

# UC Berkeley

## UC Berkeley Previously Published Works

### Title

Estimation of blood cellular heterogeneity in newborns and children for epigenome-wide association studies

### Permalink

<https://escholarship.org/uc/item/8zq6m993>

### Journal

Environmental and Molecular Mutagenesis, 56(9)

### ISSN

0893-6692

### Authors

Yousefi, Paul  
Huen, Karen  
Quach, Hong  
[et al.](#)

### Publication Date

2015-12-01

### DOI

10.1002/em.21966

Peer reviewed



Published in final edited form as:

*Environ Mol Mutagen.* 2015 December ; 56(9): 751–758. doi:10.1002/em.21966.

## Estimation of blood cellular heterogeneity in newborns and children for epigenome-wide association studies

Paul Yousefi<sup>1</sup>, Karen Huen<sup>1</sup>, Hong Quach<sup>1</sup>, Girish Motwani<sup>1</sup>, Alan Hubbard<sup>1</sup>, Brenda Eskenazi<sup>1</sup>, and Nina Holland<sup>1</sup>

<sup>1</sup>School of Public Health, University of California, Berkeley, CA, USA

### Abstract

Confounding by cellular heterogeneity has become a major concern for epigenome-wide association studies (EWAS) in peripheral blood samples from population and clinical studies. Adjusting for white blood cell percentage estimates produced by the minfi implementation of the Houseman algorithm (minfi) during statistical analysis is now an established method to account for this bias in adults. However, minfi has not been benchmarked against white blood cell counts in children that may differ substantially from the reference dataset used in its estimation. We compared estimates of white blood cell type percentages produced by two methods, minfi and differential cell count (DCC), in a birth cohort at two time points (birth and 12 years of age). We found that both minfi and DCC had similar trends as children aged, and neither count method differed by sex among newborns ( $p > 0.10$ ). However, minfi estimates did not correlate well with DCC in samples from newborns ( $\rho = -0.05$  for granulocytes;  $\rho = -0.03$  for lymphocytes). In older children, correlation improved substantially ( $\rho = 0.77$  for granulocytes;  $\rho = 0.75$  for lymphocytes), likely due to increasing similarity with minfi's adult reference data as children aged. Our findings suggest that the minfi method may provide suitable estimates of white blood cell composition for samples from adults and older children, but may not currently be appropriate for EWAS involving newborns or young children.

### Keywords

Epigenetics; minfi; differential cell count; birth cohort; 450K

### Introduction

Epigenome-wide association studies (EWAS) have increasingly been used to identify novel biological mechanisms that contribute to disease status or respond to environmental exposures. Several large-scale DNA methylation assays have been developed in recent

---

Corresponding Author: Nina Holland, PhD, 733 University Hall, School of Public Health, UC Berkeley, CA 94720-7360, Phone: 510-665-2200, Fax: 510-665-2202, ninah@berkeley.edu.

Its contents are solely the responsibility of the authors and do not necessarily represent the official views of the NIEHS and the EPA.

#### Author contributions

Dr. Holland and Mr. Yousefi conceived and designed the study. Mr. Gotwani, and Ms. Quach performed the experiments. Mr. Yousefi performed the data analysis and prepared the manuscript with important intellectual contribution from Drs. Hubbard, Holland, Huen, and Eskenazi. All authors approved the final manuscript.

years, including methylated DNA immunoprecipitation (MeDIP), reduced representation bisulfite sequencing (RRBS), and whole genome bisulfite sequencing (wgBS) [Laird, 2010; Lister et al, 2009; Meissner et al, 2008; Weber et al, 2005], but due to its reliability, relatively low cost, and broad coverage, the Illumina Infinium HumanMethylation450 BeadChip® (450K) has been widely adopted in population-based EWAS [Bibikova et al, 2011; Liu et al, 2013; Sandoval et al, 2011; Teschendorff et al, 2009].

Unlike genetics, epigenetic markers may change over time or in response to exposures. DNA methylation in particular undergoes widespread remodeling *in utero* [Foley et al, 2009; Hughes, 2014; Perera and Herbstman, 2011]. For this reason and because early life exposures have been hypothesized to contribute differential risk towards later life ill health, performing EWAS at birth or in young children has been of great interest to investigators. Several EWAS, including those for prenatal exposure to smoking and arsenic [Joubert et al, 2012; Koestler et al, 2013], have quantified DNA methylation in cord blood. This strategy to assess epigenetic perturbation as near as possible to the prenatal period remains a high priority in light of the fetal origins of human disease hypothesis [Armstrong et al, 2014; Babenko et al, 2014; Barker, 1998; Essex et al, 2013].

Whole blood is a desirable matrix to use for EWAS as it is readily available and has been obtained for many human studies with a wide variety of initial aims (including past genome-wide association studies (GWAS)) [Chadwick et al, 2014; Liang and Cookson, 2014; Lowe and Rakyan, 2014; Michels et al, 2013]. However, as EWAS are more commonly performed in blood, there is growing awareness that heterogeneous white blood cell type populations may bias results due to confounding [Liang and Cookson, 2014; Lowe and Rakyan, 2014]. Since DNA methylation may vary by cell type, analyses involving health outcomes or exposures that also covary with cell type may be confounded. The consequence of such bias has been clearly demonstrated by Jaffe and Irizarry [Jaffe and Irizarry, 2014], who found that many published associations between blood-based CpG methylation and age were no longer statistically significant after adjustment for cell composition.

Several approaches have been proposed to address confounding bias in EWAS due to varying white blood cell type composition. One method is to restrict DNA methylation measurement to isolated populations of white blood cells. In practice, this requires performing fluorescence-activated cell sorting (FACS) prior to DNA isolation and subsequently quantifying DNA methylation signal in isolated cell populations. While appealing theoretically, this approach is not feasible for large population-based studies that rely on banked samples.

One alternative involves estimation of the relative proportions of different cell types, allowing for statistical adjustment for cellular mixture during data analysis. The performance of this approach depends largely on the quality of the estimate of cell type proportions. The most reliable white blood cell count is performed either by automated hematology analyzer, as part of a complete blood count (CBC), or retrospectively by microscopic differential cell count (DCC) using histologically stained blood smear slides.

However, since many epidemiologic studies do not have direct white blood counts, there is growing interest in computational approaches that estimate cell type proportions based on DNA methylation data. In 2012, Houseman et al. were the first to develop such a computational method, using 27k BeadChip results from n=46 isolated white cell samples from an unknown number of blood donors as a reference dataset [Bibikova et al, 2009; Houseman et al, 2012]. The updated version, produced by Jaffe and Irizarry [Jaffe and Irizarry, 2014] (referred to here as the minfi method), has seen the most widespread use because it was incorporated in a popular bioinformatic software pipeline for 450K data, and made several adjustments to specifically improve performance for 450K BeadChip data including the addition of a 450K BeadChip reference dataset (see methods for details). The minfi method is appealing in the context of EWAS studies because it can be readily implemented with no additional cost or data collection. However, the cell type estimates produced by minfi have not yet been systematically validated against a gold standard cell count, such as CBC or DCC. Additionally, minfi uses a small (n=6) cell-sorted 450K dataset from middle-aged Swedish men in its estimation procedure that may not be an appropriate reference when estimating cell composition in infants and children [Reinius et al, 2012].

Here, we conduct a comparison of the estimates of the relative abundance of white blood cell types produced by two methods, minfi and DCC, with randomly selected samples from a large epidemiologic cohort followed by the Center for the Health Assessment of Mothers and Children of Salinas (CHAMACOS) study at birth and at 12 years of age with 450K BeadChip data. We report findings showing that reference data and other assumptions should be carefully considered prior to utilizing computationally derived white blood cell estimates in EWAS studies in cord samples.

## Materials and methods

### Study population

The CHAMACOS study is a longitudinal birth cohort study of the effects of exposure to pesticides and environmental chemicals on the health and development of Mexican-American children living in the agricultural region of Salinas Valley, CA. Detailed description of the CHAMACOS cohort has previously been published [Eskenazi et al, 2003; Eskenazi et al, 2004]. Briefly, 601 pregnant women were enrolled in 1999–2000 at community clinics and 527 liveborn singletons were born. Follow up visits occurred at regular intervals throughout childhood, including a visit at 12 years of age that included only male child participants. For this analysis, we include the subset of subjects that had 450K BeadChip data available at birth (n=151) and matched data for the 12-year follow up (n=60). DCC analysis included the subset of subjects with whole blood smear slides available at the birth (n=111) and 12-year visits (n=45). Both newborns and 12 year olds included in the sample were healthy at the time of blood collection according to the study protocol, and confirmed by abstracted medical records and questionnaires. All subjects included in the subset were Latino in ancestry and 94.0% had at least one Mexican-born parent. Study protocols were approved by the University of California, Berkeley Committee for Protection of Human Subjects. Written informed consent was obtained from all mothers and assent was provided at the 12-year visit.

## Blood collection and processing

Whole blood was collected in BD vacutainers (Becton, Dickinson and Company, Franklin Lakes, NJ) containing either heparin anticoagulant or no anticoagulant. Whole blood smear slides were prepared from heparinized blood using the push-wedge blood smearing technique [Turgeon, 2011] and stored at  $-20^{\circ}\text{C}$  until staining. Aliquots of blood clot were stored at  $-80^{\circ}\text{C}$  until DNA isolation.

## DNA preparation

DNA isolation was performed using QIAamp DNA Blood Maxi Kits (Qiagen, Valencia, CA) according to manufacturer's protocol with small, previously described modifications [Holland et al, 2006]. Following isolation, all samples were checked for DNA quality and quantity by Nanodrop 2000 Spectrophotometer (Thermo Scientific, Waltham, MA). Those with good quality (260/280 ratio exceeding 1.6) were normalized to a concentration of 55ng/ul.

## 450K BeadChip DNA methylation analysis

DNA samples were bisulfite converted using Zymo Bisulfite Conversion Kits (Zymo Research, Irvine, CA), whole genome amplified, enzymatically fragmented, purified, and applied to Illumina Infinium HumanMethylation450 BeadChips (Illumina, San Diego, CA) according to manufacturer protocol. 450K BeadChips were handled by robotics and analyzed using the Illumina Hi-Scan system. DNA methylation was measured at 485,512 CpG sites.

## White blood cell composition estimation

White cell composition was characterized by two different methods in whole blood:

**Differential cell counts (DCC)**—Whole blood smears were stained utilizing a DiffQuik® staining kit, a modern commercial variant of the Romanovsky stain, a histological stain used to differentiate cells on a variety of smears and aspirates. This staining highlights cytoplasmic details and neurosecretory granules, which are utilized to characterize the differential white blood count. The staining kit is composed of a fixative (3:1 methanol: acetic acid solution), eosinophilic dye (xanthene dye), basophilic dye (dimethylene blue dye) and wash (deionized water). For consistency and to ensure the best results the slides were all fixed for 15 minutes at  $23^{\circ}\text{C}$  (room temperature), stained in both the basophilic dye and eosinophilic dye for five seconds each and washed after each staining period to prevent the corruption of the dye.

Slides were scored for white blood cell type composition by Zeiss Axioplan light microscope with 100 $\times$  oil immersion lens. Scoring was conducted at the perceived highest density of white blood cells using the standard battlement track scan method, which covers the entire width of a slide examination area. The counts for each of the five cell types (lymphocytes, monocytes, neutrophils, eosinophils, and basophils) were recorded by a dedicated mechanical counter. At least 100 cells were scored for each slide. Scoring reliability was initially validated by repeated scoring of 5 sets of 100 cells from the same slide with excellent reproducibility (CV = 5%).

**Minfi cell count estimation**—Results from the 450K BeadChip analysis were stored as raw IDAT files, and read into the *minfi* (v1.10.2) Bioconductor R package [Aryee et al, 2014] using the *read.450k.exp* function. Estimation of the six (CD8+ T and CD4+ T lymphocytes, CD56+ natural killer cells, CD19+ B cells, CD14+ monocytes, and granulocytes) different white blood cell types was performed using the default implementation of the *estimateCellCounts* function. Briefly, this function takes a user-supplied target 450K BeadChip dataset, combines that with the cell-sorted Reinius reference dataset available in the *FlowSorted.Blood.450k* Bioconductor package (v1.2.0) [Jaffe; Reinius et al, 2012] and quantile normalizes the combined data. The reference dataset has  $n$  observations from  $i=6$  subjects at each of  $j=6$  different separated cell types. Six hundred cell type informative CpG sites are chosen in the reference dataset, by comparing the mean methylation in a given cell type to the mean methylation of all five remaining cell types for CpG sites assayed. One hundred CpGs are chosen to distinguish each of the  $j$  cell types. These represent the 50 CpG sites with the greatest T statistic that were hypermethylated and the 50 CpG sites that were most hypomethylated compared to other cells. This results in a  $600 \times n$  matrix of CpG site methylation subset from the full reference dataset, called  $S_0$ . Within  $S_0$ , the relationship between indicators for each cell type and DNA methylation is then estimated, producing a vector of coefficients, called  $B_0$ , of length  $j$ . Individual level predictions of cell type proportion,  $\Gamma^*$ , are then fit in the corresponding user submitted dataset,  $S_1$ , using the coefficients estimated in  $S_0$  by the following equation:

$$\Gamma^* = (\tilde{B}_0^T \tilde{B}_0)^{-1} \tilde{B}_0^T S_1$$

Detailed description of the algorithm has previously been published [Houseman et al, 2012; Jaffe and Irizarry, 2014].

### Statistical analysis

All statistical analyses were performed using R statistical computing software (v3.1.0) [Team, 2013]. Differences in means by age and sex were assessed by Mann-Whitney U-Test. The linear relationships between cell type estimates by the two methods were determined using the Spearman correlation coefficient.

### Results

Estimates of white blood cell composition by the two different methods implemented, minfi and DCC, are summarized in Table I. The minfi method estimated the relative percentages of six white blood cell types (CD8+ T and CD4+ T lymphocytes, CD56+ natural killer cells, CD19+ B cells, CD14+ monocytes, and granulocytes) in samples from  $n=151$  newborns and again for 60 of the same children at age 12. Microscopic differential cell count (DCC) of these newborns ( $n=111$ ) and 12 year olds ( $n=45$ ) used banked available whole blood smear slides to count the frequency of five types of easily visually identifiable blood cells (lymphocytes, monocytes, neutrophils, eosinophils, and basophils) (see Methods for details).

### Cell composition estimates by age and sex

By the minfi method, the mean percentage estimates of all cell types except CD4+T lymphocytes were significantly different between newborn and 12 year old samples ( $p<0.01$ ) (Table I). Estimates of granulocytes represented the largest percentage of cell types in newborn samples (mean=55.0%), while lymphocytes (CD8+ T and CD4+ T lymphocytes, and natural killer cells) were noticeably less frequent (mean=37.2%,  $p<0.01$ ) (Figure 1A). In minfi estimates from 12 year olds, granulocyte and lymphocyte populations became much more comparable, with means 49.1% and 46.1% respectively ( $p=0.21$ ).

By the DCC method, the mean percentage of all but one cell type (eosinophil granulocytes) also differed significantly between newborns and 12 year olds (Table I). For newborns, mean DCC counts were 63.6% for granulocytes and 29.2% for lymphocytes ( $p<0.01$ ; Figure 1B). By 12 years of age, the gap in the frequency of cell types narrowed (46.9% and 46.8% respectively;  $p=0.57$ ).

In newborns, there was no difference in cell type distributions by sex ( $p>0.10$ ) by either minfi estimates ( $n=58$  girls,  $n=93$  boys; Figure 2A) or DCC direct analysis ( $n=58$  girls,  $n=53$  boys; Figure 2B). At age 12 only boys were sampled ( $N_{\text{minfi}}=60$ ;  $N_{\text{DCC}}=45$ ) so the comparison by sex was not possible.

### Comparison of cell composition estimates by minfi and DCC

Three cell type populations (Figure 3) were used for more direct comparison of the two methods of assessment of white blood cell composition. For the minfi method, the frequencies for CD8+ T, CD4+ T, natural killer cells, and B cells, were summed to give an estimate of lymphocytes. For DCC, proportions of neutrophils, eosinophils, and basophils were summed to give an estimate of granulocytes.

In samples from newborns, those estimates by minfi and DCC had poor linear correspondence with one another (Figure 3). In fact, the Spearman rank correlation coefficients calculated between minfi estimates and direct DCC analysis of monocytes, granulocytes and lymphocytes ranged from  $-0.01$  to  $-0.05$  and were not statistically significant (Figure 3). In scatterplots showing comparison between the two methods, the minfi method appeared to overestimate the proportion of lymphocytes (mean = 37.2%) relative to that by DCC (mean = 29.2%) and reference levels from newborns that range from 19–29% (Figure 3A–C) [Dallman, 1977; Nathan and Oski, 1981]. However, minfi also appears to overestimate the percent of monocytes (mean = 11.1%, reference = 5–7%) and gives a smaller percentage of granulocytes (minfi mean = 55.0%, DCC mean = 63.6%, reference = 32–83%). The minfi estimates also had greater variability in newborn samples, standard deviations (SDs) ranging from 2.3 to 8.8, compared to DCC estimates with SDs from 1.8 – 3.6.

However, the two estimates of cell counts were much more consistent in older children than in newborns. At 12 years of age, the means and standard deviations of all comparable cell populations, including granulocytes, lymphocytes and monocytes, were similar by both approaches (Figure 1). The Spearman correlation values for 12-year-old subjects ranged from 0.26 to 0.77 and were significant for granulocytes and lymphocytes (both  $p<0.001$ ),

and approached significance for monocytes ( $p=0.08$ ) (Figure 3). Scatterplots comparing the estimates by minfi and DCC at 12 years of age also showed the trend between methods to be linear with comparable amounts of variance.

## Discussion

Here, we present a detailed comparison of the leading method for estimating white blood cell composition, minfi, against results from a well-established clinical cell counting procedure, DCC, in samples from Mexican-American children at two time points to produce adjustment covariates in a whole blood EWAS analysis. While both methods yielded similar results in 12 year olds, minfi estimates and DCC were quite different in newborns. Our findings suggest that the algorithms applied by minfi may not be appropriate for cell type estimation in newborns and young children.

Longitudinally, as children aged from birth to 12 years old, we observed similar trends by both minfi and DCC. Each found that granulocyte levels were higher than lymphocytes at birth, but became comparable with one another by 12 years of age (Figure 1). This change corresponds with age-specific reference values, which demonstrate high levels of granulocytes (32–83%) relative to lymphocytes (19–29%) in newborns [Dallman, 1977; Nathan and Oski, 1981]. The levels of these two cell types vary noticeably postnatally and through early childhood, reaching a peak difference at 24 months, but stabilizing to adult levels around 9–12 years of age (28–48% for lymphocytes and 33–76% for granulocytes) [Dallman, 1977; Nathan and Oski, 1981].

Among older children, we found that minfi and DCC estimates were consistent with one another and mean estimates by both methods fell within published age-specific reference values [Dallman, 1977; Nathan and Oski, 1981]. However, minfi and DCC estimates differed greatly in newborns. In fact, we saw a negative correlation and linear trend across methods for each comparable cell type (Figure 3) suggesting that the algorithm implemented by minfi may have difficulty estimating cell composition in samples from newborns.

This deviation is likely explained by the reference dataset used in the minfi prediction model, which is derived from six middle-aged Swedish men. The composition of this reference dataset is crucial to minfi's performance: it is used to both identify CpGs differentially methylated by cell type and fit a regression model to those informative sites. The coefficients estimated in the reference data,  $B_0$ , establish the linear relationship between methylation at the informative sites and cell composition, which is used for prediction in the target dataset. A key assumption of this method is that the magnitude of the elements in  $B_0$  are consistent between the reference and target data. In many situations, this may not be an unreasonable assumption to make. Since the sites used to fit  $B_0$  are chosen by their association with cell type, one may expect them to perform cell type specific functions that would be consistent over time. However, given the poor performance of the minfi estimator in newborns, it seems likely that consistent effect of  $B_0$  in an adult reference does not hold for young children and impacts the estimator's accuracy. White blood cell populations are still maturing in the early post-natal period and are known to change greatly in relative abundance [Cheng et al, 2004; Dallman, 1977; Nathan and Oski, 1981]. Further, DNA



methylation is known to vary greatly over embryogenesis and may still be changing during early life [Guo et al, 2014; Smith et al, 2014]. Should the relationship between these two factors be inconsistent between early childhood and later life, this would result in biased minfi estimates.

Similar bias could occur if the consistent effect assumption does not hold across other biological host factors, such as gender or racial/ethnic ancestry. Both leukocyte populations and DNA methylation are known to vary by such factors [Adkins et al, 2011; Hsieh et al, 2007; Lim et al, 2010; McCarthy et al, 2014]. The current minfi reference data may be particularly susceptible to these forms of bias because all subjects are men of northern European descent. However, these biases are likely not as pronounced as those introduced by age in young children since the minfi estimates in CHAMACOS boys of Mexican ancestry are relatively accurate at age 12. Filtering out sites that vary by ethnicity or sex when fitting  $B_0$  could potentially reduce bias further, resulting in more accurate estimates of cell composition. Similarly, sites that vary by age could be excluded, an approach that has been used previously to identify candidate metastable epialleles [Harris et al, 2013]. However, given the lack of variation of the current Swedish male reference dataset over any of these potentially biasing factors, it is preferable to expand or generate a new reference that would have observations from early childhood, and that vary by race and ethnicity.

In conclusion, our comparison of the minfi method for estimating white blood cell composition against a cytological differential cell count demonstrates that minfi can robustly estimate cell populations in children as young as 12 years of age. However, minfi did not perform well in samples from newborns that are important targets of future EWAS because of interest in prenatal epigenetic changes due to exposure or physiological effect on future health. We hypothesize that this is due to low generalizability of the reference dataset currently used in the minfi estimation and suggest that improvement of this dataset would likely enhance its predictions in young children. We encourage using caution when applying the minfi method in populations that deviate substantially in white cell composition and/or methylation patterns from the current minfi reference data, such as by sex, racial/ethnic ancestry, and age in particular. Future work should explore further other factors, such as environmental exposures, that may also impact the validity of the minfi estimates.

## Acknowledgments

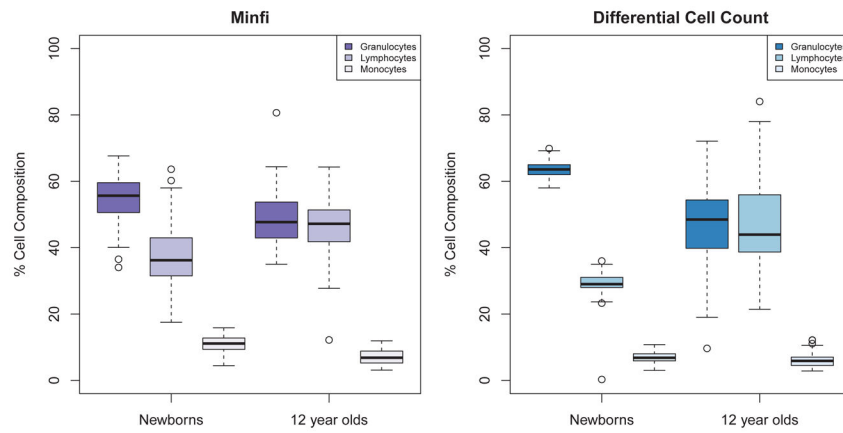
We are grateful to the laboratory and clinical staff and participants of the CHAMACOS study for their contributions. We thank Drs. Raul Aguilar Schall and Reuben Thomas for their helpful discussions regarding this work. We are also grateful to Michael Ha for technical assistance. This publication was made possible by grants RD83451301 from the U.S. Environmental Protection Agency (EPA) and PO1 ES009605 and R01ES021369 from the National Institute of Environmental Health Science (NIEHS).

## References

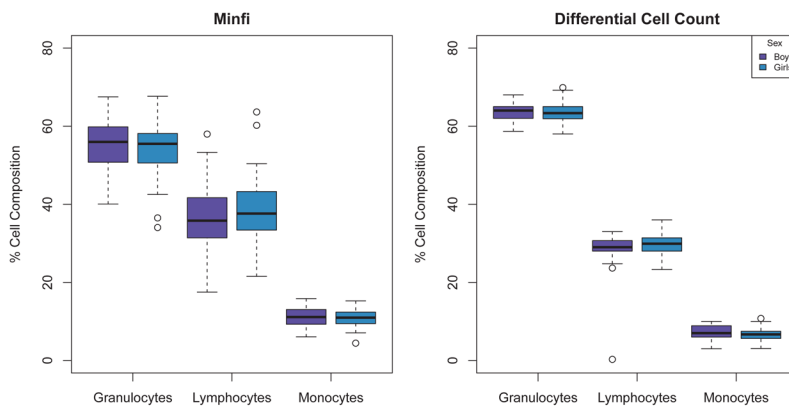
- Adkins RM, Krushkal J, Tylavsky FA, Thomas F. Racial differences in gene-specific DNA methylation levels are present at birth. *Birth Defects Res A Clin Mol Teratol.* 2011; 91(8):728–36. [PubMed: 21308978]
- Armstrong DA, Lesseur C, Conradt E, Lester BM, Marsit CJ. Global and gene-specific DNA methylation across multiple tissues in early infancy: implications for children's health research. *FASEB J.* 2014; 28(5):2088–97. [PubMed: 24478308]

- Aryee MJ, Jaffe AE, Corrada-Bravo H, Ladd-Acosta C, Feinberg AP, Hansen KD, Irizarry RA. Minfi: a flexible and comprehensive Bioconductor package for the analysis of Infinium DNA methylation microarrays. *Bioinformatics*. 2014; 30(10):1363–1369. [PubMed: 24478339]
- Babenko O, Kovalchuk I, Metz GA. Stress-induced Perinatal and Transgenerational Epigenetic Programming of Brain Development and Mental Health. *Neurosci Biobehav Rev*. 2014; 48C:70–91. [PubMed: 25464029]
- Barker DJ. In utero programming of chronic disease. *Clin Sci (Lond)*. 1998; 95(2):115–28. [PubMed: 9680492]
- Bibikova M, Barnes B, Tsan C, Ho V, Klotzle B, Le JM, Delano D, Zhang L, Schroth GP, Gunderson KL, et al. High density DNA methylation array with single CpG site resolution. *Genomics*. 2011; 98(4):288–95. [PubMed: 21839163]
- Bibikova M, Le J, Barnes B, Saedinia-Melnyk S, Zhou L, Shen R, Gunderson KL. Genome-wide DNA methylation profiling using InfiniumR assay. *Epigenomics*. 2009; 1(1):177–200. [PubMed: 22122642]
- Chadwick LH, Sawa A, Yang I, Baccarelli A, Breakefield X, Deng H, Dolinoy DC, Fallin MD, Holland N, Houseman EA, et al. New insights and updated guidelines for epigenome-wide association studies. *Neuroepigenetics*. 2014
- Cheng CK-W, Chan J, Cembrowski GS, van Assendelft OW. Complete blood count reference interval diagrams derived from NHANES III: stratification by age, sex, and race. *Laboratory hematology : official publication of the International Society for Laboratory Hematology*. 2004; 10(1):42–53. [PubMed: 15070217]
- Dallman, PR. Blood and blood-forming tissues. In: Rudolph, A., editor. *Paediatrics*. New York: Appleton-Century-Crofts; 1977. p. 1109–1112.
- Eskenazi B, Bradman A, Gladstone E, Jaramillo S, Birch K, Holland N. CHAMACOS, a longitudinal birth cohort study: lessons from the fields. *J Childrens Healt*. 2003; 1:3–27.
- Eskenazi B, Harley K, Bradman A, Weltzien E, Jewell NP, Barr DB, Furlong CE, Holland NT. Association of in utero organophosphate pesticide exposure and fetal growth and length of gestation in an agricultural population. *Environ Health Perspect*. 2004; 112(10):1116–24. [PubMed: 15238287]
- Essex MJ, Boyce WT, Hertzman C, Lam LL, Armstrong JM, Neumann SM, Kobor MS. Epigenetic vestiges of early developmental adversity: childhood stress exposure and DNA methylation in adolescence. *Child Dev*. 2013; 84(1):58–75. [PubMed: 21883162]
- Foley DL, Craig JM, Morley R, Olsson CJ, Dwyer T, Smith K, Saffery R. Prospects for epigenetic epidemiology. *Am J Epidemiol*. 2009; 169(4):389–400. [PubMed: 19139055]
- Guo H, Zhu P, Yan L, Li R, Hu B, Lian Y, Yan J, Ren X, Lin S, Li J, et al. The DNA methylation landscape of human early embryos. *Nature*. 2014; 511(7511):606–10. [PubMed: 25079557]
- Harris RA, Nagy-Szkal D, Kellermayer R. Human metastable epiallele candidates link to common disorders. *Epigenetics*. 2013; 8(2):157–163. [PubMed: 23321599]
- Holland N, Furlong C, Bastaki M, Richter R, Bradman A, Huen K, Beckman K, Eskenazi B. Paraoxonase polymorphisms, haplotypes, and enzyme activity in Latino mothers and newborns. *Environ Health Perspect*. 2006; 114(7):985–91. [PubMed: 16835048]
- Houseman EA, Accomando WP, Koestler DC, Christensen BC, Marsit CJ, Nelson HH, Wiencke JK, Kelsey KT. DNA methylation arrays as surrogate measures of cell mixture distribution. *BMC Bioinformatics*. 2012; 13:86. [PubMed: 22568884]
- Hsieh MM, Everhart JE, Byrd-Holt DD, Tisdale JF, Rodgers GP. Prevalence of neutropenia in the U.S. population: age, sex, smoking status, and ethnic differences. *Ann Intern Med*. 2007; 146(7):486–92. [PubMed: 17404350]
- Hughes V. Epigenetics: The sins of the father. *Nature*. 2014; 507(7490):22–4. [PubMed: 24598623]
- Jaffe AE. FlowSorted.Blood.450k: Illumina HumanMethylation data on sorted blood cell populations. R package version 1.2.0.
- Jaffe AE, Irizarry RA. Accounting for cellular heterogeneity is critical in epigenome-wide association studies. *Genome Biol*. 2014; 15(2):R31. [PubMed: 24495553]
- Joubert BR, Håberg SE, Nilsen RM, Wang X, Vollset SE, Murphy SK, Huang Z, Hoyo C, Midttun Ø, Cupul-Uicab LA, et al. 450K Epigenome-Wide Scan Identifies Differential DNA Methylation in

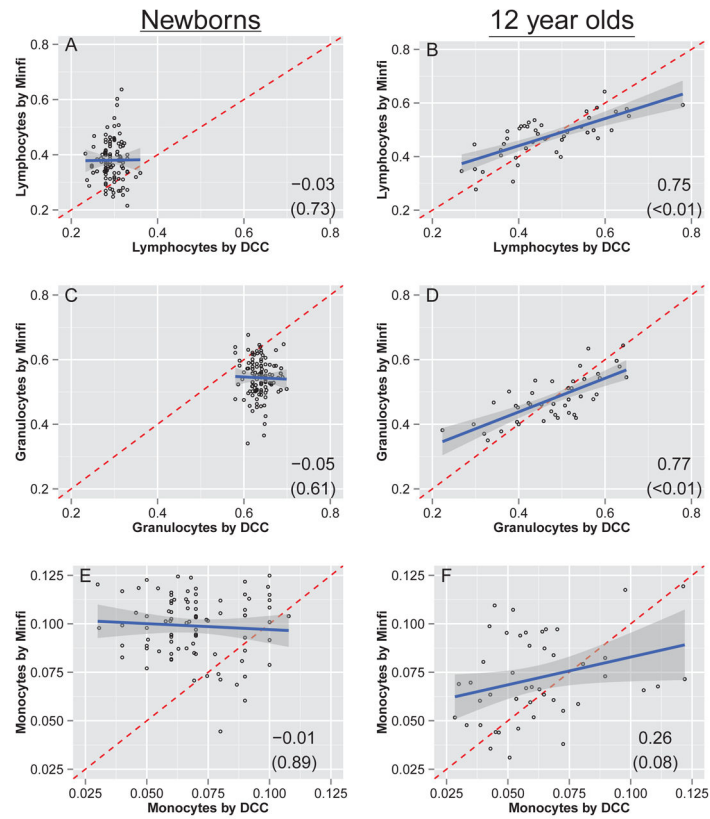
- Newborns Related to Maternal Smoking during Pregnancy. *Environmental Health Perspectives*. 2012; 120(10):1425–1431. [PubMed: 22851337]
- Koestler DC, Avissar-Whiting M, Houseman EA, Karagas MR, Marsit CJ. Differential DNA methylation in umbilical cord blood of infants exposed to low levels of arsenic in utero. *Environ Health Perspect*. 2013; 121(8):971–977. [PubMed: 23757598]
- Laird PW. Principles and challenges of genomewide DNA methylation analysis. *Nat Rev Genet*. 2010; 11(3):191–203. [PubMed: 20125086]
- Liang L, Cookson WOC. Grasping nettles: cellular heterogeneity and other confounders in epigenome-wide association studies. *Human Molecular Genetics*. 2014; 23(R1):ddu284–R88.
- Lim EM, Cembrowski G, Cembrowski M, Clarke G. Race-specific WBC and neutrophil count reference intervals. *Int J Lab Hematol*. 2010; 32(6 Pt 2):590–7. [PubMed: 20236184]
- Lister R, Pelizzola M, Downen RH, Hawkins RD, Hon G, Tonti-Filippini J, Nery JR, Lee L, Ye Z, Ngo QM, et al. Human DNA methylomes at base resolution show widespread epigenomic differences. *Nature*. 2009; 462(7271):315–22. [PubMed: 19829295]
- Liu Y, Aryee MJ, Padyukov L, Fallin MD, Hesselberg E, Runarsson A, Reinius L, Acevedo N, Taub M, Ronninger M, et al. Epigenome-wide association data implicate DNA methylation as an intermediary of genetic risk in rheumatoid arthritis. *Nat Biotechnol*. 2013; 31(2):142–147. [PubMed: 23334450]
- Lowe R, Rakyan VK. Correcting for cell-type composition bias in epigenome-wide association studies. *Genome Medicine*. 2014; 6(3):1–2. [PubMed: 24433494]
- McCarthy NS, Melton PE, Cadby G, Yazar S, Franchina M, Moses EK, Mackey DA, Hewitt AW. Meta-analysis of human methylation data for evidence of sex-specific autosomal patterns. *Bmc Genomics*. 2014:15. [PubMed: 24405808]
- Meissner A, Mikkelsen TS, Gu H, Wernig M, Hanna J, Sivachenko A, Zhang X, Bernstein BE, Nusbaum C, Jaffe DB, et al. Genome-scale DNA methylation maps of pluripotent and differentiated cells. *Nature*. 2008; 454(7205):766–770. [PubMed: 18600261]
- Michels KB, Binder AM, Dedeurwaerder S, Epstein CB, Grealley JM, Gut I, Houseman EA, Izzi B, Kelsey KT, Meissner A, et al. Recommendations for the design and analysis of epigenome-wide association studies. *Nat Methods*. 2013; 10(10):949–955. [PubMed: 24076989]
- Nathan, DG.; Oski, FA. *Hematology of Infancy and Childhood*. 2. 1981. p. 1552-1574.
- Perera F, Herbstman J. Prenatal environmental exposures, epigenetics, and disease. *Reprod Toxicol*. 2011; 31(3):363–73. [PubMed: 21256208]
- Reinius LE, Acevedo N, Joerink M, Pershagen G, Dahlen SE, Greco D, Soderhall C, Scheynius A, Kere J. Differential DNA methylation in purified human blood cells: implications for cell lineage and studies on disease susceptibility. *PLoS One*. 2012; 7(7):e41361. [PubMed: 22848472]
- Sandoval J, Heyn H, Moran S, Serra-Musach J, Pujana MA, Bibikova M, Esteller M. Validation of a DNA methylation microarray for 450,000 CpG sites in the human genome. *Epigenetics*. 2011; 6(6):692–702. [PubMed: 21593595]
- Smith ZD, Chan MM, Humm KC, Karnik R, Mekhoubad S, Regev A, Eggan K, Meissner A. DNA methylation dynamics of the human preimplantation embryo. *Nature*. 2014; 511(7511):611–5. [PubMed: 25079558]
- Team RC. R: A language and environment for statistical computing. R Foundation for Statistical Computing; Vienna, Austria: Vienna, Austria: 2013. <http://www.R-project.org/>
- Teschendorff AE, Menon U, Gentry-Maharaj A, Ramus SJ, Gayther SA, Apostolidou S, Jones A, Lechner M, Beck S, Jacobs IJ, et al. An Epigenetic Signature in Peripheral Blood Predicts Active Ovarian Cancer. *PLoS One*. 2009; 4(12)
- Turgeon, ML. *Clinical Hematology Theory and Procedures*. 5. 2011. p. 40-45.
- Weber M, Davies JJ, Wittig D, Oakeley EJ, Haase M, Lam WL, Schübeler D. Chromosome-wide and promoter-specific analyses identify sites of differential DNA methylation in normal and transformed human cells. *Nat Genet*. 2005; 37(8):853–862. [PubMed: 16007088]



**Figure 1.** Box plots of percent cell composition estimated by minfi and differential cell count (DCC) in samples from newborns and 12 year olds. The minfi estimates are taken from  $n=151$  newborns and  $n=60$  12 year old boys, and have summed estimates of CD8+ T, CD4+ T, natural killer cells, and B cells into a single category of lymphocytes for comparison. The DCC estimates are taken from  $n=111$  newborns and  $n=45$  12 year olds, and have summed proportions of neutrophils, eosinophils, and basophils into category of granulocytes.



**Figure 2.** Box plots of percent cell composition estimated by minfi and differential cell count (DCC) for girls and boys in samples from newborns. The minfi estimates are taken from n=151 newborns (n=58 girls and n=93 boys). They have summed estimates of CD8+ T, CD4+ T, natural killer cells, and B cells into a single category of lymphocytes for comparison. The DCC estimates are taken from n=111 newborns (n=58 girls and n=53 boys). They have summed proportions of neutrophils, eosinophils, and basophils into category of granulocytes.



**Figure 3.**

Scatter plot of cell type percentages by minfi and differential cell count (DCC) methods in cord samples for lymphocytes (A), granulocytes (C), and monocytes (E). Also, plots of 12 year samples for lymphocytes (B), granulocytes (D), and monocytes (F). Estimate of linear trend by regression shown in blue with 95% confidence interval in gray. Exact linear correlation, slope=1 and intercept=0, shown in dotted red for reference. Spearman rank correlation,  $\rho$ (P-value), shown in bottom right corner for each comparison.

Summary of white blood cell type percentage estimates by two methods, differential cell count (DCC) and minfi in newborn and 12 year old child blood samples.

**Table I**

	Newborns			12 year olds			P Value
	N	Mean, %	Min Max	N	Mean, %	Min Max	
<b>Minfi</b>							
CD8+ T	151	7.1	1.6 13.5	60	14.3	4.6 24.7	<0.01*
CD4+ T	151	13.7	3.5 25.6	60	13.9	2.4 27.1	0.90
NK cells	151	4.2	0 17.4	60	8.5	0 21.4	<0.01*
B Cells	151	12.2	1.9 24.1	60	9.4	5.1 16.7	<0.01*
Monocytes	151	11.1	4.4 15.8	60	7.1	3.1 11.9	<0.01*
Granulocytes	151	55.0	34.1 67.7	60	49.1	35 80.7	<0.01*
<b>DCC</b>							
Lymphocytes	111	29.2	0.3 36	45	46.8	21.4 84.1	<0.01*
Monocytes	111	6.9	3 10.8	45	6.2	2.8 12.2	<0.01*
Neutrophils	111	60.4	54 67.2	45	42.7	9.7 68.5	<0.01*
Eosinophils	111	2.9	0.5 5.3	45	4.1	0 15	0.20
Basophils	111	0.2	0 2	45	0.1	0 3.9	0.04*

\* Cell types with mean percentages significantly different ( $p < 0.05$ ) between ages by Mann-Whitney U test.