

- ⁵⁴ McGowan PO, Suderman M, Sasaki A *et al.* Broad epigenetic signature of maternal care in the brain of adult rats. *PLoS ONE* 2011;**6**:e14739.
- ⁵⁵ Galobardes B, Shaw M, Lawlor DA, Lynch JW, Davey Smith G. Indicators of socioeconomic position (part 2). *J Epidemiol Community Health* 2006;**60**:95–101.
- ⁵⁶ Galobardes B, Shaw M, Lawlor DA, Lynch JW, Davey Smith G. Indicators of socioeconomic position (part 1). *J Epidemiol Community Health* 2006;**60**:7–12.
- ⁵⁷ Rakyan VK, Down TA, Thorne NP *et al.* An integrated resource for genome-wide identification and analysis of human tissue-specific differentially methylated regions (tDMRs). *Genome Res* 2008;**18**:1518–29.
- ⁵⁸ Atherton K, Fuller E, Shepherd P, Strachan DP, Power C. Loss and representativeness in a biomedical survey at age 45 years: 1958 British birth cohort. *J Epidemiol Community Health* 2008;**62**:216–23.

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/3.0>), which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.

Published by Oxford University Press on behalf of the International Epidemiological Association

International Journal of Epidemiology 2012;**41**:74–78

© The Author 2012; all rights reserved. Advance Access publication 23 January 2012

doi:10.1093/ije/dyr225

Commentary: The seven plagues of epigenetic epidemiology

Bastiaan T Heijmans^{1,2*} and Jonathan Mill³

¹Molecular Epidemiology, Leiden University Medical Center, Leiden, The Netherlands, ²Netherlands Consortium for Healthy Ageing, Leiden, The Netherlands and ³Institute of Psychiatry, King's College London, London, UK

*Corresponding author. Molecular Epidemiology, Leiden University Medical Center, Postal Zone S-5-P, PO Box 9600, 2300 RC, Leiden, The Netherlands. E-mail: bas.heijmans@lumc.nl

Accepted 1 December 2011

Epigenetics is being increasingly combined with epidemiology to add mechanistic understanding to associations observed between environmental, genetic and stochastic factors and human disease phenotypes. Currently, epigenetic epidemiological studies primarily focus on exploring if and where the epigenome (i.e. the overall epigenetic state of a cell) is influenced by specific environmental exposures like prenatal nutrition,¹ sun exposure² and smoking.³ In this issue of the *IJE*, Nada Borghol *et al.*⁴ report an association between childhood social-economic status (SES) and differential DNA methylation in adulthood. Low SES may integrate diverse and heterogeneous environmental influences, and knowing which epigenetic changes are associated with low SES may provide clues about the biological processes underlying its health consequences. The authors stress that their study is preliminary. This statement is, in fact, to a greater or lesser extent applicable to the entire first wave of studies currently being published that likewise aim to discover associations between epigenetic variation measured on a genome-wide scale and environmental exposures or disease phenotypes. When executing such epigenome-wide association studies (EWASs),⁵

every epigenetic epidemiologist is struggling with the same biological, technical and methodological issues. It is important to take these into consideration when designing a study and interpreting the results. Let us consider seven of those issues, taking the current study on SES as a starting point.

We do not really know where to look, or what to look for

Most epigenetic epidemiological studies focus on DNA methylation for various practical and biological reasons, neglecting other layers of the epigenome-like histone modifications that are also likely to be important in influencing disease phenotypes. Our basic understanding of the methylome (i.e. the whole of DNA methylation marks on the genome) is in its infancy, and we are still learning about the specific localization of the features that, when differentially methylated, regulate gene expression and are thus relevant for epigenetic epidemiologists to study. The current study, like many others, evaluated promoter regions, in this case defined as 1000 bp upstream to

250 bp downstream of transcription start sites. Although these features are often enriched for DNA methylation marks influencing the expression of genes, recent work suggests that other regions of the methylome outside of promoters, including inter-genic CpG island shores⁶ and intra-genic CpG islands,⁷ may ultimately be more important for regulating phenotypic variation.

For any differentially methylated region identified in EWASs it will be important to demonstrate functionality. Promoter methylation in the current study was integrated with public gene expression data and, as expected, highly expressed genes were more commonly flanked by less methylated promoters and vice versa. A limitation is that this observation is for groups of promoters, whereas information is needed about this relationship for individual promoters. Mining the reference epigenomes and transcriptomes that are being generated for different cell types under the umbrella of initiatives such as the National Institutes of Health (NIH) Epigenomics Roadmap⁸ and the International Human Epigenome Consortium⁹ may contribute to such information. Additional *in vitro* experiments will be required to evaluate the transcriptional effects of differential DNA methylation at a specific locus independent of its genomic context.¹⁰

We have to rely on imperfect technology

The good news is that recent advances in genomic technology mean that genome-scale studies of DNA methylation across multiple samples are now feasible. In practice, however, one has to compromise between coverage and precision in epidemiological studies, which likely incorporate a large number of samples. A large (and growing) number of methods exist for assessing DNA methylation both genome wide and at specific CpG sites,¹¹ and one problem relates to our inability to compare results across studies that have used different platforms. On the one hand there are methods such as that used in the current study in which the methylated portion of the genome is captured using antibodies against methylated DNA and subsequently quantified using microarrays or next-generation sequencing. These approaches can provide coverage across most of the genome and may be optimally suited to discriminate low from high methylation, but have lower reliability for smaller differences and are biased by factors such as CG density.^{12,13} On the other hand, there are methods based on the bisulphite conversion of DNA combined with next-generation sequencing that provide higher accuracy and single nucleotide resolution. Although whole-genome bisulphite sequencing is currently unfeasible to use across large epidemiological cohorts, the method can be adapted to target a reduced representation of the genome (approximately 3 million out of approximately 28 million CG dinucleotides in the

human genome).^{12,13} The recently launched Illumina 450k Methylation Beadchip may offer a balance between coverage and precision, which will be attractive for epidemiological EWASs executed during the next few years.⁵ It interrogates DNA methylation at over 480 000 CG dinucleotides, is high-throughput and relatively affordable. The precision of this platform appears to compare well with some of the other platforms,^{12,13} but these results should be interpreted with caution. Although correlation coefficients reported across the various platform comparisons are high, they are mainly driven by the fact that the large majority of the genome is either unmethylated or fully methylated, and substantial discrepancies between platforms may exist for intermediate level methylation.^{12,14} Therefore, the technological validation of findings using an independent method remains important. This will be feasible for a small number of 'top hits', like the three proadherin promoters assessed in the current study. However, validating the outcomes of the complex pathway analyses performed to implicate either entire biological processes (such as extra- and intra-cellular signalling in the current study) or genomic features with a specific function in gene regulation [e.g. promoters, enhancers, inter/intragenic CG island (shores) etc.], is more demanding and currently not realized. Validating the results of such gene-set testing methods will entail the re-assessment of DNA methylation across large sets of loci.

We may be limited by available sample sizes that are optimal for epigenetic epidemiology

The current study investigated only 40 individuals. Investigators will be able to secure budgets for larger studies as empirical data increasingly highlight the value of epigenetic epidemiology, and high-throughput, economical laboratory approaches become more widely adopted. Nevertheless, it is unlikely that the simple brute-force approach that has been used relatively successfully in genome-wide association studies (GWASs) is valid for EWASs. In genetics, many of the epidemiological principles about designing studies with respect to selection biases, confounding, batch effects and appropriateness of controls could largely be replaced by the simple rule 'bigger-is-better'. This is not true for epigenetic epidemiology, because the epigenome is not a static entity like the genome, which necessitates the use of more conventional epidemiological approaches.¹⁵ Further complicating matters is the fact that, for the most powerful study designs in epigenetic epidemiology (including studies of discordant monozygotic twins¹⁶ particularly when longitudinally sampled,¹⁷ early exposure studies with long-term follow-up,¹ and studies of specific cell types¹⁸), the number of eligible individuals for whom relevant biological

materials were stored in existing epidemiological cohorts were often limited, and it will be difficult to scale-up analyses to include the thousands of samples that may be required for establishing robust associations with disease phenotypes. Moving forward, it will be important to establish cause and effect in epigenetic epidemiology; disease-associated differentially methylated regions may arise prior to illness and contribute to the disease phenotype or could be a secondary effect of the disease process, or the medications used in treatment.¹⁹ Furthermore, maximum information will be obtained from epidemiological studies that are able to integrate epigenomic information with genomic, transcriptomic and proteomic data obtained from the same samples.

Whatever we do, it may never be enough to fully account for epigenetic differences between tissues and cells

In many respects, large comprehensively phenotyped and longitudinally sampled epidemiological studies, like the 1958 British birth cohort used in the current study, are an ideal resource for epigenetic epidemiology. In nearly all of these studies, however, whole blood is the only biological material that has been archived. Blood is a heterogeneous tissue and any DNA methylation difference between groups could be confounded by differences in the cellular composition of whole blood samples, for example, resulting from the immune response to sub-clinical infection. The good news is that fewer than perhaps expected DNA methylation differences exist between leucocyte types, and controlling for cellular heterogeneity may be possible in biobanks with a simple blood cell count.²⁰ Whether the latter is sufficient (and under which circumstances it is not), however, remains to be established. Epigenomic studies of separate cell types such as those being undertaken by the NIH Epigenomic Roadmap Initiative and the European Union Blueprint consortium are currently generating reference epigenomes of haematopoietic cells that will be of great utility in this regard.⁸ When moving beyond associations with environmental exposures to epigenetic associations with phenotypes, a key question for epigenetic epidemiology concerns the extent to which easily accessible peripheral tissues (such as blood) can be used to ask questions about inter-individual phenotypic variation manifest in inaccessible tissues such as the brain, visceral fat and other internal organs and tissues. Cross-tissue comparisons of the methylome within the same individual are currently underway to establish the relationship between epigenetic patterns in blood with other tissues. Although these analyses are crucial, the results may not be generally applicable; higher inter-tissue concordance may be present for DNA methylation

changes induced early in development (and potentially propagated soma-wide) than for changes occurring during ageing that are more likely to remain tissue specific.^{19,21} Efforts to obtain biopsies (subcutaneous fat, muscle, etc.) and post-mortem material in subsets of longitudinal biobanks will greatly increase their value for epigenetic studies, despite the problems associated with cellular heterogeneity that also hold for such samples.

We may be trying to detect inherently small effect sizes using these sub-optimal methods and sample cohorts

The main findings in the current study concerned DNA methylation differences at three procadherin promoters.⁴ The extent of the difference at these promoters was similar to those commonly observed in other recent studies, namely ~5%,⁵ and was most apparent for a single, nominally statistically significant CG dinucleotide in each region. The biological implications of such small alterations in DNA methylation in terms of gene expression and function are unknown. Although DNA methylation is recognized as one of the most stable epigenetic marks, it is still relatively dynamic and this has important implications for epigenetic epidemiology. The randomness of maintaining and mitotically transmitting DNA methylation patterns may potentially dilute the putative epigenetic signatures of an adverse exposure early in life (e.g. to low SES in childhood) observed decades later. Of note, recent studies indicate that DNA methylation patterns in leucocytes undergo considerable changes during the first years of life.²² Thus on top of the previously discussed question of whether DNA methylation at a specific locus actually influences transcriptional activity, researchers should also aim to establish whether the small DNA methylation differences often observed between groups—either expressed as absolute difference, relative difference or relative to the variation in the population—translate into differences in gene expression in the relevant tissue. It will be of particular interest to see whether the effects of such modest differences, while perhaps of little consequence individually, may shift transcription of a biological process or functional network when they co-occur with other changes to the methylome.²³ Little is known about the actual scale and extent of between-individual variation in DNA methylation across the genome. In this regard, public genome-scale resources need to be created that document *inter-individual differences* in DNA methylation and gene expression, in addition to the reference epigenomes that are currently being generated.

We lack a framework for the analysis of genome-wide epigenetic data

The results of GWASs are relatively easy to judge. Quality-control steps are well-defined and reported, individually testing every genetic variant [i.e. single nucleotide polymorphism (SNP)] is straightforward, and levels of genome-wide statistical significance are clear. For EWASs, the analytical methodology is very much under construction. For example, in the current study it was not possible to attain genome-wide levels of significance, which is acceptable for an exploratory study, but makes it difficult to fully interpret the reported differences. Because of the vast range of methods currently being used to assess DNA methylation, meta-analyses across different studies are difficult. The adoption of a common technology platform, such as the new Illumina 450k Methylation Beadchip, across multiple studies would provide an excellent opportunity to converge on widely accepted guidelines for the analysis and integration of EWAS data. Apart from pre-processing procedures (quality control, normalization, handling different probe types, accounting for genetic variation, etc.), elements of these guidelines should deal with the analysis of individual CG dinucleotides vs groups of (correlated) adjacent CGs, the use of genome annotations in the analysis (histone states, promoter types, CG content, etc.), and levels of epigenome-wide significance for various analyses. An important aspect will be the exploration of the previously mentioned gene-set testing methods in the context of DNA methylation since they will be vital to obtain meaningful interpretations of genome-wide data in terms of underlying biological processes or genomic functions [e.g. promoters, enhancers, inter/intragenic CG island (shores), etc.]. For example, commonly used enrichment methods assume independence within a gene set and, apart from consistency in biological signal in a gene set, statistical significance may reflect consistency in other characteristics such as GC content, coverage or other sequence features.²⁴ Alternative implementations of gene-set testing methods include global testing approaches.²⁵ Finally, it will be important to adopt an integrative paradigm based on the combination of genetic and epigenetic epidemiological data.²⁶ Of particular relevance in this respect is evidence for the widespread occurrence of allele-specific DNA methylation (ASM) across the genome. Recent studies have shown that there are considerable inter-individual differences in ASM, which are frequently associated with genetic variation but can also be mediated by genomic imprinting (i.e. the parent-of-origin dependent silencing of expression by epigenetic mechanisms), environmental influences and apparently stochastic factors in the cell.^{27,28} ASM can mask the effect of risk alleles by silencing their expression, and also provides a potential mechanism underlying gene-environment

interactions.²⁶ Furthermore, ASM may contribute towards the apparent 'missing heritability' of many complex diseases and the low penetrance often reported for SNPs identified by GWASs.²⁹

We have to manage high expectations

There is a considerable interest in epigenetic research in the popular press. The current study is a vivid illustration: even though the authors deem it preliminary, it was widely covered by the media.³⁰ Epigenetics should avoid some of the hype that surrounded the early days of genetic epidemiology. After the draft human genome sequence was announced in 2001, it was widely perceived that we would soon understand the causes of most common diseases and how to treat them. This expectation was not realistic, but not always renounced by geneticists. Currently, many scientists outside the field are disappointed by results of human genetics, and in particular GWASs, despite their overall considerable success. Genetic epidemiology has proven to be harder than expected despite the favourable starting point of thousands of Mendelian diseases and the high heritabilities associated with most traits to be explained. Very much like genetics, epigenetics will not be able to deliver the miracles it is sometimes claimed it will.

In conclusion, epigenetic epidemiology is early in its development and susceptible to new ideas and approaches. Only a few years ago empirical papers were greatly outnumbered by reviews. Now, reference epigenomes are produced at great pace (see <http://epigenomeatlas.org>).^{8,9} Moreover, furthered by pilot studies like the one from Nada Borghol *et al.*,⁴ the outline of the infrastructure required for EWASs is emerging. Crucial elements include optimal study designs, benchmarking technology and data analysis approaches that are statistically and biologically sound. An additional key aspect to the successful design and interpretation of epigenetic epidemiological studies will be the creation of public genome-scale resources focusing on inter-individual variation incorporating epigenomic, DNA sequence and transcriptomic data. Education, hard work and a certain degree of luck will get us there—not very different to the remedy against low SES.

Funding

NGI/NWO (#93518027, to B.T.H.); NGI/NWO-funded Netherlands Consortium for Healthy Ageing (NCHA) (#05060810, B.T.H.); NIH grant (AG036039, to J.M.).

Acknowledgement

We thank Elmar Tobi for his comments.

Conflict of interest: None declared.

References

- ¹ Tobi EW, Lumey LH, Talens RP *et al.* DNA methylation differences after exposure to prenatal famine are common and timing- and sex-specific. *Hum Mol Genet* 2009;**18**: 4046–53.
- ² Gronniger E, Weber B, Heil O *et al.* Aging and chronic sun exposure cause distinct epigenetic changes in human skin. *PLoS Genet* 2010;**6**:e1000971.
- ³ Breitling LP, Yang R, Korn B, Burwinkel B, Brenner H. Tobacco-smoking-related differential DNA methylation: 27K discovery and replication. *Am J Hum Genet* 2011;**88**: 450–57.
- ⁴ Borghol N, Suderman M, McArdle W *et al.* Associations with early-life socio-economic position in adult DNA methylation. *Int J Epidemiol* 2012;**41**:62–74.
- ⁵ Rakyan VK, Down TA, Balding DJ, Beck S. Epigenome-wide association studies for common human diseases. *Nat Rev Genet* 2011;**12**:529–41.
- ⁶ Irizarry RA, Ladd-Acosta C, Wen B *et al.* The human colon cancer methylome shows similar hypo- and hypermethylation at conserved tissue-specific CpG island shores. *Nat Genet* 2009;**41**:178–86.
- ⁷ Deaton AM, Webb S, Kerr AR *et al.* Cell type-specific DNA methylation at intragenic CpG islands in the immune system. *Genome Res* 2011;**21**:1074–86.
- ⁸ Bernstein BE, Stamatoyannopoulos JA, Costello JF *et al.* The NIH Roadmap Epigenomics Mapping Consortium. *Nat Biotechnol* 2010;**28**:1045–8.
- ⁹ Anonymous. Moving AHEAD with an international human epigenome project. *Nature* 2008;**454**:711–5.
- ¹⁰ Klug M, Rehli M. Functional analysis of promoter CpG methylation using a CpG-free luciferase reporter vector. *Epigenetics* 2006;**1**:127–30.
- ¹¹ Laird PW. Principles and challenges of genomewide DNA methylation analysis. *Nat Rev Genet* 2010;**11**:191–203.
- ¹² Bock C, Tomazou EM, Brinkman AB *et al.* Quantitative comparison of genome-wide DNA methylation mapping technologies. *Nat Biotechnol* 2010;**28**:1106–14.
- ¹³ Harris RA, Wang T, Coarfa C *et al.* Comparison of sequencing-based methods to profile DNA methylation and identification of monoallelic epigenetic modifications. *Nat Biotechnol* 2010;**28**:1097–105.
- ¹⁴ Sandoval J, Heyn HA, Moran S *et al.* Validation of a DNA methylation microarray for 450,000 CpG sites in the human genome. *Epigenetics* 2011;**6**:692–702.
- ¹⁵ Relton CL, Davey Smith G. Epigenetic epidemiology of common complex disease: prospects for prediction, prevention, and treatment. *PLoS Med* 2010;**7**:e1000356.
- ¹⁶ Dempster EL, Pidsley R, Schalkwyk LC *et al.* Disease-associated epigenetic changes in monozygotic twins discordant for schizophrenia and bipolar disorder. *Hum Mol Genet* 2011. doi:10.1093/hmg/ddr416 [Epub 9 September 2011].
- ¹⁷ Wong CC, Caspi A, Williams B *et al.* A longitudinal study of epigenetic variation in twins. *Epigenetics* 2010;**5**: 516–26.
- ¹⁸ Ollikainen M, Smith KR, Joo EJ *et al.* DNA methylation analysis of multiple tissues from newborn twins reveals both genetic and intrauterine components to variation in the human neonatal epigenome. *Hum Mol Genet* 2010;**19**: 4176–88.
- ¹⁹ Heijmans BT, Tobi EW, Lumey LH, Slagboom PE. The epigenome: archive of the prenatal environment. *Epigenetics* 2009;**4**:526–31.
- ²⁰ Talens RP, Boomsma DI, Tobi EW *et al.* Variation, patterns, and temporal stability of DNA methylation: considerations for epigenetic epidemiology. *FASEB J* 2010;**24**: 3135–44.
- ²¹ Thompson RF, Atzmon G, Gheorghie C *et al.* Tissue-specific dysregulation of DNA methylation in aging. *Aging Cell* 2010;**9**:506–18.
- ²² Martino DJ, Tulic MK, Gordon L *et al.* Evidence for age-related and individual-specific changes in DNA methylation profile of mononuclear cells during early immune development in humans. *Epigenetics* 2011;**6**: 1085–94.
- ²³ Stoger R. The thrifty epigenotype: an acquired and heritable predisposition for obesity and diabetes? *Bioessays* 2008;**30**:156–66.
- ²⁴ Goeman JJ, Buhlmann P. Analyzing gene expression data in terms of gene sets: methodological issues. *Bioinformatics* 2007;**23**:980–87.
- ²⁵ Goeman JJ, van de Geer SA, de Kort F, van Houwelingen HC. A global test for groups of genes: testing association with a clinical outcome. *Bioinformatics* 2004;**20**:93–99.
- ²⁶ Meaburn EL, Schalkwyk LC, Mill J. Allele-specific methylation in the human genome: implications for genetic studies of complex disease. *Epigenetics* 2010;**5**: 578–82.
- ²⁷ Shoemaker R, Deng J, Wang W, Zhang K. Allele-specific methylation is prevalent and is contributed by CpG-SNPs in the human genome. *Genome Res* 2010;**20**:883–89.
- ²⁸ Schalkwyk LC, Meaburn EL, Smith R *et al.* Allelic skewing of DNA methylation is widespread across the genome. *Am J Hum Genet* 2010;**86**:196–212.
- ²⁹ Kong A, Steinthorsdottir V, Masson G *et al.* Parental origin of sequence variants associated with complex diseases. *Nature* 2009;**462**:868–74.
- ³⁰ Coghlan A. Childhood poverty leaves its marks on adult genetics. *New Scientist* 2011. [Epub 26 October 2011]; <http://www.newscientist.com/article/dn20255-childhood-poverty-leaves-its-mark-on-adult-genetics.html> (15 November 2011, date last accessed).