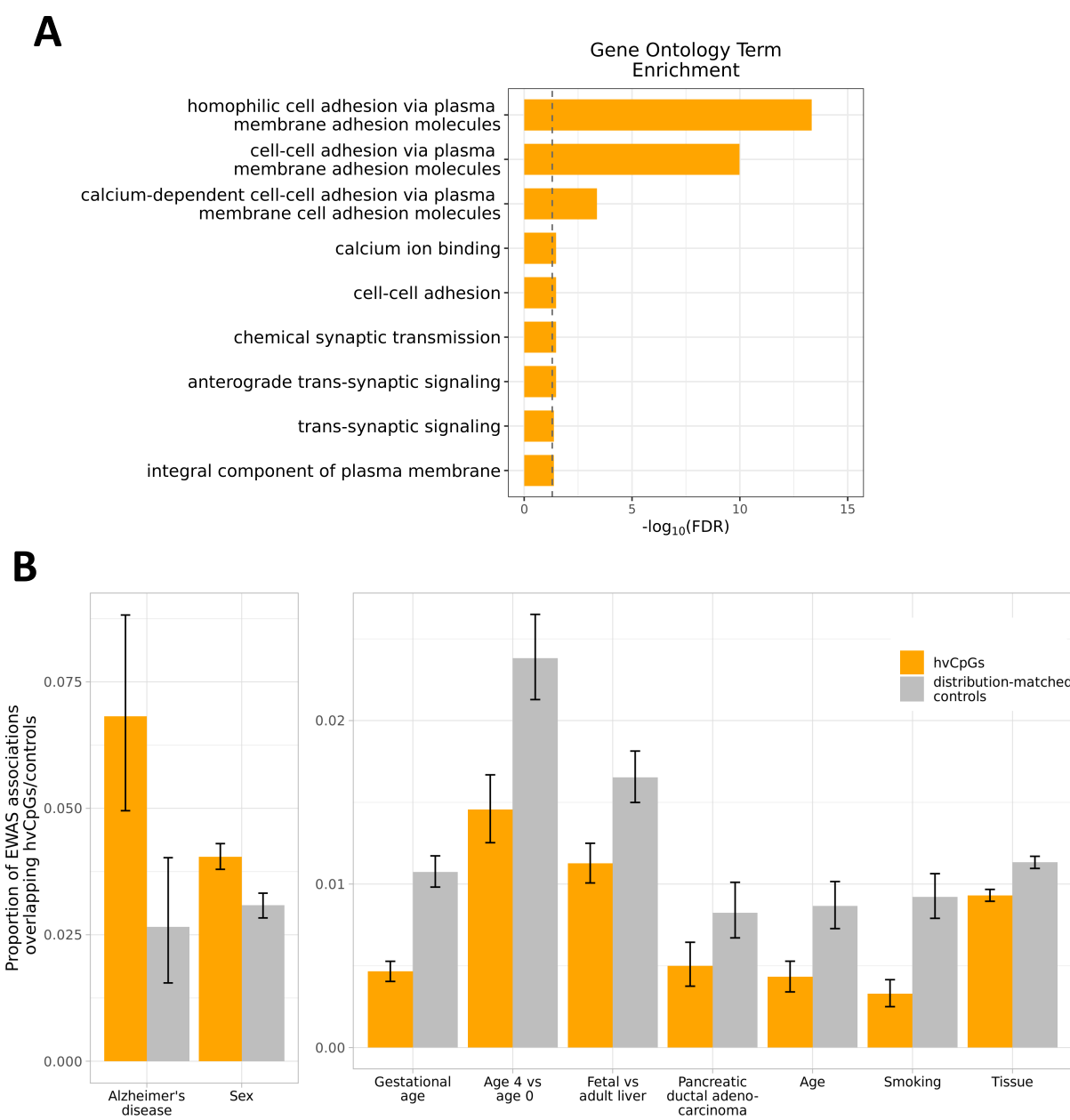
FIGURE 3



**Figure 3. hvCpGs are enriched for loci and genomic features linked to variable methylation establishment in early development.** **A)** The proportion of 3,566 hvCpGs (y-axis) vs corresponding distribution-matched controls (x-axis) covered in the 'Blood\_Cauc' dataset that overlap 3,089 SIV-CpGs, 1,217 ESS CpGs identified by van Baak *et al.* (2018), 728 'MZ twinning' CpGs identified by van Dongen *et al.* (2021) and 732 PofOm CpGs identified by Zink *et al.* (2018). **B)** The proportion of SIV-CpGs, ESS CpGs, MZ twinning CpGs and PofOm CpGs that are hvCpGs. SIV-CpGs identified in at least two or three independent screens were also included in this plot. **C)** Inter-germ layer correlations at hvCpGs using a fetal multi-tissue dataset that comprises methylation data from 10 individuals with endoderm- and ectoderm-derived tissues, 9 individuals with endoderm- and mesoderm- derived tissues and 8 individuals with mesoderm- and ectoderm- derived tissues (see Supplementary Table 7). **Left:** The distribution of average inter-germ layer correlations at 3,878 hvCpGs (orange) and 372,571 array background CpGs (excluding previously identified SIV CpGs and hvCpGs) (dark grey) covered in the fetal multi-tissue dataset. **Top Right:** Interindividual variation at 3,878 hvCpGs (orange), 4,076 previously identified SIV loci (blue) covered in the fetal multi-tissue dataset, and 372,571 array background CpGs (see 'Methods' for definition of interindividual variation). **Bottom Right:** Comparison of average inter-germ layer correlations at hvCpGs, SIV-CpGs and array background CpGs, stratified by interindividual variation. Each point indicates the median average inter-germ layer correlation for those CpGs with interindividual variation falling within each bound specified on the x-axis. **C)** The proportion of 3,566 hvCpGs, distribution-matched controls and array background CpGs that are  $\leq x$  bp from the nearest ERV1 and ERVK transposable elements determined by RepeatMasker. Error bars in all panels are bootstrapped 95% confidence intervals. SIV = systemic interindividual variation, ESS = epigenetic supersimilarity, PofOm = parent-of-origin-specific methylation.



FIGURE 4



**Figure 4. Functional annotation of hvCpGs. A)** Gene ontology term enrichment analysis at hvCpGs. Vertical line indicates a significance threshold of  $\text{FDR} < 0.05$ . **B)** Enrichment (left) and depletion (right) of hvCpGs amongst EWAS trait associations relative to blood distribution-matched controls. Y-axis gives the proportion of EWAS trait associations that comprise hvCpGs and controls. Only traits overlapping at least 1% of hvCpGs were considered (see ‘Methods’ and Supplementary Table 9 for further details). Error bars are bootstrapped 95% confidence intervals.

## Tables

**Table 1. Main CpG sets used in this study.**

CpG set	n	Notes
hvCpGs	4143	CpGs within top 5% methylation Beta variance in at least 65% datasets in which the CpG is covered, requiring the hvCpG to be covered in at least 15 datasets and to be reported as technically reliable.
array background	406306	CpGs covered in at least 15 of the 30 datasets used in this study.
distribution-matched controls	3566	Array background CpGs with similar methylation Beta distributions to hvCpGs in the 'Blood_Cauc' dataset, requiring each control CpG to be technically reliable.
de-clustered hvCpGs	2640	A set of hvCpGs in which no CpGs is within 4kb of another CpG.
mQTL-matched controls	3722	CpGs reported by the GoDMC meta-GWAS <sup>41</sup> with the same number of mQTL associations and similar mean % variance explained by an mQTL, requiring each control CpG to be present in at least as many datasets as the hvCpG.

**Table 2. Published CpG sets used in this study**

CpG set	Description	n. CpGs overlapping array background	Reference
SIV	Interindividual methylation variation with concordant methylation across tissues derived from different germ layers within a given individual. See Supplementary Table 8.	3089	Harris <i>et al.</i> , van Baak <i>et al.</i> , Kessler <i>et al.</i> , Gunasekara <i>et al.</i> <sup>24-27</sup>
ESS	Greater-than-expected methylation similarity between MZ co-twins	1217	van Baak <i>et al.</i> <sup>25</sup>
MZ twinning CpGs	Probes differentially methylated between MZ and DZ twins.	728	van Dongen <i>et al.</i> <sup>28</sup>
evCpGs	MZ co-twin methylation discordance that is equivalent to methylation discordance between unrelated individuals in whole blood. A subset of these replicated in adipose tissue.	317 (blood) 145 (blood & adipose)	Planterose Jiménez <i>et al.</i> <sup>45</sup>
SoC	CpGs at which methylation is associated with season of conception in Gambian children.	242	Silver <i>et al.</i> <sup>29</sup>
PofOm	Regions of parent-of-origin-specific methylation identified in peripheral blood from Icelandic individuals.	732 CpGs in 116 PofOm regions	Zink <i>et al.</i> <sup>68</sup>

SIV = systemic interindividual variation, ESS = epigenetic supersimilarity, evCpGs = equivalently variable CpGs, SoC = season-of-conception, PofOm = parent-of-origin-specific methylation.