# Data processing choices can affect findings in differential methylation analyses: an investigation using data from the LIMIT RCT

Jennie Louise [Corresp., 1, 2] , Andrea R Deussen [3] , Jodie M Dodd [3, 4]

[1] Discipline of Obstetrics & Gynaecology and The Robinson Research Institute, The University of Adelaide, Adelaide, Australia

[2] Adelaide Health Technology Asseessment, The University of Adelaide, Adelaide, Australia

[3] Discipline of Obstetrics & Gynaecology and The Robinson Research Institute, The University of Adelaide, Adelaide, South Australia, Australia

[4] Department of Perinatal Medicine, Women's and Babies Division, The Women's and Children's Hospital, Adelaide, South Australia, Australia

Corresponding Author: Jennie Louise
Email address: jennie.louise@adelaide.edu.au

Objective: A wide array of methods exist for processing and analysing DNA methylation data. We aimed to perform a systematic comparison of the behaviour of these methods, using cord blood DNAm from the LIMIT RCT, in relation to detecting hypothesised effects of interest (intervention and pre-pregnancy maternal BMI) as well as effects known to be spurious, and known to be present. Methods: DNAm data, from 645 cord blood samples analysed using Illumina 450K BeadChip arrays, were normalised using three different methods (with probe filtering undertaken pre- or post- normalisation). Batch effects were handled with a supervised algorithm, an unsupervised algorithm, or adjustment in the analysis model. Analysis was undertaken with and without adjustment for estimated cell type proportions. The effects estimated included intervention and BMI (effects of interest in the original study), infant sex and randomly assigned groups. Data processing and analysis methods were compared in relation to number and identity of differentially methylated probes, rankings of probes by p value and log-fold-change, and distributions of p values and log-fold-change estimates. Results: There were differences corresponding to each of the processing and analysis choices. Importantly, some combinations of data processing choices resulted in a substantial number of spurious 'significant' findings. We recommend greater emphasis on replication and greater use of sensitivity analyses.

1 # Data Processing Choices Can Affect Findings in Differential
2 # Methylation Analyses: An Investigation Using Data from the
3 # LIMIT RCT

4 **Authors**
5 Jennie Louise[1,2]
6 Andrea R Deussen[1]
7 Jodie M Dodd[1,3]
8

9 **Affiliations**
10 1.  Discipline of Obstetrics & Gynaecology and The Robinson Research Institute, The University of
11 Adelaide, Adelaide, South Australia, AUSTRALIA
12 2. Adelaide Health Technology Asseessment, The University of Adelaide, Adelaide, South Australia,
13 AUSTRALIA
14 3. Department of Perinatal Medicine, Women's and Babies Division, The Women's and Children's
15 Hospital, Adelaide, South Australia, AUSTRALIA.
16

17 **Corresponding Author**
18 Dr Jennie Louise
19 The University of Adelaide
20 Women's and Children's Hospital
21 72 King William Rd
22 North Adelaide, South Australia, AUSTRALIA 5006
23 Email: jennie.louise@adelaide.edu.au
24

25

## Abstract

**Objective:** A wide array of methods exist for processing and analysing DNA methylation data. We aimed to perform a systematic comparison of the behaviour of these methods, using cord blood DNAm from the LIMIT RCT, in relation to detecting hypothesised effects of interest (intervention and pre-pregnancy maternal BMI) as well as effects known to be spurious, and known to be present.

**Methods:** DNAm data, from 645 cord blood samples analysed using Illumina 450K BeadChip arrays, were normalised using three different methods (with probe filtering undertaken pre- or post-normalisation). Batch effects were handled with a supervised algorithm, an unsupervised algorithm, or adjustment in the analysis model. Analysis was undertaken with and without adjustment for estimated cell type proportions. The effects estimated included intervention and BMI (effects of interest in the original study), infant sex and randomly assigned groups. Data processing and analysis methods were compared in relation to number and identity of differentially methylated probes, rankings of probes by p value and log-fold-change, and distributions of p values and log-fold-change estimates.

**Results:** There were differences corresponding to each of the processing and analysis choices. Importantly, some combinations of data processing choices resulted in a substantial number of spurious 'significant' findings. We recommend greater emphasis on replication and greater use of sensitivity analyses.

**Clinical Trials Registration:** ACTRN12607000161426



**Word Count:** 6298


## Acknowledgements

## Introduction and Background

With the advent of high-throughput assays, epigenome-wide DNA methylation studies have become more popular, and researchers are now investigating the effects on DNA methylation (DNAm) of a wide range of environmental exposures and physiological conditions, with particular interest in the contribution of epigenetic mechanisms such as DNAm to the early life origins of health and disease. The ability to perform EWAS is particularly useful in relation to conditions where associated differences in DNAm are likely to be fairly modest (Marabita et al., 2013). However, DNAm data – as with high-dimensional 'omics' data generally – requires substantial pre-processing prior to analysis, including probe and sample filtering, normalisation to remove variation due to technological factors, and correction for other factors which may confound effects of interest, such as batch effects or differences in cell type proportions between samples. Numerous methods exist to perform these processing steps, and many articles have been published which provide useful guidance for the use of analysis pipelines(Marabita et al., 2013; Yousefi et al., 2013; Lehne et al., 2015; Morris & Beck, 2015; Maksimovic, Phipson & Oshlack, 2017), or comparing some alternatives for individual steps in the overall processing pipeline, including probe filtering (Heiss & Just, 2019), normalisation (Wang et al., 2012, 2015; Fortin et al., 2014; Wu et al., 2014; Hicks & Irizarry, 2015), or correction for / avoidance of batch effects. These have led to some general conclusions regarding the need to account for batch effects, the importance of correcting for estimated cell type proportion, and perhaps the greater suitability of within-array normalisation methods compared to between-array methods when global methlation differences are not expected (Maksimovic, Phipson & Oshlack, 2017), but there is no clear overall consensus on the best processing or analysis approach (Price & Robinson, 2018; Zindler et al., 2020), or of the overall advantages and disadvantages of different combinations of processing choices.

We recently investigated the effect of an antenatal diet and lifestyle intervention, and of maternal early pregnancy BMI, on neonatal cord blood DNA methylation in infants of mothers who were overweight or obese in early pregnancy(Louise et al., 2022). We were unable to replicate findings from previous studies  which reported a range of loci to be significantly differentially methylated in relation to maternal BMI or diet and lifestyle in pregnancy (Gemma et al., 2009; Sharp et al., 2015, 2017; Thakali et al., 2017; Hjort et al., 2018)and indeed did not find any significant differences In methylation corresponding to either BMI or intervention effects. We were aware of literature suggesting that use of supervised batch-correction algorithms may produce spurious findings (Nygaard, Rødland & Hovig, 2016a; Price & Robinson, 2018; Zindler et al., 2020), and that the number of statistically significant findings may differ according to normalisation method (Wu et al., 2014), adjustment for estimated cell type proportion (Sharp et al., 2017) or stringency of type I error control (Wu et al., 2014), which led to the hypothesis that the discrepancy in findings may be due in part to differences in data processing and analysis methods.  Following the common practice in clinical trials of conducting sensitivity analyses to assess robustness of results to various assumptions and decisions (Thabane et al., 2013), we performed several re-analyses with different normalisation methods, methods for batch effect handling, type I error control, and presence vs absence of adjustment for estimated cell type proportions. This confirmed that our findings indeed differed under different data-processing and analysis choices.

While previous studies comparing different methods have also produced different findings, these have tended to consider only one element of the processing and analysis pipeline (such as normalisation, or batch correction) in isolation. In addition, they have tended to concentrate on the tendency for some methods to produce results which are likely to be spurious (false positives), while often being unable to

102    definitively confirm this due to the lack of known truth regarding the presence and magnitude of
103    differential methylation effects.
104
105    We therefore set out to investigate the impact of different data-processing choices in a more systematic
106    way, looking at the effect of combinations of data processing choices on findings specifically regarding
107    statistically significant differentially methylated probes (DMPs), and of the behaviour of these
108    combinations in cases where effects are known to be either absent or present, as well as their behaviour
109    in relation to our effects of interest (maternal BMI and lifestyle intervention).  We were able to create a
110    scenario in which effects were known to be absent by randomly assigning samples to groupings.  We
111    could not similarly ensure a scenario where effects were known to be present (as the truth regarding the
112    existence, location and magnitude of any effects in our samples is not known); however we investigated
113    effects of infant sex as a rough proxy, since infant sex is known to have substantial effects on DNAm
114    which can be detected by the 450K array (Yousefi et al., 2015).
115

## Data and Methods

### The LIMIT Randomised Controlled Trial

118
119    The LIMIT study was a randomised, controlled trial of an antenatal diet and lifestyle intervention for
120    women with early pregnancy BMI $\geq$25.0 kg/m$^2$.  The study, and its primary and main secondary
121    outcomes, have been extensively reported elsewhere (Dodd et al., 2014b). Inclusion criteria were early
122    pregnancy BMI $\geq$25.0 kg/m$^2$ and pregnancy between $10^{+0}$ and $20^{+0}$ weeks' gestation, with exclusion
123    criteria of multiple gestation or previously existing diabetes.  The study randomised 2212 women in total
124    to one of two groups: a comprehensive diet and lifestyle intervention (Lifestyle Advice; n=1108) or
125    antenatal care delivered according to local guidelines (Standard Care; n=1104) which did not include
126    information on diet or physical activity.  The study was reviewed by the ethics committee of each
127    participating institution including the Women's and Children's Health Network Human Research Ethics
128    Committee (1839 & 2051); the Central and Northern Adelaide Health Network Human Research Ethics
129    Committee (2008033) and the Southern Adelaide Local Health Network Human Research Ethics
130    Committee (128/08).  Informed written consent was obtained for all participants to participate in the
131    LIMIT study, and additional written consent was obtained to collect samples of umbilical cord blood at
132    delivery for the purposes of gene expression research related to weight and to the diet and lifestyle
133    intervention.
134
135    The primary outcome of the LIMIT study was birth of a Large for Gestational Age (LGA) infant. There
136    were no significant differences observed between the groups in relation to this outcome; however, a
137    significantly lower incidence of high birthweight (>4kg) was observed in the Lifestyle Advice group, with
138    a Relative Risk of 0.82 (95% CI: 0.68, 0.99, p=0.04). Additionally, measures of diet quality and physical
139    activity were improved in women in the Lifestyle Advice group compared to those in the Standard Care
140    group (Dodd et al., 2014a).
141
142    As previously outlined in the companion paper (Louise et al., 2022), Cord Blood DNA for a range of
143    secondary studies was collected at the time of birth from consenting participants, and was frozen as
144    whole blood preserved with EDTA. Funding was available to perform DNA methylation analysis for a
145    total of 649 samples, which were randomly selected from the total number of available samples,
146    balanced between the Lifestyle Advice and Standard Care groups (Supplementary Table S1).  After DNA

147    extraction, genome-wide DNA methylation was performed using the Illumina Infinium

148    HumanMethylation 450K Bead-Chip array. Results were supplied as raw probe intensities (IDAT files).

149

150    For the additional analyses investigating known spurious effects, artificial (fake) groups were created by

151    assigning samples based on random draws from binomial distributions.  The first grouping

152    ('Tortoiseshell' vs'Tabby') was generated using a binomial distribution with 50% probability of

153    assignment to each group. The second grouping ('Long'- vs 'Short-Haired') was created to mimic

154    stratified randomisation as well as unequal proportions in each group: within each level of the first fake

155    group, samples were assigned to Long-Haired with 40% probability and Short-Haired with 60%

156    probability.

157

158    All data processing and analyses were undertaken using R version 4.0 (R Core Team, 2018).

159

## Probe and Sample Filtering

160

161

162    The *minfi* package (Aryee et al., 2014) was used to read in the raw *idats* (without normalisation), and to

163    calculate both probe-wise and sample-wise 'detection p values'. Samples were identified as 'faulty' if

164    they had a detection p-value $\geq$0.05. 13 such samples were excluded; however these were due to a

165    known chip failure, and had subsequently been rerun. A further four samples were excluded because

166    the correct corresponding study identifier could not be ascertained, leaving 645 samples for analysis.

167

168    Probes were filtered using multiple criteria. Firstly, probes were excluded if they had a detection p-value

169    $\geq$0.001 in more than 25% of the 645 samples, indicating that their signal could not be accurately

170    detected for a large proportion of samples (Dedeurwaerder et al., 2014; Maksimovic, Phipson & Oshlack,

171    2017). Secondly, probes were excluded if they were on a list of those previously identified as cross-

172    reactive (Chen et al., 2013); i.e. there was a high probability they may hybridize to locations on the

173    genome different to those for which the probe was designed (Dedeurwaerder et al., 2014; Naeem et al.,

174    2014). Thirdly, probes with an identified SNP within 3 nucleotides of the CpG site and minor allele

175    frequency >1%, and probes on the X and Y chromosomes were excluded. This was done in order to avoid

176    spurious methylation 'differences' due either to SNPs within the CpG targets, or due to X and Y

177    chromosomes (Dedeurwaerder et al., 2014; Naeem et al., 2014). Filtering of cross-reactive probes,

178    probes with a nearby SNP and probes on the X and Y chromosomes was performed using the *DMRCate*

179    package (Peters et al., 2015). This left 426,572 probes available for analysis.

180

181    Probe filtering was performed either after normalisation (post-filtered) or prior to normalisation (pre-

182    filtered). The one exception for pre-filtering was when normalising using the BMIQ method, where

183    probes on the X and Y chromosome were retained as this was required in order for the function to run.

184

185

## Normalisation

186

187

188    Normalisation involves making changes to the raw data in order to remove artifactual variation. In the

189    case of Illumina 450K BeadChip arrays, this requires correcting for the presence of two different probe

190    types. Infinium I probes use the same colour signal for methylated and unmethylated CpG and are often

191    used for regions of high CpG density, while Infinium II probes use different colours to differentiate

192    between methylated and unmethylated states (Pidsley et al., 2013; Wang et al., 2015).  Normalisation is

193    performed on $\beta$ values, or the ratio of methylated to total intensity, defined as $\frac{M}{M + U + offset}$ . Here, M

194    is the methylated intensity and U is the unmethylated intensity; the offset is a constant added to
195    regularize the $\beta$ value where both methylated and unmethylated intensities are low.  The distribution of
196    $\beta$ values  is bimodal, with peaks corresponding to methylated and unmethylated states, but the
197    distribution of Infinium II probes differs from that of Infinium I, being more compressed towards 0.5
198    (Pidsley et al., 2013) and hence having a smaller 'dynamic range' (Teschendorff et al., 2013;
199    Dedeurwaerder et al., 2014).
200
201    Numerous methods exist for normalising Illumina BeadChip array data, but there is little consensus or
202    guidance on which should be employed in a given context. The main advice is that 'between-array'
203    methods, which normalise across samples, are preferable when global differences between samples are
204    expected, while 'within-array' methods, which normalise probes within each sample, are better suited
205    to effects in which the majority of genes will not be differentially expressed. (Maksimovic, Phipson &
206    Oshlack, 2017) The latter is the context in which many EWAS studies, including the present one, are
207    conducted; as noted above, only modest differences in a small proportion of genes are expected for
208    most early-life exposures. The methods chosen for the present investigation have all been used in the
209    context of studies such as this: Categorical-Subset Quantile Normalisation (SQN) (Wu et al., 2014; Wang
210    et al., 2015), Beta-Mixture Quantile Normalisation (BMIQ) (Teschendorff et al., 2013), and Subset-
211    Quantile Within-Array Normalisation (SWAN) (Maksimovic, Gordon & Oshlack, 2012).  While numerous
212    other methods exist, a comparison of all available normalisation methods was beyond the scope of this
213    paper. Further details on the  methods are given in the Supplementary Information.
214
215    Both Subset Quantile Normalisation and Subset-Within-Array-Normalisation were performed using
216    functions in the *minfi* package (*preprocessQuantile* and *preprocessSWAN* respectively), on raw intensity
217    data. Beta-Mixture Quantile normalisation was performed using the *champ.norm* function in the *ChAMP*
218    package after converting intensities to $\beta$ values.
219

## Batch Effects

221
222    Batch effects arise when samples are processed in separate groups, creating unwanted variation due,
223    for example, to different reagents, different plates or different scanner settings. (Morris & Beck, 2015;
224    Nygaard, Rødland & Hovig, 2016a; Price & Robinson, 2018)
225    There are 12 Illumina 450K arrays (samples) per chip (this is reduced to 8 arrays per chip for the more
226    recent 850K array); thus most studies involving large numbers of samples must be run on multiple chips.
227    This introduces extra variability to the data, and may also confound the actual effects of interest, if
228    samples from different groups are not evenly distributed between the batches. These effects must be
229    accounted for in order to obtain valid estimates of the effects of interest.
230
231    Unlike probe filtering and normalisation, batch effects can be handled at the analysis stage, by adjusting
232    for batch in the analysis model.  However, it is also common to address batch effects at the data-
233    processing stage, using a batch-correction algorithm, with the resulting data considered to be free of
234    batch effects (Nygaard, Rødland & Hovig, 2016a). The ComBat algorithm has been widely used and
235    considered the most effective method (Zindler et al., 2020) of removing batch effects in DNAm data; it
236    has been incorporated into various analysis pipelines. Until recently, ComBat could be implemented only
237    as a supervised method, in which the biological factors of interest had to be specified along with the
238    batch variable (Price & Robinson, 2018) (Fortin, Triche & Hansen, 2016); it can now also be implemented
239    as an unsupervised method, in which only the batch variable is specified.
240

241  For each of the normalised datasets (i.e. SQN, BMIQ and SWAN normalised datasets, each with probes
242  filtered either before or after normalisation), we handled batch effects in three ways: firstly, by
243  adjusting for batch in the analysis model (BatchAdjust); secondly, implementing the supervised ComBat
244  algorithm (sCB); and thirdly, implementing the unsupervised ComBat algorithm (uCB).  For the
245  supervised ComBat algorithm, it was necessary to run the process twice: once with the effects of
246  interest specified as maternal early pregnancy BMI, antenatal intervention group, and their interaction;
247  and again with the effects of interest specified as Fake Group 1, Fake Group 2, their interaction, and
248  infant sex.
249
250  ## Cell Type Proportions
251  Cord blood, like whole blood, contains a mixture of different cell types, which have different DNA
252  methylation profiles.(Jaffe & Irizarry, 2014; Teschendorff & Zheng, 2017) If samples differ in the
253  proportion of these different cell types, this may confound effects of interest, either hiding true
254  differences in DNAm, or giving rise to spurious differences. Most studies of the effect of BMI, lifestyle
255  interventions, or similar factors on cord blood DNA methylation have not attempted (or have not
256  documented an attempt) to correct for potential differences in cell type composition, perhaps because
257  reference profiles for cord blood were not available until more recently (Bakulski et al., 2016), and the
258  mix of cell types and DNAm profiles may differ in cord blood compared to whole blood, making it
259  inappropriate to apply reference profiles from whole blood to cord blood data.(Cardenas et al., 2016)
260
261  We estimated the proportion of B cells, CD4+T, CD8+T, granulocytes, monocytes, natural killer, and
262  nucleated RBCs in the raw data using the *estimateCellCounts()* function in the *minfi* package, with the
263  Cord Blood reference panel. The estimated proportions were then added to the metadata for use as
264  adjustment variables in the analyses. We then undertook analyses either adjusted or not adjusted for
265  estimated cell type proportion.
266
267  Figure 1 depicts the combinations of data-processing and analysis choices that were undertaken. In
268  brief, there were six normalised datasets (three different normalisation methods, with probe filtering
269  performed before normalisation or after normalisation). These datasets were either used immediately
270  for analysis, or processed using the ComBat algorithm (in both supervised and unsupervised form) prior
271  to analysis.  Non-ComBat-processed data were analysed with three different models: an unadjusted
272  model (containing only the effects of interest), a model adjusted for batch, and a model adjusted for
273  batch and estimated cell type proportions.  ComBat-processed data were analysed with two different
274  models: one containing no other adjustment variables (but assumed to be 'pre-adjusted' for batch), and
275  one adjusted for cell type proportion.
276
277  ## Statistical Analysis
278
279  Differential methylation was investigated probe-wise using linear models with empirical Bayes variance
280  correction as implemented in the *limma* package (Ritchie et al., 2015; Smyth).  For effects of BMI and
281  intervention, models specified BMI (as continuous and mean-centred), intervention (Lifestyle Advice vs
282  Standard Care), and their interaction. Contrasts were specified to estimate the effect of the intervention
283  and of maternal BMI.  Because of the presence of the intervention-by-BMI interaction term, this
284  required specification of the BMI values at which the intervention effect was to be estimated, and the
285  intervention groups in which the effect of BMI was to be estimated.  For estimating intervention effects,
286  we chose the mean BMI of the cohort (i.e. value of 0 for the mean-centred variable, corresponding to an
287  actual BMI of approximately 33 kg/m$^2$), and at 5 kg/m$^2$ above the mean (a value of 5, corresponding to

288    an actual BMI of approximately 38 kg/m$^2$.  For the effect of BMI, we estimated the effect of an increase
289    of 5 kg/m$^2$ in BMI in each of the intervention groups (Standard Care, Lifestyle Advice) respectively. For
290    effects of fake groups and infant sex, the models specified sex (Female vs Male), Fake Group 1
291    (Tortoiseshell vs Tabby), Fake Group 2 (Long-Haired vs Short-Haired), and their interaction. Contrasts
292    were specified for infant sex, and for the effect of each fake group separately within levels of the other
293    fake group (i.e. effect of Tortoiseshell in Long-Haired and in Short-Haired; and effect of Short-Haired in
294    Tortoiseshell and Tabby). The model matrix and contrast matrices are shown in the Supplementary
295    Table S4.
296
297    For each contrast in each model, the number and identity (where applicable) of any differentially
298    methylated probes (DMPs) were obtained. For detection of DMPs, *limma*'s default method of multiple-
299    comparisons correction (Benjamini-Hochberg) and default alpha of 0.05 was used; this method controls
300    the False Discovery Rate, or the proportion of statistically significant results not corresponding to true
301    effects.  Where DMPs were obtained, a comparison was made using the Holm method (retaining alpha
302    of 0.05), which controls the Family-Wise-Error Rate (the probability that at least one statistically
303    significant result does not correspond to a true effect).  The Holm method can be considered more
304    stringent than Benjamini-Hochberg, but is less stringent than Bonferroni correction, which is known to
305    be too conservative even outside the context of high-dimensional data and is therefore not generally
306    appropriate for EWAS studies (and has not been used in other studies investigating cord blood DNAm in
307    relation to maternal BMI or diet and lifestyle).  The full set of p values and estimated log-fold-changes
308    (for all 426572 probes) corresponding to each contrast were also obtained, in order to compare probe
309    rankings and overall distributions. To make the comparison more tractable, probe rankings were
310    investigated using only those probes ranked in the top 10 (i.e. the probes with the smallest p value, or
311    largest estimated logFC, in a given model).
312
313    The findings from different data-processing choices were then compared along five dimensions:
314        1.  Number and identity of differentially methylated probes (DMPs); for infant sex, the direction of
315            differential methylation ('down', corresponding to negative t-statistics or lower methylation in
316            females, versus 'up', corresponding to higher methylation in females) was also examined.  For
317            BMI and intervention effects, the number and identity of statistically significant DMPs allows us
318            to see differences in detection of effects, and whether the different analysis pipelines produce
319            consistent results regarding the identity of any DMPs, though the truth is not known.  For the
320            fake groupings, the number of statistically significant DMPs is an indication of the tendency to
321            produce spurious findings.  For infant sex, while we do not know the actual number and identity
322            of truly differentially methylated sites, differences in the number and identity of DMPs
323            demonstrate that there must be either false positives or false negatives.
324        2.  Consistency of rankings by p value for 'top 10' probes.  This gives an indication of whether
325            different methods will give the same results for the probes with the largest differences.
326        3.  Consistency of rankings by logFC for 'top 10' probes, as well as the consistency of the logFC
327            estimates. This gives an indication of whether the estimates of effect are similar between
328            methods.
329        4.  Overall distribution of p values. Under the null hypothesis of no effect, p values should have a
330            uniform distribution between 0 and 1; if effects are present, there will be more p values at the
331            lower end of the distribution, the extent of which will depend on how many DMPs there are and
332            the strength of the effects.
333        5.  Overall distribution of logFC estimates. Under the null hypothesis of no effect, logFC estimates
334            should be roughly normally distributed around 0.  If effects are present, there will be more
335            estimates far away from 0 (with the direction depending on whether the effect is one of

336        hypomethylation or hypermethylation, and the distance depending on the strength of the
337        effect).

338

## Results

340   All dimensions of data processing choices had some impact on downstream analysis results, in terms of
341   the number (and identity) of differentially methylated probes, rankings of probes (by p value and logFC),
342   estimates of logFC, and overall distribution of p values and logFC, corresponding to both real and
343   spurious effects of interest. In some cases a consistent impact of a particular choice was observed, while
344   in others there was no consistent pattern, or this pattern varied according to the other choices with
345   which it was combined.

346

347   Tables 1-3 give information about the number of significantly differentially methylated probes in each of
348   the models fitted for the combinations of filtering, normalisation, batch correction and cell adjustment
349   approaches, for infant sex, maternal BMI (in the Standard Care group) and 'Tortoiseshell' (in the 'Tabby'
350   group) respectively. Supplementary Figures S1-S4 show the degree of overlap in the actual probes found
351   to be significantly differently methylated between models for infant sex, intervention (at the mean BMI
352   of the cohort), BMI (in the Standard Care group), and the effect of 'Short-Haired' in the 'Tabby' group.
353   Figures 2 and Figures 3-5 show the differences in ranking of probes (those which were in the top 10 in
354   any model) by p value and log-Fold-Change for the same set of effects, and Supplementary Tables S4-S6
355   gives Spearman Rank Correlation matrices for these rankings.  The overall distribution of p values, and of
356   log-Fold-Change estimates, for the same set of effects is shown in Figures 6 and 7.

357

358   Below we discuss the effect of each dimension (probe filtering, normalisation, batch effects, cell type
359   correction) on results.

360

361

### Effect of probe filtering pre-normalisation vs post-normalisation

363

364   Filtering probes prior to normalisation, compared to filtering after normalisation, led to modest
365   differences in number of DMPs, rankings of probes by logFC and p value, and overall distributions of p
366   values and logFC estimates. Filtering pre-normalisation produced different numbers of DMPs for infant
367   sex, but the nature of the effect differed by normalisation method: in SWAN data there was a consistent
368   pattern of fewer significant probes both negative and positive, while in BMIQ data there were fewer
369   negative but more positive probes, and in SQN data there were more negative but fewer positive
370   probes. In relation to effects of BMI, intervention, and fake groups, differences were harder to discern
371   due to the lack of any DMPs for many models; however, when DMPs were present for an effect, there
372   was a tendency for there to be a greater number of them in the pre-filtered data.

373

374   Probe rankings, by logFC and p value, tended to be relatively consistent between pre-filtered and post-
375   filtered data, with some cases of larger discrepancies in rankings for individual probes. The discrepancies
376   were more common, and larger, for fake group, BMI and intervention effects than for infant sex.
377   Similarly, there were no dramatic differences in distributions of p values or logFC estimates for infant
378   sex; there were differences in distribution between pre- and post-filtered data for fake group,
379   intervention and BMI effects, but there was no consistent pattern to these differences.

380

381    The question of whether probe filtering should be carried out before or after normalisation is one which
382    has received surprisingly little attention in the literature, but our results suggest that it can make a
383    difference to findings in some contexts.  In particular there may be a higher risk of spurious findings in
384    pre-filtered data, but there may also be a risk of failing to detect true differences – either any
385    differences, or specifically hypomethylated or hypermethylated loci, depending upon the normalisation
386    method employed.
387
388    ## Effect of Normalisation Method
389    Normalisation method had a substantial influence on number and identity of DMPs, rankings of probes
390    and p values, and distributions of p values and logFC estimates. For infant sex, SQN data consistently
391    had the highest number of significant negative probes and the lowest number of significant positive
392    probes, while SWAN data always had the lowest number of significant positive probes. For BMI and
393    intervention effects, only BMIQ data produced DMPs where no ComBat processing was used; in data
394    processed using supervised ComBat, all three normalisation methods resulted in some DMPs, but the
395    number and identity of these probes differed. In fake group data, SQN data produced a large number of
396    significant probes in non-ComBat-processed and supervised-ComBat data, while BMIQ and SWAN data
397    produced a small number of probes in supervised-ComBat data only; again, the number of significant
398    probes differed between the normalisation methods (see Tables 1-3 and Supplementary Figures S1-S4).
399
400    There was a fair degree of consistency in rankings of probes by p value for infant sex, but some large
401    discrepancies in rankings for BMI, intervention, and fake group effects.  The rankings were less
402    consistent for highest-ranked probes by logFC, with some quite large differences in both rankings and
403    effect estimates (including different directions of effect) for infant sex, BMI, intervention and fake
404    groups. BMIQ estimates tended to be more extreme (further from 0) than the other two methods.
405
406    Distributions of p values and logFC estimates also differed between normalisation methods. For p values
407    the differences were not consistent across models and effects, but for logFC there was a clear difference
408    between BMIQ and the other two methods, with the range of estimates in BMIQ data being much more
409    widely dispersed; SQN and SWAN data had more similar distributions, but SQN was moderately
410    narrower than SWAN across all effects and models.
411
412    Overall, there was little difference between SQN and SWAN methods when adjusting for batch in the
413    model.  There is some evidence that SQN would result in fewer significant DMPs than SWAN for known
414    effects (particularly when using supervised ComBat), but (many) more spuriously significant DMPs than
415    either SWAN or BMIQ where effects are absent. The behaviour of BMIQ was more variable depending
416    on other dimensions of the pipeline, but had a wider dispersion of logFC estimates than the other
417    methods, particularly when adjusting for batch in the model.  This tended to result in more DMPs in
418    some scenarios, but in general will lead effect estimates derived from BMIQ data to be more extreme
419    (and probably overestimates of the true effect).
420
421    ## Effect of Batch Correction Method
422    There were clear differences in all dimensions between batch correction methods. For all effects (infant
423    sex, BMI, intervention and fake groups), supervised ComBat processing produced a larger number of
424    DMPs compared to either unsupervised ComBat processing or adjustment for batch in the analysis
425    model. The difference between unsupervised ComBat and batch-adjustment was less consistent for
426    infant sex effects, but for BMI, intervention and fake group effects, there were no DMPs in unsupervised
427    ComBat models, whereas there were a few for batch-adjusted models.

428

429  Rankings of top probes by p value were relatively consistent between batch-adjustment methods for
430  infant sex, but there were some large discrepancies particularly for BMI and intervention effects, and
431  especially in BMIQ data.  The same phenomenon was observed for logFC rankings, which also showed a
432  tendency for logFC estimates in unsupervised-ComBat data to be smaller in absolute magnitude (closer
433  to 0).

434

435  The distribution of p values showed clear and consistent differences between batch-correction methods,
436  with the distribution in supervised ComBat data shifted substantially towards 0 relative to both
437  unsupervised ComBat and batch-adjusted models, for all effects. For logFC estimates, supervised
438  ComBat and batch-adjusted data were generally fairly similar, but unsupervised ComBat data generally
439  resulted in a narrower range.  This means that supervised ComBat will tend to produce more statistically
440  significant probes, regardless of the presence or absence of an effect.  Conversely, effect estimates from
441  unsupervised ComBat may be underestimated; at least, they will tend to be smaller in magnitude than
442  those derived from data where batch is handled differently.

443

444  Of particular note is the combination of SQN normalisation and either adjustment for batch in the
445  model, or use of supervised ComBat.  These combinations produced an extremely large number of
446  significant DMPs for fake group effects; this was more extreme in the case of supervised ComBat
447  (producing over 6000 DMPs) than when adjusting for batch (somewhat over 2000 DMPs).  Additional
448  adjustment for cell type proportion ameliorated this effect, as discussed below, but in the case of
449  supervised ComBat data, did not eliminate spurious findings. This suggests that batch adjustment may
450  be particularly ill-advised in the context of SQN normalisation; since SQN involves between-array as well
451  as within-array normalisation, additional adjustment for batch may be over-correcting.

452

453

454  ## Effect of Adjustment for Estimated Cell Type Proportion

455  Adjustment for cell type proportion affected results, but the impact was not consistent across the
456  different types of effects studied. Adjustment for batch resulted in a substantially larger number of
457  DMPs (both negative and positive) for infant sex, but reduced the number of DMPs for fake groups (for
458  models where there were DMPs for fake groups effects). For BMI and intervention, the effect of cell
459  type adjustment was mostly but not entirely to produce more DMPs.

460

461  The effect of cell type adjustment on top probe rankings was fairly modest, although some quite large
462  discrepancies were observed for p value rankings, logFC rankings, and logFC estimates. The effect on
463  distribution of p values depended on the effect: for infant sex, adjustment for cell type proportion
464  consistently (for all normalisation and batch-correction methods) shifted the distribution downwards
465  towards 0 (i.e. more statistically significant probes), whereas the differences were less consistent and
466  smaller in BMI, intervention and fake group effects.  There were no large or consistent differences in
467  distribution of logFC estimates between cell-type-adjusted and non-adjusted models.

468

469  Overall, adjustment for cell type proportion tended to improve model behaviour regarding spurious
470  results: the number of significantly differentially methylated probes decreased with adjustment for cell
471  type proportion, though they were not always eliminated.  The number of differentially methylated
472  probes for infant sex was increased, which may reflect either improvement (greater ability to detect
473  true effects due to removal of noise due to cell type differences) or harm (greater number of spurious
474  effects) depending on whether the extra probes are in fact differentially methylated between males and
475  females; without knowing the true number and identity of DMPs, we cannot be certain.  Similarly,

476  adjustment for cell type proportion increased the number of DMPs for BMI effects in BMIQ and SWAN
477  data, in one scenario (BMIQ with Supervised ComBat) by a substantial amount (from 99 to 2017 DMPs)
478  and these are most likely to be false positives.
479
480

481  ## Discussion

482  Different choices in probe filtering, normalisation, batch handling, and adjustment for cell types resulted
483  in different findings regarding the presence and identity of differentially methylated probes, rankings of
484  probes by p value and log-fold-change, and different overall distributions of p values and log-fold-
485  change estimates. While some of these differences were relatively modest, our results nevertheless
486  show that particular combinations of data processing and analysis choices may result in spurious false
487  positive findings, and/or potentially the failure to detect true effects. Additionally, while the magnitude
488  of effect estimates is often not considered in differential methylation studies, some pipelines may result
489  in an overestimate or underestimate of the true effect.  Importantly, the results tended to depend on
490  *combinations* of choices rather than individual elements of the analysis pipeline.
491
492  The results of our analyses are consistent with other investigations which have been undertaken into
493  different data-processing and analysis choices. As noted above, the potential for 'false positives' to
494  result from supervised batch-correction methods specifying effects of interest has been previously
495  identified by a number of authors.(Nygaard, Rødland & Hovig, 2016a; Price & Robinson, 2018; Zindler et
496  al., 2020).  Our finding that the distribution of p values in the supervised ComBat algorithm tends to shift
497  the p value distribution downward is consistent with the finding of Nygaard et al that, in contexts where
498  the effects of interest are not evenly spread between batches, the distribution of F-statistics will be
499  biased upwards (Nygaard, Rødland & Hovig, 2016a).  While implementation as an unsupervised method
500  may be preferable, our findings suggest that this may create a different problem, with the estimates of
501  log-fold-change corresponding to effects of interest biased towards zero.
502
503  Wu et al's study (Wu et al., 2014) comparing a variety of normalisation approaches noted a tendency for
504  more statistically significant differences to arise in SQN data, which they hypothesise may be due to
505  reduced overall variance.  In our investigation, the main context in which SQN data produced a large
506  number of spurious differentially methylated probes was when supervised ComBat, or adjustment for
507  batch in the model, was used; we additionally found that adjustment for cell type proportion reduced
508  the number of spurious findings (while not necessarily eliminating them).  Thus, SQN is not universally
509  more prone to producing spurious findings than other normalisation methods.
510
511  Our findings do not suggest that there is one particular combination of methods which can be
512  guaranteed to 'work' in all contexts, and some of the recommendations which have been made by
513  others may need to be modified somewhat.  For example, Nygaard et al conclude that adjustment for
514  batch in the model is preferable to the use of batch-correction algorithms,(Nygaard, Rødland & Hovig,
515  2016b), but our results suggest that this is inadvisable for data that have been normalised using SQN; in
516  our data, this combination resulted in a large number of spurious findings. In general, while our results
517  support others' findings that supervised bach-correction algorithms should not be used,  there does not
518  appear to be much difference between unsupervised batch-correction and adjustment for batch in the
519  model.  The only caveat here is that some of our results (particularly regarding effects of infant sex)
520  suggest that unsupervised ComBat may underestimate the magnitude of effects, as the distribution of
521  logFC estimates was substantially narrower than other methods. The use of a more stringent method of
522  Type I error control may also help to reduce the number of spurious findings: the use of FDR correction

523    methods such as Benjamini-Hochberg, while very common (Maksimovic, Phipson & Oshlack, 2017), may
524    not be sufficient to deal with higher rates of spurious results (Nygaard, Rødland & Hovig, 2016a). In our
525    data, the use of the Holm method (which controls the Family-Wise Error Rate) reduced, but did not
526    eliminate, spurious findings associated with fake group effects. Investigation of DNA regions, rather than
527    probe-wise analysis, may also help to differentiate true methylation differences from spurious ones
528    (Wang et al., 2015): the statistically significant DMPs for fake group effects (as well  as for BMI and
529    intervention effects) tended to be isolated rather than being grouped in the same region, and in our
530    companion paper, we found no significant differences in methylation for groups of probes on candidate
531    genes (Louise et al., 2022).
532
533    One limitation of our study is our inability to compare model behaviour in relation to known effects.  It
534    was relatively simple to create fake groups to study behaviour of models for effects which were known
535    *not* to exist, but as we do not know the truth about which effects actually exist in our data, we could not
536    compare behaviour of models in their ability to detect these known effects.  Simulated data could
537    potentially be used for this purpose; however, the effects in the simulated data would have to be
538    biologically plausible.  This was beyond the scope of our study; however, it is a good subject for future
539    research.  We used infant sex as the nearest proxy to a known effect, as we knew at least that some
540    effects existed. However, we cannot say whether, and to what extent, the differences observed in
541    relation to infant sex reflect spurious findings versus the failure to detect true effects.
542
543    Overall, as many other authors have noted, researchers working with DNAm data should better
544    understand the methods built into standard pipelines (Price & Robinson, 2018; Zindler et al., 2020), and
545    should better document the specific data-processing methods used (Nygaard, Rødland & Hovig, 2016a;
546    Zindler et al., 2020).  It is also important, in our view, to pay more attention to the context in which a
547    particular epigenome-wide analysis is performed. For example, a less stringent method of Type I error
548    control may often be chosen because the study is exploratory (hypothesis-generating) rather than
549    confirmatory, and it is considered more important not to miss potential findings than to rule out
550    spurious ones. In this case, the results from such studies should be interpreted accordingly: as
551    suggestive findings which cannot be confidently accepted until they are validated in new data. The
552    validation of existing findings should be treated as a high priority in epigenetics research (Price &
553    Robinson, 2018).
554
555    Additionally, the degree of confidence that can be placed in any new discoveries could be enhanced by
556    performing sensitivity analyses – re-performing analyses using different normalisation methods, batch
557    correction methods, or models -  which we believe should become standard in this area.

## References

560  Aryee MJ, Jaffe AE, Corrada-Bravo H, Ladd-Acosta C, Feinberg AP, Hansen KD, Irizarry RA. 2014. Minfi: a

561  flexible and comprehensive Bioconductor package for the analysis of Infinium DNA methylation

562  microarrays. *Bioinformatics* 30:1363–1369. DOI: 10.1093/bioinformatics/btu049.

563  Bakulski KM, Feinberg JI, Andrews SV, Yang J, Brown S, L McKenney S, Witter F, Walston J, Feinberg AP,

564  Fallin MD. 2016. DNA methylation of cord blood cell types: Applications for mixed cell birth

565  studies. *Epigenetics* 11:354–362. DOI: 10.1080/15592294.2016.1161875.

566  Cardenas A, Allard C, Doyon M, Houseman EA, Bakulski KM, Perron P, Bouchard L, Hivert M-F. 2016.

567  Validation of a DNA methylation reference panel for the estimation of nucleated cells types in

568  cord blood. *Epigenetics* 11:773–779. DOI: 10.1080/15592294.2016.1233091.

569  Chen Y, Lemire M, Choufani S, Butcher DT, Grafodatskaya D, Zanke BW, Gallinger S, Hudson TJ,

570  Weksberg R. 2013. Discovery of cross-reactive probes and polymorphic CpGs in the Illumina

571  Infinium HumanMethylation450 microarray. *Epigenetics* 8:203–209. DOI: 10.4161/epi.23470.

572  Dedeurwaerder S, Defrance M, Bizet M, Calonne E, Bontempi G, Fuks F. 2014. A comprehensive

573  overview of Infinium HumanMethylation450 data processing. *Briefings in bioinformatics* 15:929–

574  941. DOI: 10.1093/bib/bbt054.

575  Dodd JM, Cramp C, Sui Z, Yelland LN, Deussen AR, Grivell RM, Moran LJ, Crowther CA, Turnbull D,

576  McPhee AJ, Wittert G, Owens JA, Robinson JS, For the LIMIT Randomised Trial Group. 2014a.

577  The effects of antenatal dietary and lifestyle advice for women who are overweight or obese on

578  maternal diet and physical activity: the LIMIT randomised trial. *BMC Medicine* 12:161. DOI:

579  10.1186/s12916-014-0161-y.

580  Dodd JM, Turnbull D, McPhee AJ, Deussen AR, Grivell RM, Yelland LN, Crowther CA, Wittert G, Owens JA,

581  Robinson JS, for the LIMIT Randomised Trial Group. 2014b. Antenatal lifestyle advice for women

582          who are overweight or obese: LIMIT randomised trial. *BMJ* 348:g1285–g1285. DOI:

583          10.1136/bmj.g1285.

584     Fortin J-P, Labbe A, Lemire M, Zanke BW, Hudson TJ, Fertig EJ, Greenwood CM, Hansen KD. 2014.

585          Functional normalization of 450k methylation array data improves replication in large cancer

586          studies. *Genome Biology* 15. DOI: 10.1186/s13059-014-0503-2.

587     Fortin J-P, Triche T, Hansen K. 2016. Preprocessing, normalization and integration of the Illumina

588          HumanMethylationEPIC array. DOI: 10.1101/065490.

589     Gemma C, Sookoian S, Alvariñas J, García SI, Quintana L, Kanevsky D, González CD, Pirola CJ. 2009.

590          Maternal Pregestational BMI Is Associated With Methylation of the *PPARGC1A* Promoter in

591          Newborns. *Obesity* 17:1032–1039. DOI: 10.1038/oby.2008.605.

592     Heiss JA, Just AC. 2019. Improved filtering of DNA methylation microarray data by detection p values

593          and its impact on downstream analyses. *Clinical Epigenetics* 11:15. DOI: 10.1186/s13148-019-

594          0615-3.

595     Hicks SC, Irizarry RA. 2015. quantro: a data-driven approach to guide the choice of an appropriate

596          normalization method. *Genome Biology* 16:117. DOI: 10.1186/s13059-015-0679-0.

597     Hjort L, Martino D, Grunnet LG, Naeem H, Maksimovic J, Olsson AH, Zhang C, Ling C, Olsen SF, Saffery R,

598          Vaag AA. 2018. Gestational diabetes and maternal obesity are associated with epigenome-wide

599          methylation changes in children. *JCI Insight* 3:e122572. DOI: 10.1172/jci.insight.122572.

600     Jaffe AE, Irizarry RA. 2014. Accounting for cellular heterogeneity is critical in epigenome-wide

601          association studies. *Genome biology* 15:R31–R31. DOI: 10.1186/gb-2014-15-2-r31.

602     Lehne B, Drong AW, Loh M, Zhang W, Scott WR, Tan S-T, Afzal U, Scott J, Jarvelin M-R, Elliott P,

603          McCarthy MI, Kooner JS, Chambers JC. 2015. A coherent approach for analysis of the Illumina

604          HumanMethylation450 BeadChip improves data quality and performance in epigenome-wide

605          association studies. *Genome biology* 16:37–37. DOI: 10.1186/s13059-015-0600-x.

606    Louise J, Deussen AR, Koletzko B, Owens J, Saffery R, Dodd JM. 2022. Effect of an antenatal diet and

607        lifestyle intervention and maternal BMI on cord blood DNA methylation in infants of overweight

608        and obese women: The LIMIT Randomised Controlled Trial. *PLOS ONE* 17:e0269723. DOI:

609        10.1371/journal.pone.0269723.

610    Maksimovic J, Gordon L, Oshlack A. 2012. SWAN: Subset-quantile Within Array Normalization for

611        Illumina Infinium HumanMethylation450 BeadChips. *Genome Biology* 13:R44. DOI: 10.1186/gb-

612        2012-13-6-r44.

613    Maksimovic J, Phipson B, Oshlack A. 2017. A cross-package Bioconductor workflow for analysing

614        methylation array data. *F1000Research* 5:1281. DOI: 10.12688/f1000research.8839.3.

615    Marabita F, Almgren M, Lindholm ME, Ruhrmann S, Fagerström-Billai F, Jagodic M, Sundberg CJ, Ekström

616        TJ, Teschendorff AE, Tegnér J, Gomez-Cabrero D. 2013. An evaluation of analysis pipelines for

617        DNA methylation profiling using the Illumina HumanMethylation450 BeadChip platform.

618        *Epigenetics* 8:333–346. DOI: 10.4161/epi.24008.

619    Morris TJ, Beck S. 2015. Analysis pipelines and packages for Infinium HumanMethylation450 BeadChip

620        (450k) data. *Methods* 72:3–8. DOI: 10.1016/j.ymeth.2014.08.011.

621    Naeem H, Wong NC, Chatterton Z, Hong MKH, Pedersen JS, Corcoran NM, Hovens CM, Macintyre G.

622        2014. Reducing the risk of false discovery enabling identification of biologically significant

623        genome-wide methylation status using the HumanMethylation450 array. *BMC Genomics* 15:51.

624        DOI: 10.1186/1471-2164-15-51.

625    Nygaard V, Rødland EA, Hovig E. 2016a. Methods that remove batch effects while retaining group

626        differences may lead to exaggerated confidence in downstream analyses. *Biostatistics* 17:29–39.

627        DOI: 10.1093/biostatistics/kxv027.

628    Nygaard V, Rødland EA, Hovig E. 2016b. Methods that remove batch effects while retaining group

629        differences may lead to exaggerated confidence in downstream analyses. *Biostatistics* 17:29–39.

630        DOI: 10.1093/biostatistics/kxv027.

631    Peters TJ, Buckley MJ, Statham AL, Pidsley R, Samaras K, V Lord R, Clark SJ, Molloy PL. 2015. De novo

632        identification of differentially methylated regions in the human genome. *Epigenetics &*

633        *Chromatin* 8:6. DOI: 10.1186/1756-8935-8-6.

634    Pidsley R, CC YW, Volta M, Lunnon K, Mill J, Schalkwyk LC. 2013. A data-driven approach to

635        preprocessing Illumina 450K methylation array data. *BMC Genomics* 14:293–293. DOI:

636        10.1186/1471-2164-14-293.

637    Price EM, Robinson WP. 2018. Adjusting for Batch Effects in DNA Methylation Microarray Data, a Lesson

638        Learned. *Frontiers in Genetics* 9. DOI: 10.3389/fgene.2018.00083.

639    R Core Team. 2018. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R

640        Foundation for Statistical Computing.

641    Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, Smyth GK. 2015. limma powers differential

642        expression analyses for RNA-sequencing and microarray studies. *Nucleic acids research* 43:e47–

643        e47. DOI: 10.1093/nar/gkv007.

644    Sharp GC, Lawlor DA, Richmond RC, Fraser A, Simpkin A, Suderman M, Shihab HA, Lyttleton O, McArdle

645        W, Ring SM, Gaunt TR, Smith GD, Relton CL. 2015. Maternal pre-pregnancy BMI and gestational

646        weight gain, offspring DNA methylation and later offspring adiposity: Findings from the Avon

647        Longitudinal Study of Parents and Children. *International Journal of Epidemiology* 44:1288–

648        1304. DOI: 10.1093/ije/dyv042.

649    Sharp GC, Salas LA, Monnereau C, Allard C, Yousefi P, Everson TM, Bohlin J, Xu Z, Huang RC, Reese SE, Xu

650        CJ, Baïz N, Hoyo C, Agha G, Roy R, Holloway JW, Ghantous A, Merid SK, Bakulski KM, Küpers LK,

651        Zhang H, Richmond RC, Page CM, Duijts L, Lie RT, Melton PE, Vonk JM, Nohr EA, Williams-

652      DeVane CL, Huen K, Rifas-Shiman SL, Ruiz-Arenas C, Gonseth S, Rezwan FI, Herceg Z, Ekström S,

653      Croen L, Falahi F, Perron P, Karagas MR, Quraishi BM, Suderman M, Magnus MC, Jaddoe VWV,

654      Taylor JA, Anderson D, Zhao S, Smit HA, Josey MJ, Bradman A, Baccarelli AA, Bustamante M,

655      Håberg SE, Pershagen G, Hertz-Picciotto I, Newschaffer C, Corpeleijn E, Bouchard L, Lawlor DA,

656      Maguire RL, Barcellos LF, Smith GD, Eskenazi B, Karmaus W, Marsit CJ, Hivert MF, Snieder H,

657      Fallin MD, Melén E, Munthe-Kaas MC, Arshad H, Wiemels JL, Annesi-Maesano I, Vrijheid M,

658      Oken E, Holland N, Murphy SK, Sørensen TIA, Koppelman GH, Newnham JP, Wilcox AJ, Nystad

659      W, London SJ, Felix JF, Relton CL. 2017. Maternal BMI at the start of pregnancy and offspring

660      epigenome-wide DNA methylation: Findings from the pregnancy and childhood epigenetics

661      (PACE) consortium. *Human Molecular Genetics* 26:4067–4085. DOI: 10.1093/hmg/ddx290.

662   Smyth GK. limma: Linear Models for Microarray Data. *Bioinformatics and Computational Biology*

663      *Solutions Using R and Bioconductor*:397–420. DOI: 10.1007/0-387-29362-0_23.

664   Teschendorff AE, Marabita F, Lechner M, Bartlett T, Tegner J, Gomez-Cabrero D, Beck S. 2013. A beta-

665      mixture quantile normalization method for correcting probe design bias in Illumina Infinium 450

666      k DNA methylation data. *Bioinformatics* 29:189–196. DOI: 10.1093/bioinformatics/bts680.

667   Teschendorff AE, Zheng SC. 2017. Cell-type deconvolution in epigenome-wide association studies: a

668      review and recommendations. *Epigenomics* 9:757–768. DOI: 10.2217/epi-2016-0153.

669   Thabane L, Mbuagbaw L, Zhang S, Samaan Z, Marcucci M, Ye C, Thabane M, Giangregorio L, Dennis B,

670      Kosa D, Debono VB, Dillenburg R, Fruci V, Bawor M, Lee J, Wells G, Goldsmith CH. 2013. A

671      tutorial on sensitivity analyses in clinical trials: the what, why, when and how. *BMC Medical*

672      *Research Methodology* 13:92. DOI: 10.1186/1471-2288-13-92.

673   Thakali KM, Faske JB, Ishwar A, Alfaro MP, Cleves MA, Badger TM, Andres A, Shankar K. 2017. Maternal

674      obesity and gestational weight gain are modestly associated with umbilical cord DNA

675      methylation. *Placenta* 57:194–203. DOI: 10.1016/j.placenta.2017.07.009.

676    Wang T, Guan W, Lin J, Boutaoui N, Canino G, Luo J, Celedón JC, Chen W. 2015. A systematic study of

677         normalization methods for Infinium 450K methylation data using whole-genome bisulfite

678         sequencing data. *Epigenetics* 10:662–669. DOI: 10.1080/15592294.2015.1057384.

679    Wang D, Zhang Y, Huang Y, Li P, Wang M, Wu R, Cheng L, Zhang W, Zhang Y, Li B, Wang C, Guo Z. 2012.

680         Comparison of different normalization assumptions for analyses of DNA methylation data from

681         the cancer genome. *Gene* 506:36–42. DOI: 10.1016/j.gene.2012.06.075.

682    Wu MC, Joubert BR, Kuan P, Håberg SE, Nystad W, Peddada SD, London SJ. 2014. A systematic

683         assessment of normalization approaches for the Infinium 450K methylation platform.

684         *Epigenetics* 9:318–329. DOI: 10.4161/epi.27119.

685    Yousefi P, Huen K, Aguilar Schall R, Decker A, Elboudwarej E, Quach H, Barcellos L, Holland N. 2013.

686         Considerations for normalization of DNA methylation data by Illumina 450K BeadChip assay in

687         population studies. *Epigenetics* 8:1141–1152. DOI: 10.4161/epi.26037.

688    Yousefi P, Huen K, Davé V, Barcellos L, Eskenazi B, Holland N. 2015. Sex differences in DNA methylation

689         assessed by 450 K BeadChip in newborns. *BMC Genomics* 16:911. DOI: 10.1186/s12864-015-

690         2034-y.

691    Zindler T, Frieling H, Neyazi A, Bleich S, Friedel E. 2020. Simulating ComBat: how batch correction can

692         lead to the systematic introduction of false positive results in DNA methylation microarray

693         studies. *BMC bioinformatics* 21:271. DOI: 10.1186/s12859-020-03559-6.

694
695

# Figure 1

Flowchart of Data Processing and Analysis

Combinations of data-processing and analysis choices, consisting of six normalised datasets (SQN, BMIQ or SWAN, with probe filering before or afterwards), use or non-use of ComBat processing (supervised or unsupervised), and analysis with either an unadjusted model, a model adjusted for batch, and a model adjusted for batch and cell type proportion.
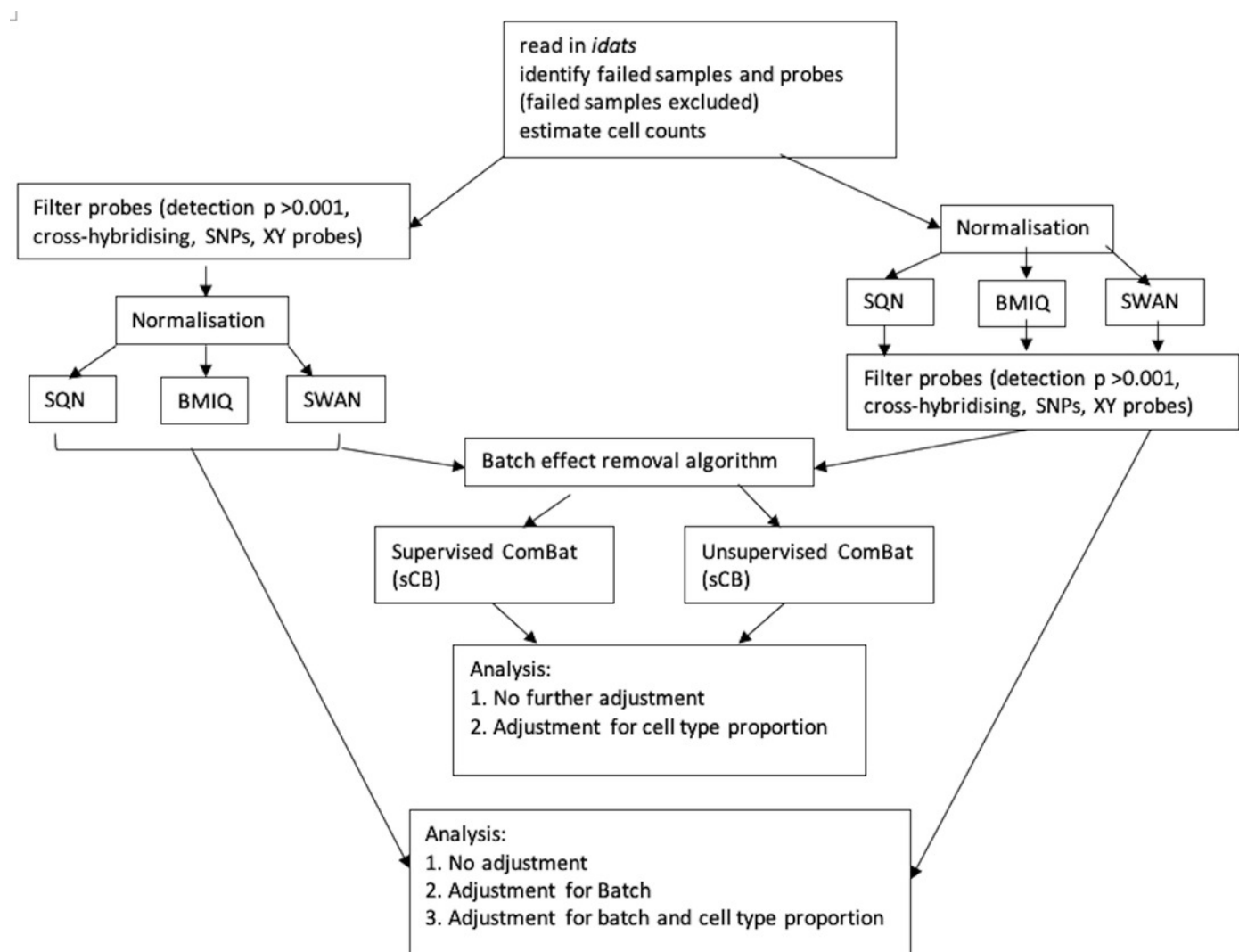
PeerJ

# Figure 2

Probes ranked in top 10 by p-value in Batch+Cell Adjusted Model, for (a) Infant Sex, (b) BMI in Standard Care, (c) Short-Haired in Tabby

For each probe the rank is given by pre- vs post-filtering, normalisation method, and batch-handling method. The model is one adjusting for batch (either explicitly in the model or via batch-correction algorithm) and cell type proportion. Adjust=adjusted for batch in the model; SCB=Supervised ComBat; UCB=Unsupervised ComBat.
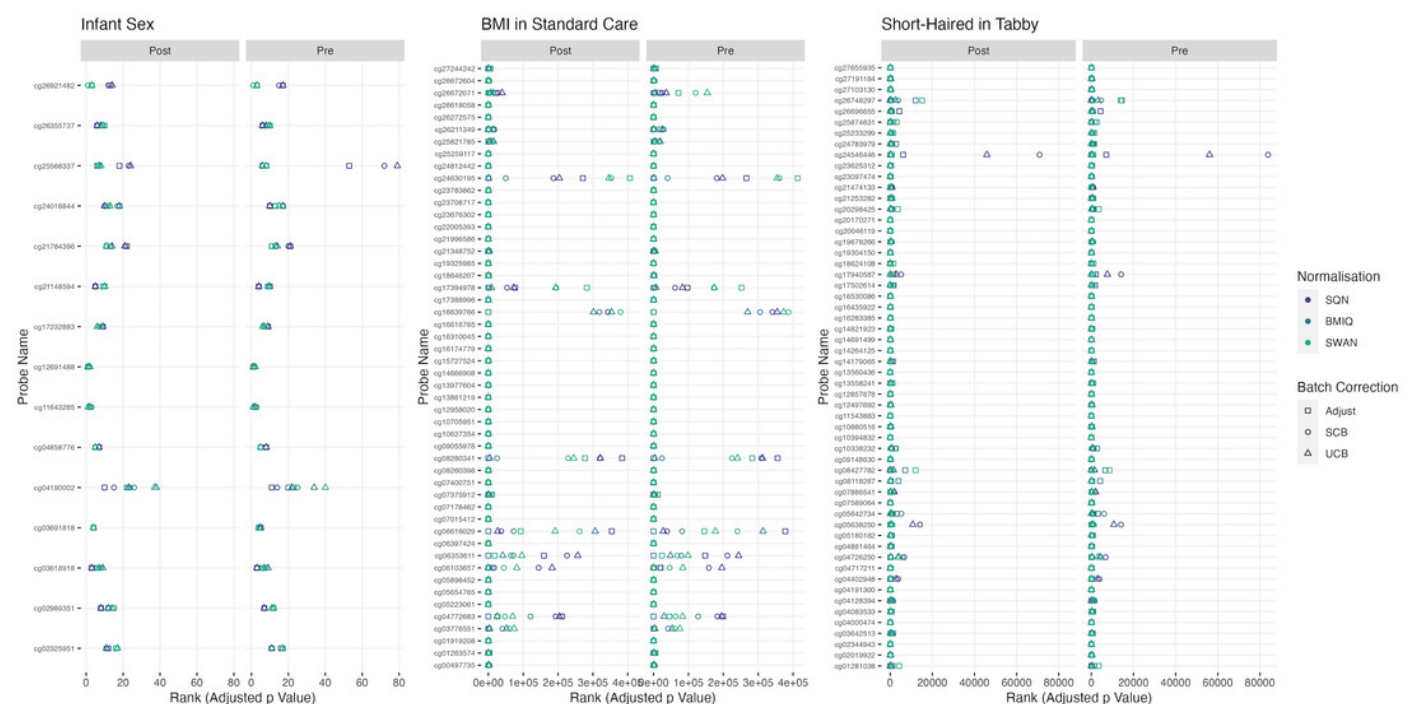
# Figure 3

Top 10 Probes by LogFC: Infant Sex

Largest LogFC for Infant Sex (female), by normalisation and batch correction method
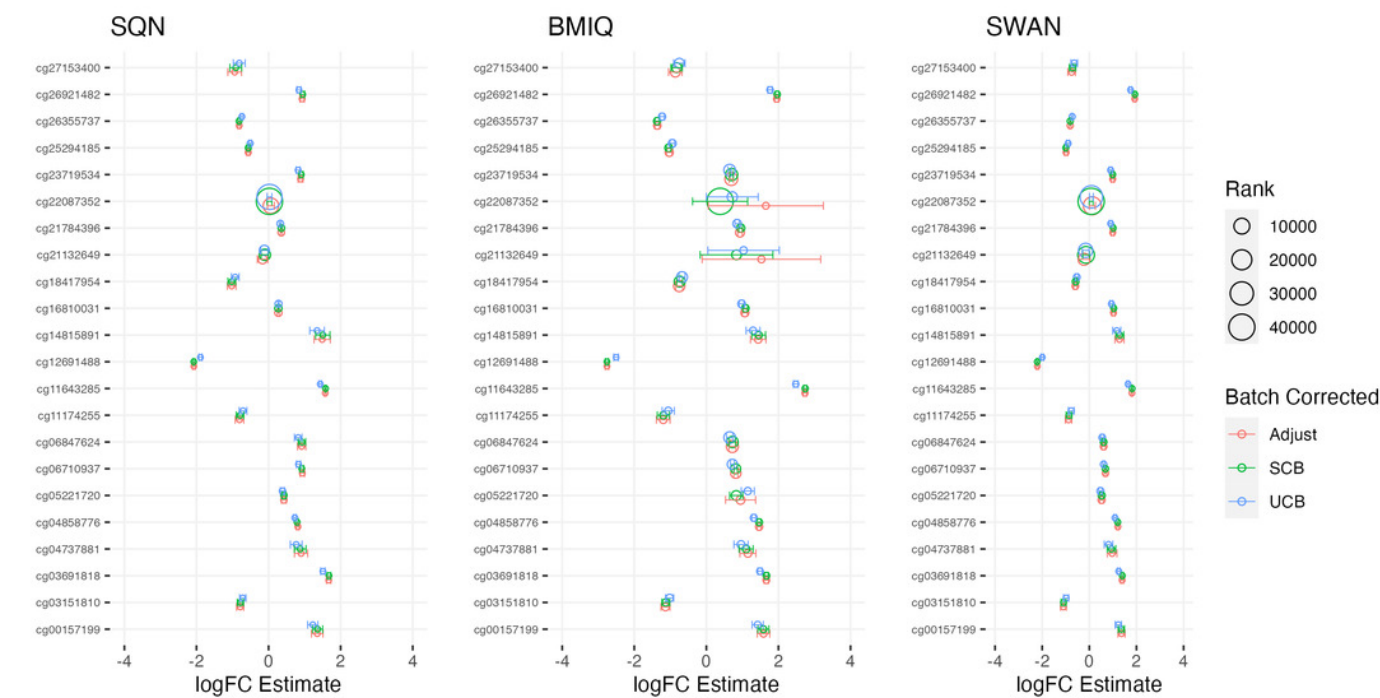
# Figure 4

Top 10 Probes by LogFC: BMI in Standard Care

Largest LogFC for effect of BMI in Standard Care group, by normalisation and batch correction method
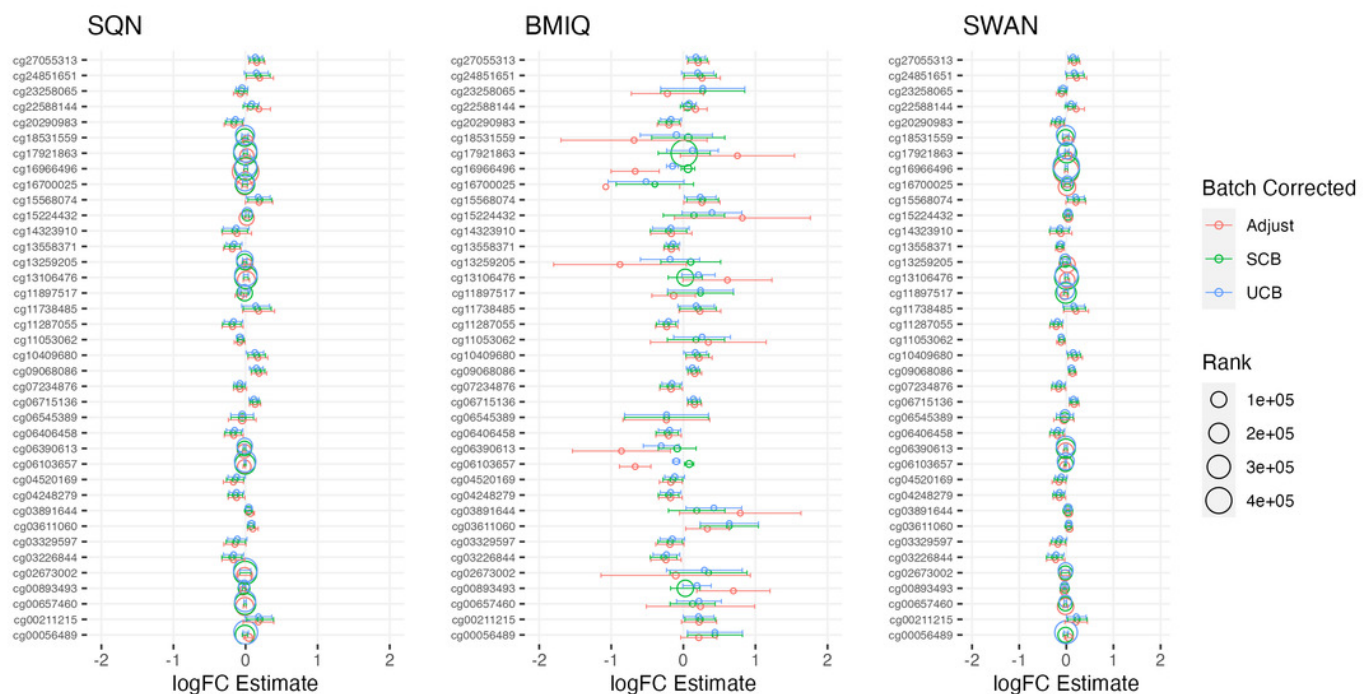
# Figure 5

Top 10 Probes by LogFC: 'Short-Haired' in 'Tabby'

Largest LogFC for effect of 'Short-Haired' in 'Tabby' group, by normalisation and batch correction method.
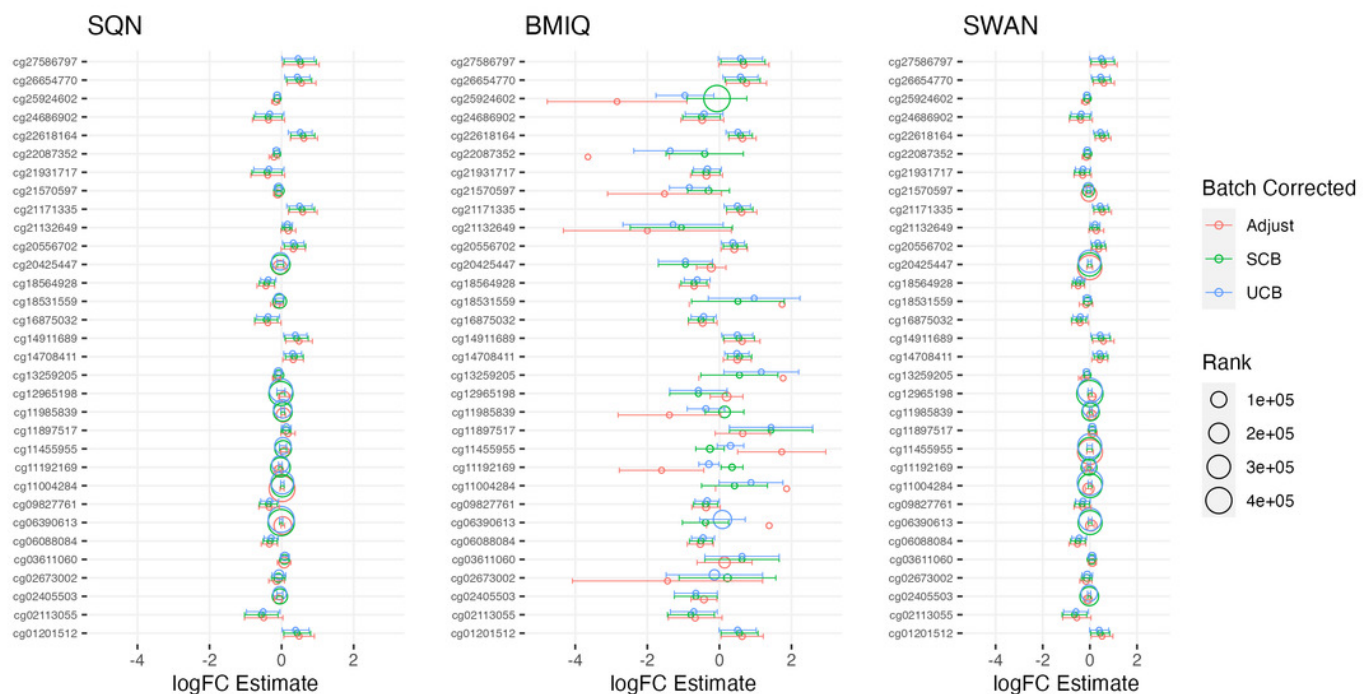
# Figure 6

Distribution of Unadjusted P Values by Normalisation and Batch Correction Method, for Batch and Cell Adjusted Models

Only models from data where probe filtering was performed post-normalisation are included.

The model is one adjusting for batch (either explicitly in the model or via batch-correction algorithm) and cell type proportion.
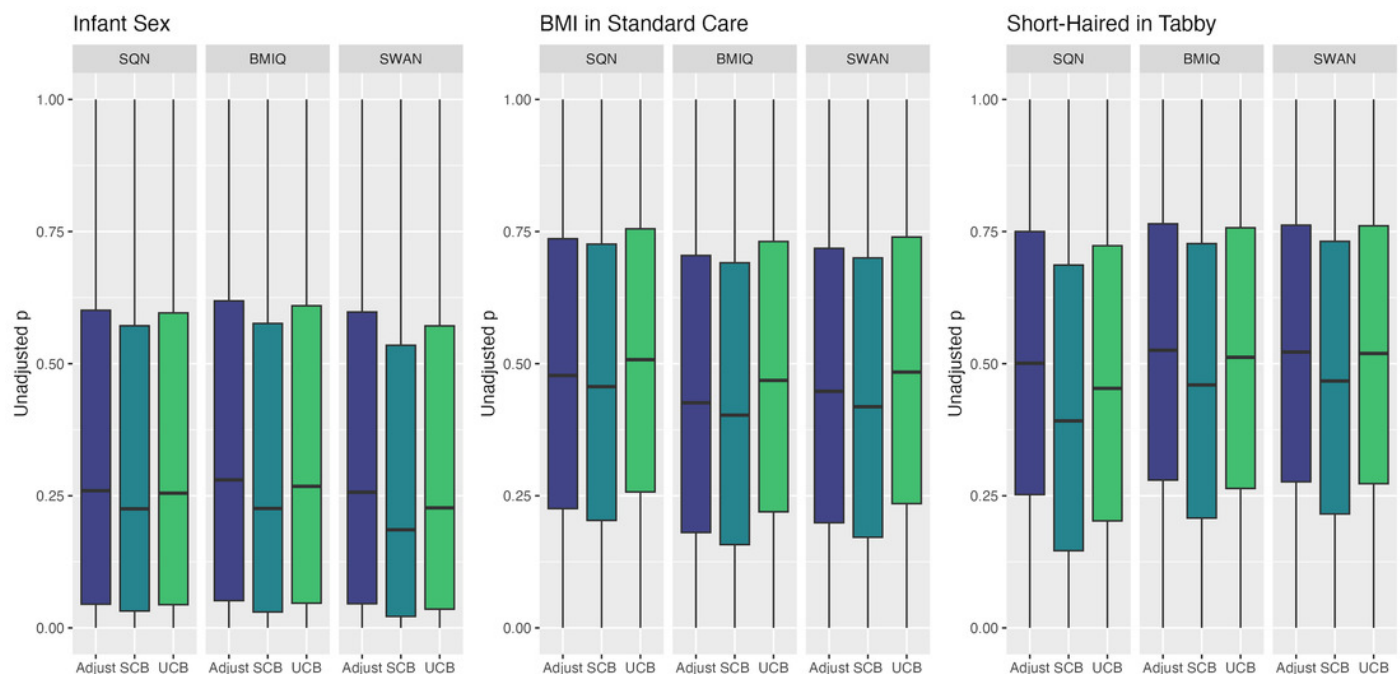
# Figure 7

Distribution of Log-Fold-Change Estimates by Normalisation and Batch Correction Method, for Batch and Cell Adjusted Models

Only models from data where probe filtering was performed post-normalisation are included. The model is one adjusting for batch (either explicitly in the model or via batch-correction algorithm) and cell type proportion.
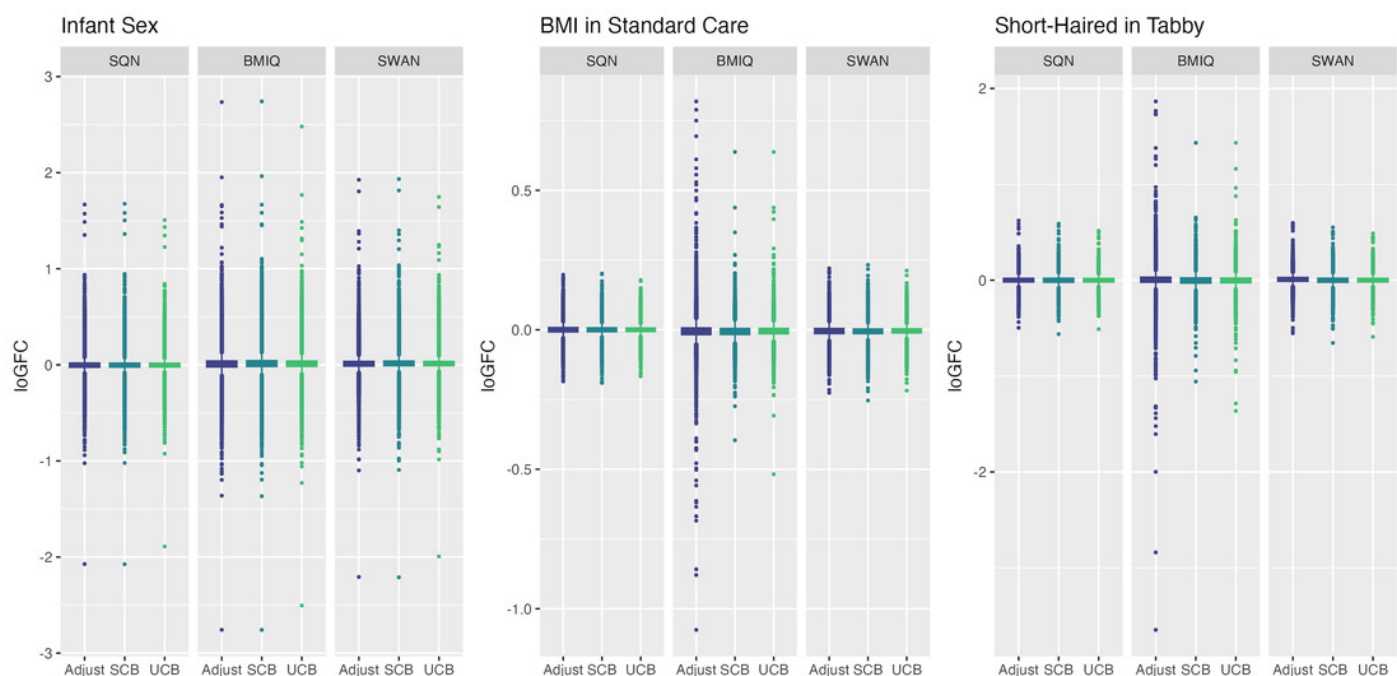
# Table 1(on next page)

Number of DMPs for Infant Sex (Female), by Probe Filtering Method, Batch Correction Method, Normalisation Method and Cell Type Method

* No adjustment beyond the correction for batch as implemented in the ComBat algorithm

1 Table 1. Number of DMPs for Infant Sex (Female), by Probe Filtering Method, Batch Correction Method, Normalisation Method and Cell Type
2 Method
3

| Model | | SQN | | BMIQ | | SWAN | |
|---|---|---|---|---|---|---|---|
| | | Post Filtered | Pre Filtered | Post Filtered | Pre Filtered | Post Filtered | Pre Filtered |
| No ComBat | | | | | | | |
| - Unadjusted | Down | 15088 | 14878 | 7935 | 7554 | 7581 | 6890 |
| | Up | 20587 | 20441 | 28004 | 30741 | 25784 | 25962 |
| - Adjusted for Batch | Down | 13132 | 13215 | 6602 | 6112 | 7100 | 6140 |
| | Up | 20225 | 19956 | 39482 | 34408 | 30362 | 29780 |
| - Adjusted for Batch + Cell | Down | 28406 | 28633 | 15855 | 14900 | 14239 | 10408 |
| | Up | 31973 | 32204 | 39719 | 45255 | 43155 | 40709 |
| Supervised ComBat | | | | | | | |
| - Unadjusted* | Down | 20967 | 21230 | 11235 | 10180 | 10772 | 9252 |
| | Up | 25690 | 25320 | 41518 | 48191 | 45193 | 43237 |
| - Adjusted for Cell | Down | 35559 | 36022 | 18036 | 16512 | 16423 | 12198 |
| | Up | 37068 | 36972 | 56521 | 64851 | 68769 | 65109 |
| UnSupervised ComBat | | | | | | | |
| - Unadjusted* | Down | 14603 | 14763 | 7634 | 6836 | 7344 | 6377 |
| | Up | 21336 | 21041 | 30961 | 34882 | 31892 | 31170 |
| - Adjusted for Cell | Down | 28012 | 28286 | 14060 | 12916 | 13030 | 9560 |
| | Up | 32478 | 32447 | 43370 | 49520 | 52037 | 49102 |

4 * No adjustment beyond the correction for batch as implemented in the ComBat algorithm
5

**Table 2**(on next page)

DMPs for Intervention and BMI Effects

1    Table 2 DMPs for Effect of Maternal BMI in the Standard Care group

2

| | SQN | | BMIQ | | SWAN | |
|---|---|---|---|---|---|---|
| Model | Post Filtered | Pre Filtered | Post Filtered | Pre Filtered | Post Filtered | Pre Filtered |
| No ComBat | | | | | | |
| - Unadjusted | 0 | 0 | 5 | 6 | 0 | 0 |
| - Adjusted for Batch | 0 | 0 | 6 | 0 | 0 | 0 |
| - Adjusted for Batch + Cell | 0 | 0 | 8 | 0 | 6 | 0 |
| Supervised ComBat | | | | | | |
| - Unadjusted | 0 | 0 | 0 | 10 | 0 | 0 |
| - Adjusted for Cell | 0 | 0 | 99 | 207 | 0 | 0 |
| UnSupervised ComBat | | | | | | |
| - Unadjusted | 0 | 0 | 0 | 0 | 0 | 0 |
| - Adjusted for Cell | 0 | 0 | 0 | 0 | 0 | 6 |

3

**Table 3**(on next page)

DMPs for Fake Groups

PeerJ

1  **Table** 3 **DMPs for Fa**ke Groups: 'Short-Haired' **in** 'Tabby'

2

|  | SQN | | BMIQ | | SWAN | |
|---|---|---|---|---|---|---|
| Model | Post Filtered | Pre Filtered | Post Filtered | Pre Filtered | Post Filtered | Pre Filtered |
| No ComBat | | | | | | |
| - Unadjusted | 0 | 0 | 0 | 0 | 0 | 0 |
| - Adjusted for Batch | 2180 | 2574 | 0 | 0 | 0 | 0 |
| - Adjusted for Batch + Cell | 0 | 0 | 0 | 0 | 0 | 0 |
| Supervised ComBat | | | | | | |
| - Unadjusted | 6768 | 7007 | 3 | 6 | 8 | 8 |
| - Adjusted for Cell | 123 | 133 | 1 | 1 | 0 | 0 |
| UnSupervised ComBat | | | | | | |
| - Unadjusted | 0 | 0 | 0 | 0 | 0 | 0 |
| - Adjusted for Cell | 0 | 0 | 0 | 0 | 0 | 0 |

3