

# The Epigenomics of Human Ageing

**Academic Unit:** Human Development and Health

**Supervisors:** Chris Bell, Karen Lillycrop and Cyrus Cooper

ORCID: [0000-0002-2574-9611](https://orcid.org/0000-0002-2574-9611)

Richard J. Acton

2020-09-17



# Contents

<b>Research Thesis: Declaration of Authorship</b>	<b>5</b>
<b>Abstract</b>	<b>7</b>
<b>Acknowledgements</b>	<b>9</b>
<b>1 Introduction</b>	<b>11</b>
1.1 The Ageing Population and The Burden of Disease . . . . .	11
1.2 The biology of Ageing . . . . .	15
1.3 Epigenomics - Overview . . . . .	18
1.4 Fundamentals of DNA Methylation . . . . .	20
1.5 DNA Methylation and Ageing . . . . .	31
1.6 Aims . . . . .	45
<b>2 Methods</b>	<b>47</b>
2.1 Illumina DNA Methylation arrays . . . . .	47
2.2 MeDIP-seq . . . . .	53
2.3 Targeted Bisulfite sequencing . . . . .	57
<b>3 Epigenome Wide Association Studies for Bone Health Outcomes in Umbilical Cord Blood and Tissue</b>	<b>61</b>
3.1 Abstract . . . . .	61
3.2 Introduction . . . . .	62
3.3 Methods . . . . .	66
3.4 Results . . . . .	69
3.5 Discussion . . . . .	117
<b>4 The Genomic Loci of Specific Human tRNA Genes Exhibit Ageing-Related DNA Hypermethylation</b>	<b>119</b>
4.1 Abstract . . . . .	119
4.2 Introduction . . . . .	121
4.3 Methods . . . . .	128
4.4 Results . . . . .	135
4.5 tRNA Gene Methylation in Cancer . . . . .	156
4.6 Discussion . . . . .	162
<b>5 DNA methylation clocks in Alu Elements</b>	<b>167</b>
5.1 Abstract . . . . .	167
5.2 Introduction . . . . .	167
5.3 Methods . . . . .	172
5.4 Results . . . . .	177
5.5 Discussion . . . . .	190
<b>6 Discussion</b>	<b>191</b>
<b>References</b>	<b>193</b>

**Appendices**

6.1 <i>BioRxiv</i> manuscript: The Genomic Loci of Specific Human tRNA Genes Exhibit Ageing-Related DNA Hypermethylation . . . . .	<b>195</b>
--	------------

# Research Thesis: Declaration of Authorship

Name: Richard J. Acton

Title of thesis: The Epigenomics of Human Ageing

I declare that this thesis and the work presented in it is my own and has been generated by me as the result of my own original research.

I confirm that:

1. This work was done wholly or mainly while in candidature for a research degree at this University;
2. Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated;
3. Where I have consulted the published work of others, this is always clearly attributed;
4. Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work;
5. I have acknowledged all main sources of help;
6. Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself;
7. Parts of this work have been published as: a *BioRxiv* pre-print [1]

Signature: \_\_\_\_\_

Date: 2020-09-17



# **Abstract**

...



# Acknowledgements

I would like to acknowledge Nevena Krstic for her work extracting the DNA for the MAVIDOS EPIC array analysis and Dr Millie Parsons for her assistance with the MAVIDOS sample metadata, as well as Dr Beth Curtis and Professor Nick Harvey of the MRC-LEU for their assistance with framing the research questions for the vitamin D and bone development outcomes work with the MAVIDOS samples. The MRC-LEU is supported by the Medical Research Council (MRC). I gratefully acknowledge the individuals from TwinsUK, Mavidos and the Hertfordshire cohort. I would like to acknowledge Nikki Graham for her assistance identifying suitable samples for targeted bisulfite sequencing work and her work extracting DNA from those samples. TwinsUK received funding from the Wellcome Trust (Ref: 081878/Z/06/Z), European Community's Seventh Framework Programme (FP7/2007-2013), the National Institute for Health Research (NIHR)-funded BioResource, Clinical Research Facility and Biomedical Research Centre based at Guy's and St Thomas' NHS Foundation Trust in partnership with King's College London. Further funding support for the EpiTwin project was obtained from the European Research Council (project number 250157) and BGI. SNP Genotyping was performed by The Wellcome Trust Sanger Institute and National Eye Institute via NIH/CIDR. I acknowledge the use of the IRIDIS High-Performance Computing Facility, and associated support services at the University of Southampton, in the completion of this work. I would like to thank the MRC Doctoral fund (1820097) for supporting this work. Finally, I would like to acknowledge the extensive help and support of my primary supervisor Dr Chris Bell.



# Chapter 1

## Introduction

### 1.1 The Ageing Population and The Burden of Disease

The global population is ageing as is evident from the changing shape of the world population pyramid (figure 1.1). The increase in average lifespan underlying this demographic shift is projected to continue in industrialised nations with a probability of >65% for women and 85% for men [2] (figure 1.2). Ageing is the main risk factor for cancer, cardiovascular, neurodegenerative diseases, and many other conditions including osteoarthritis [3]. Furthermore, mortality from these conditions increases according to a logistic function with age [4] (figure 1.3). In addition to the non-infectious diseases, immunosenescence and frailty in old age contribute to increased mortality from infectious diseases, as exemplified by dramatically higher hospital admissions and deaths for respiratory disease during seasonal flu epidemics in the over 75s [5] (figure 1.4). Ageing is a common underlying risk factor for many conditions and as this would lead one to expect the number of concomitant disorders and the proportion of persons with multiple morbidities increases with age [6] (figure 1.5).

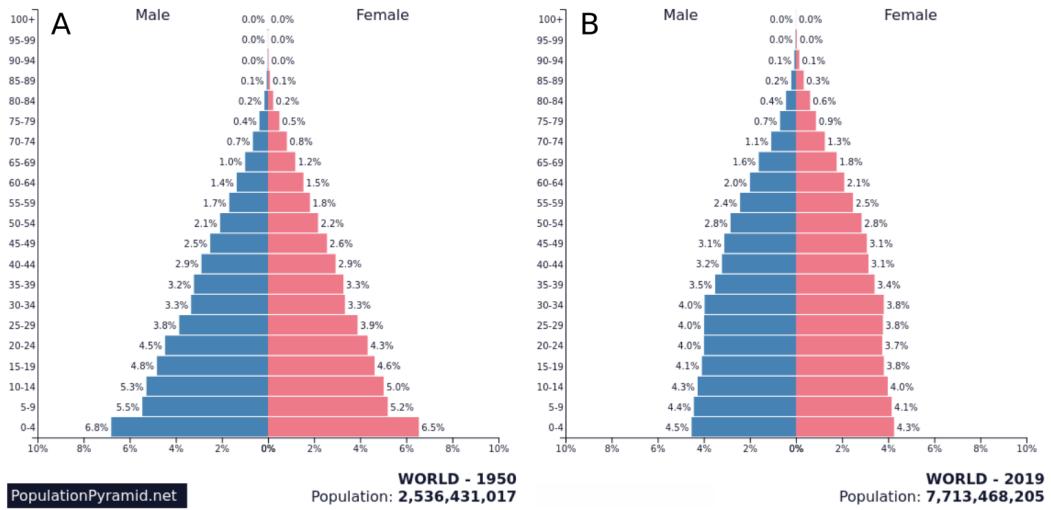


Figure 1.1: **The Population is Getting Older** Population pyramids for global population for A) 1950 and B) 2019 from [populationpyramid.net](http://populationpyramid.net) [7]

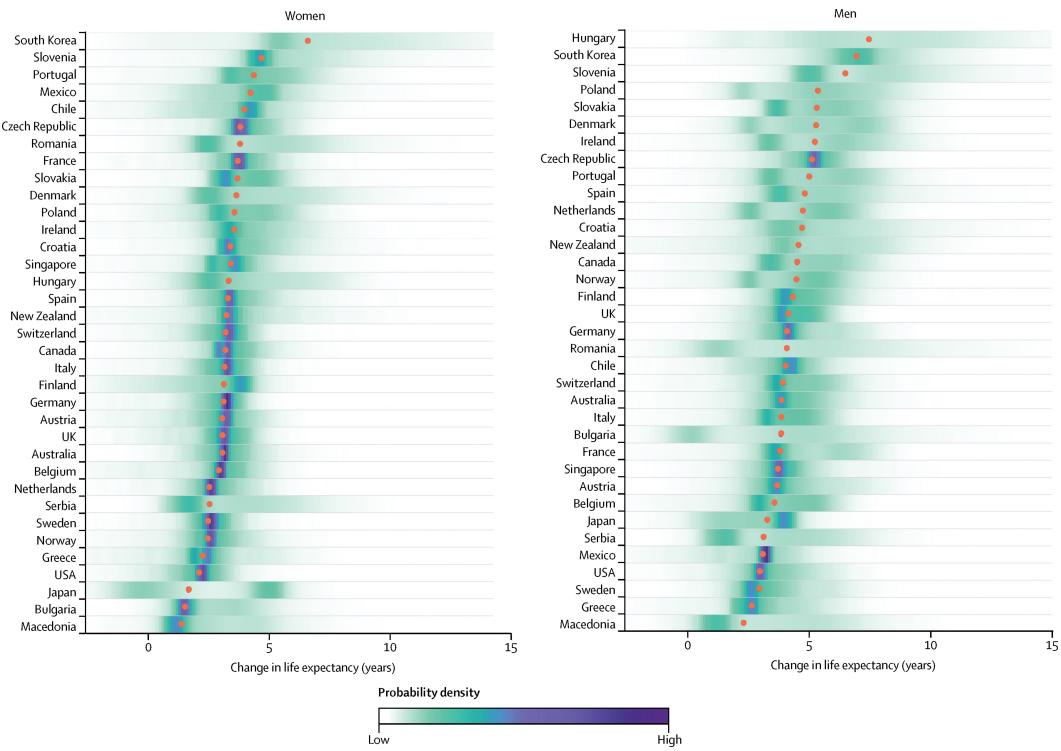
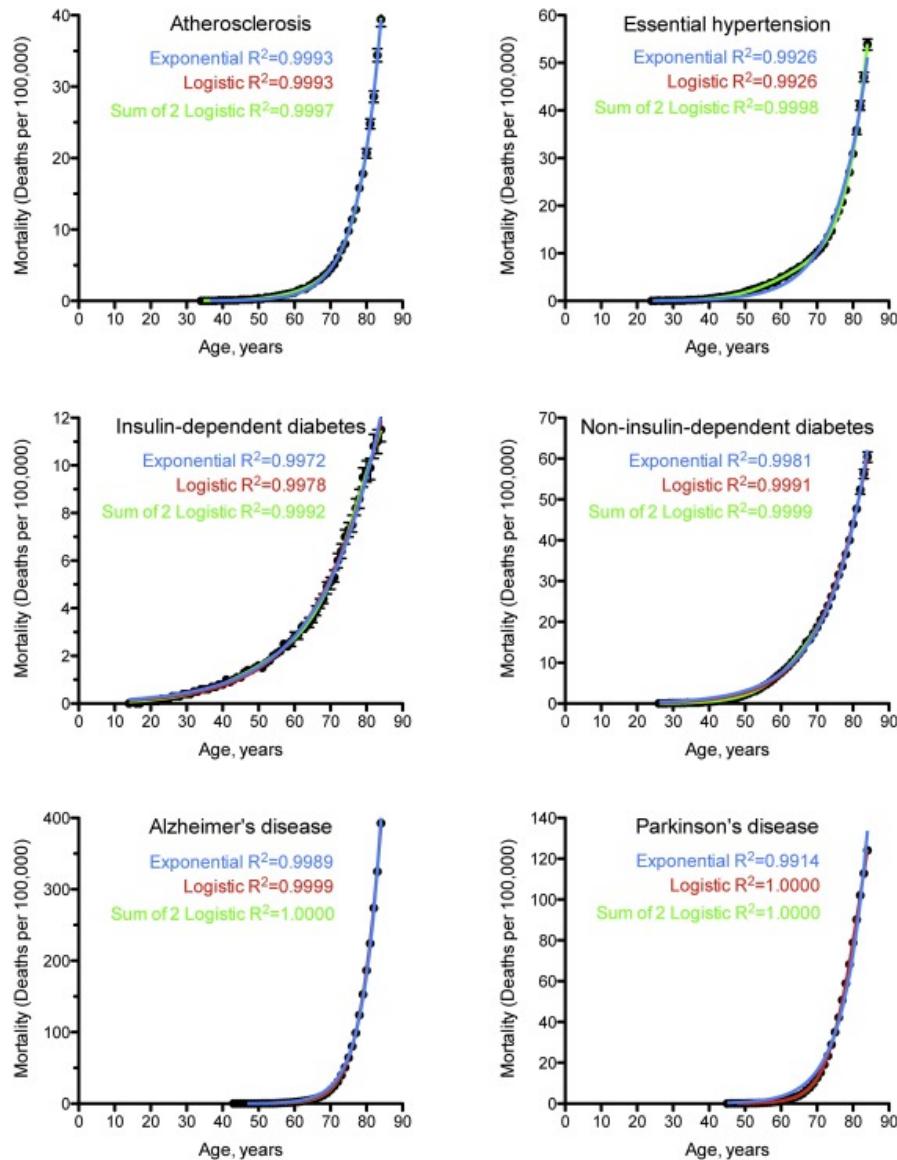


Figure 1.2: **The Population is Likely to Continue Getting Older** Median projected increase in life expectancy from birth in the period 2010 to 2030 for 35 industrial nations reproduced from Kontis et al. 2017 figure 1 [2].



**Figure 1.3: Mortality from Age-Related Conditions Increases According to a Logistic Function with Age** Mortality rate in deaths per 100,000 at different ages for Atherosclerosis, Essential Hypertension, Insulin-dependent diabetes, Non-insulin-dependent diabetes, Alzheimer's disease, and Parkinson's disease. Showing Model fits for exponential, logistic and the sum of 2 logistic functions. Underlying Cause of Death data (1999-2015) were form the Centers for Disease Control and Prevention, Wide-ranging OnLine Data for Epidemiologic Research (CDC WONDER) database. Reproduced from Belikov 2019 figure 1 [4].

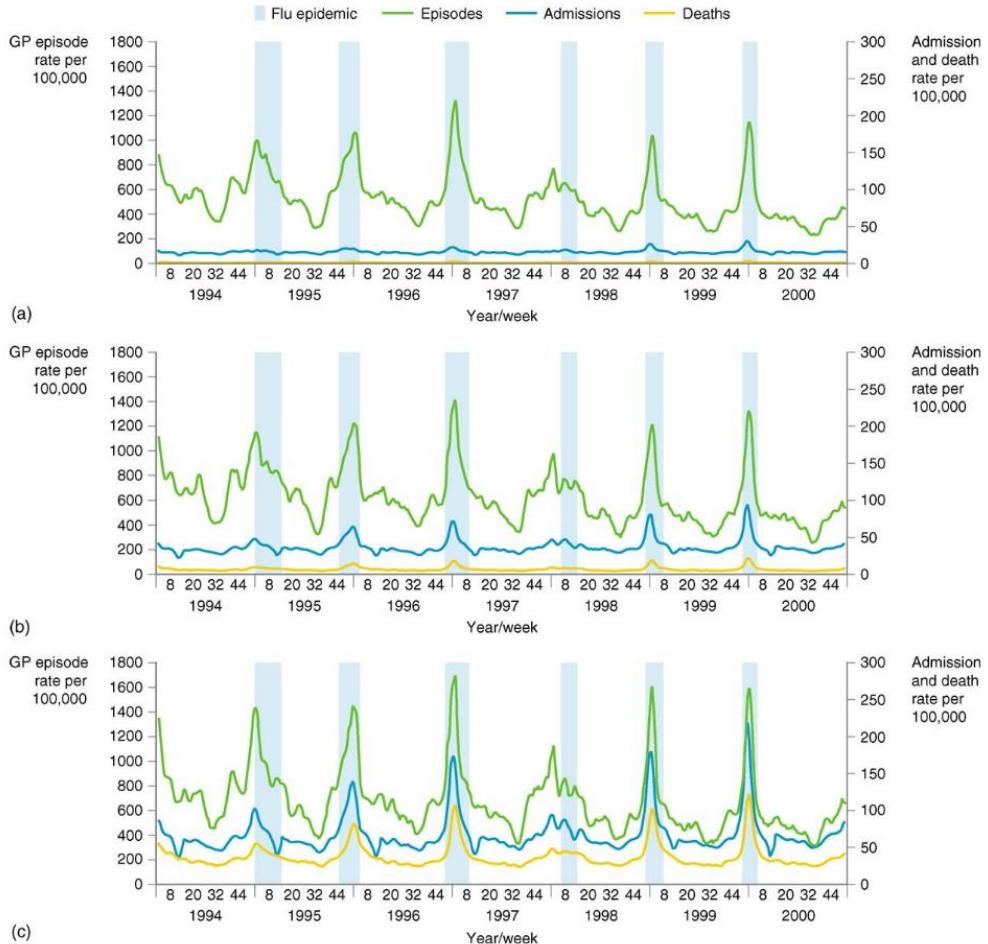


Figure 1.4: **Mortality from Infectious Diseases Increases With Age** Visits to General practitioners, hospital admissions and deaths from respiratory disease in persons aged: (a) 45-64, (b) 65-74 and (c) 75 years and over. Reproduced from Fleming et al. 2005 figure 3 [5].

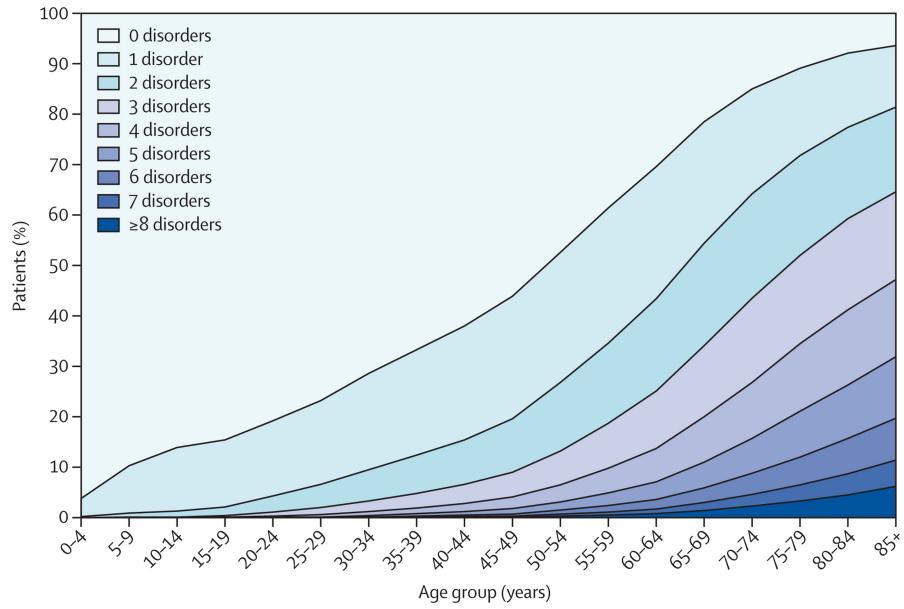


Figure 1.5: **Number and Frequency of Comorbidities increase with Age** Number of chronic disorders by age group. reproduced from Barnett et al. 2012 figure 1 [6].

Given the rapid increases in the rates of numerous diseases and the rate of death from those diseases as we age, there is a clear impetus to develop interventions to prevent and treat age-related diseases. The common factor among all these conditions is their relationship to age, but age is not merely a shared risk factor but frequently the largest; outside of rare genetic conditions which predispose to early-onset [3]. Thus understanding the underlying biological mechanisms of the ageing process and intervening in those mechanisms, has the potential to be vastly more effective and economical [8] than attempting to tackle diseases individually. The discovery in 1993 by Cynthia Kenyon and others [9] that a single mutation in the *daf-2* gene of *C. elegans* could double the lifespan of this organism is often cited as an inflection point in the confidence of the ageing research field that interventions in the process of ageing itself may be a realistic therapeutic target. What then do we know about the underlying biology of ageing?

## 1.2 The biology of Ageing

Ageing is characterised by a progressive deterioration of physiological integrity with time, resulting in an increased risk of mortality and morbidity. Theories of how and why this occurs are numerous with a history in the scientific literature stretching back to the late 1800s [10], setting aside, for now, theoretical frameworks for the process of ageing we will first consider the nature of these changes. Changes at the organismal level have their origins in changes at the molecular and cellular scales. Lopez-Otin et al. 2013 [11] identified nine hallmarks of ageing in their 2013 review. These hallmarks meet the criteria of: Manifesting during normal ageing; Accelerating ageing when aggravated; Retarding ageing when ameliorated. The hallmarks are further subdivi-

vided and hierarchically arranged into primary, antagonistic and integrative categories. Primary hallmarks are of unequivocally negative effect. Antagonistic hallmarks arise as a response to the primary and are initially protective against these changes but become problems in their own right as they occur at higher levels over time. Lastly, the integrative hallmarks directly impact on tissue function and homoeostasis. The hallmarks are as follows:

### **Primary Hallmarks - Causes of Damage**

1. *Genomic instability*: Damage to the genomes of our cells occurs over time and accumulates as we age, this damage ranges from point mutations to large structural mutations, as well as mutations in the mitochondrial genome.
2. *Telomere attrition*: Telomere shortening occurs during the ageing process, somatic cells do not generally express telomerase so with each division telomeres grow shorter. Telomere exhaustion limits the replicative capacity of cells inducing replicative senescence, limiting regenerative potential but also protecting from the uncontrolled growth of cancer.
3. **Epigenetic alterations**: Patterns of DNA methylation and of a number of histone modifications exhibit age-related changes. These contribute to the dysregulation of gene expression and derepression of heterochromatin regions which also increases the risk of genomic instability.
4. *Loss of proteostasis*: Production of protein folding chaperones in response to stress decreases with age, as does the activity of the two major proteolytic systems, the ubiquitin-proteasome system and the autophagy-lysosome system. This leads to increased prevalence of misfolded, aberrant and aggregated proteins with age.

### **Antagonistic Hallmarks - Responses to damage**

5. *Deregulated nutrient sensing*: The IIS (insulin & insulin-like growth signalling) and TOR (Target Of Rapamycin) pathways signal nutrient abundance, promoting anabolism. The AMPK (5' adenosine monophosphate-activated protein kinase) and Sirtuins signal nutrient scarcity, promoting catabolism. Broadly speaking inhibiting the pro-anabolic pathways and activating the pro-catabolic pathways can extend lifespan.
6. *Mitochondrial dysfunction*: Mitochondrial ATP production drops off with and ROS (Reactive Oxygen Species) production increases. Elevated ROS do not appear to accelerate ageing nor do elevated antioxidants retard it. Indicating the primary function of ROS in ageing is not straightforwardly a result of oxidative damage, but may result from its role in stress signalling. Dysfunctional mitochondria exhibit an increased propensity to polarise in response to stress, impacting on apoptotic signalling and inflammatory responses. This contributes to the increases in senescent cells and their chronic inflammatory phenotype.
7. *Cellular senescence*: Senescent cells are in stable cell-cycle arrest induced by shortened

telomeres, other forms of DNA damage, INK4/ARF derepression, or a variety of mitogenic/oncogenic signals. Induction of a senescent state in pre-cancerous cells protects against cancer and senescent cells are initially cleared by the immune system. However, these cells accumulate in tissue as we age without contributing to function and with a pro-inflammatory secretory phenotype that promotes chronic inflammation [12].

### **Integrative Hallmarks - Proximal Causes of the Ageing phenotype**

8. *Stem cell exhaustion:* Adult stem cells capable of producing new cells for the regeneration of tissues undergo asymmetric divisions to replace themselves and produce replicative progenitor cells to replenish tissue. Stem cells may accumulate damage and cease replicating or leave their quiescent state becoming overly replicative and potentially becoming senescent. This disrupts the supply of new cells for the renewal of tissue.
9. *Altered intercellular communication:* The aforementioned cell-autonomous alterations alter external signalling behaviour of cells across intercellular signalling paradigms. This can lead to feedback loops that exacerbate the problem. The high background of pro-inflammatory signals from senescent cells dilutes already declining immune response and fewer senescent cells are cleared.

This work and therefore the following review of the literature focuses on the 3rd hallmark, epigenetic alterations, specifically the role of DNA methylation in ageing. Booth & Brunet [13] contend that epigenetic processes are a hub through which all of the other hallmarks of ageing are mediated and feed-back on one another (Figure 1.6). One of these interactions is between DNA damage and chromatin state including DNA methylation. Sinclair and Oberdoerffer [14] make the case that DNA damage and repair processes are disruptive to the epigenome and, with Hayano et al., link double-stranded DNA breaks to accelerated epigenetic ageing in a recent pre-print [15]. Kane & Sinclair [16] make the case that reprogramming of the epigenome to a ‘younger’ state is a promising mode of intervention in the ageing process. Whilst loss of genetic information is effectively irreversible, barring gene therapy-like interventions, loss of epigenetic information is not necessarily so. Dysregulation from the loss of epigenetic information is arguably proximate to age-related changes and may run ahead of significant genetic information loss such that intervening in the former may have beneficial effects in its own right and stave off the latter.

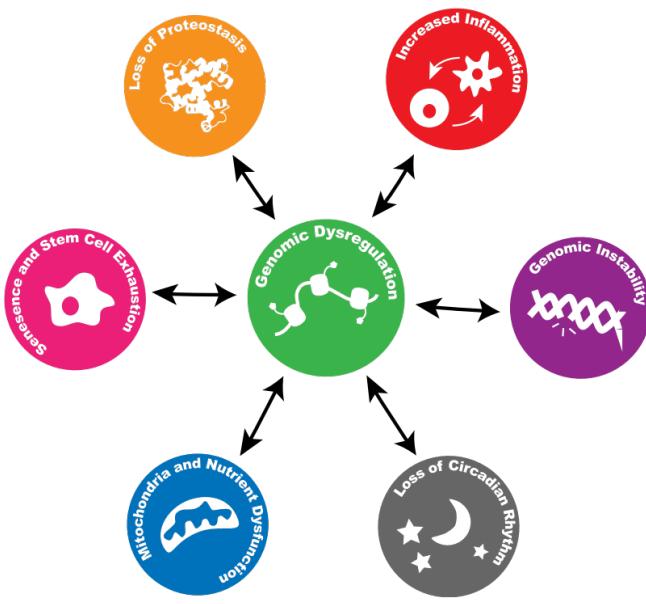


Figure 1.6: Epigenetic Changes as a hub for the hallmarks of ageing. (Reproduced from Booth & Brunet [13], figure 3)

Prior to examining the specific relationship of epigenetic changes to ageing we take a step back to define and outline epigenomics and examine the fundamentals of DNA methylation biology.

### 1.3 Epigenomics - Overview

Epigenetics generally refers to modifications to DNA and chromatin which do not affect the primary sequence of DNA bases [17], but which are to varying degrees stable and heritable. The term derives from epigenesis and genetics originating with Conrad Waddington [18]. Epigenesis refers to the idea that organisms develop through the progressive differentiation of cells from the egg into adult tissues. As modern genetics revealed that every cell contained a complete copy of the genome which was differentially utilised by the cells of adult organisms, the terms were fused to reflect the study of how this occurs. The usage has evolved a little further since molecular biology began to elucidate the mechanisms involved in this process and the term now frequently refers to the study of these mechanisms and their effects in less explicitly developmental contexts. Arthur Riggs et al. considered that heritability should be a criterion for a mark to be considered epigenetic, but this excludes many phenomena now commonly referred to using the term [19]. Requiring heritability results in a further definitional dispute over degrees of heritability; mitotic, meiotic, intergenerational, transgenerational and to what degree of fidelity? Adrian Bird proposed the definition: “the structural adaptation of chromosomal regions so as to register, signal or perpetuate altered activity states.” as a useful compromise [17]. Epigenomics refers to the

totality of the epigenetic modifications present in a particular cell, tissue type or genome. This definition appears to have caught on with at least some as: ‘These days, “epigenetics” more generally refers to all non-genomic information storage in cells including gene networks, chromatin structure and post-translational modifications to histones’, according to Alice Kane and David Sinclair [16].

Whilst an organism can generally be thought of as having a single genome, with exceptions such as Somatic mutations and Chimerism, it will have at least as many epigenomes as it has cell types [20]. There are on the order of  $10^{13}$  cells in the human body [21]. Estimates of the number of cell types vary with the resolution at which one deems cells functionally distinct [22]. At this point in time, however, there is no clear definition of what constitutes a distinct cell type, indeed they are beginning to be defined by the distinct patterns of gene expression and epigenetic modifications they exhibit [23]. Consequently, the number epigenomes that could be considered distinct will likely ultimately approximate to the number of epigenomes which can usefully distinguish between sub-populations of cells. A consortium has been established to produce a Human Cell Atlas [22] which aims to define all human cell types, and a cell type ontology [24] exists.

There are a number of epigenetic modifications which can be subdivided into four broad categories:

### 1. DNA modifications

In mammals DNA methylation primarily takes place on Cytosine residues. Cytosine methylation occurs principally in an mCpG sequence context but can also occur in mCpH, mCHG and mCHH (H=A/C/T) contexts; particularly in cells of the nervous system [25] and embryonic stem cells where as many as 25% of all cytosines can be methylated in non-canonical contexts. In contrast to methylation in a differentiated non-neuronal tissue such as foetal lung fibroblasts which is 99.98% in CpG context [26,27]. 5-methylcytosine (5mc) can be oxidised to produce another modified DNA base 5-hydroxymethylcytosine (5hmC) [28,29]. 5hmC may act as an intermediate to demethylation and potentially has regulatory functions in its own right [30]. 5hmC can be oxidised further to formylcytosine (5fC) and carboxylcytosine (5caC) [31]. The extent to which these modifications are stable and functional is still being explored. Other DNA bases can be methylated such as N6-methyladenine, but this occurs at a substantially lower frequency than 5mC and much less is known about their potential functions [32]. N6-methyladenine may play a role in condensing chromatin [33].

### 2. Histone tail modifications

Histone proteins form disk-shaped octamers around which ~150bp of DNA can be wrapped to form a nucleosome. The ‘tails’ are generally the N termini of the histone proteins, outside

of the core globular domains, which protrude from the nucleosome structure [34]. Histone tails can be subject to post-translational modification which commonly takes place at lysine residues, other residues are subject to modifications but lysines are among the best characterised. Examples of modifications include methylation with between one and three methyl groups per lysine, Acetylations, Ubiquitylation, SUMOylations, and phosphorylations [35]. Many of these marks can be generally classified as permissive or repressive but are frequently found in combinations of marks of opposing effect, rendering the interpretation of the ‘Histone Code’ extremely challenging [36–38]. One approach to interpreting chromatin state is segmentation. Segmentation makes use of pattern discovery techniques to divide the genome into discrete sections assigning these sections to a set of categories which can then be compared to existing annotations to provide functional labels [39,40]. When successful the pattern discovery algorithm independently recapitulates our existing ontology of functional elements and hopefully provides new insights by for example highlighting regions not previously considered to be in a given functional category. This is a “top-down” approach making use of the “sum” of available chromatin state data rather than trying to interpret the significance of a single type of mark.

### 3. Histone Variants

Alternatives to the canonical histones can be substituted in the nucleosome, examples of such alternative histones include H2A.Z and H3.3 [41]. Alternative histones can alter chromatin structure and dynamics by altering nucleosome stability, binding different factors, and presenting a different substrate for histone modifications, thus altering their pattern [42].

### 4. Some non-coding RNAs (ncRNA)

A subset of long ncRNAs, particularly those which persist in close association with chromatin can fall within the definition of epigenetic features [43]. An example of such a long non-coding RNA (lncRNA) is Xist which is a key regulator of X inactivation [44]. The categorisation of these as epigenetic is somewhat controversial, as they are not direct chemical modifications to chromatin. However Xist, for example, facilitates the mitotically heritable repression of the inactive X chromosome satisfying other criteria for being considered an epigenetic modification.

## 1.4 Fundamentals of DNA Methylation

### 1.4.1 Structure and Context

DNA methylation is the most well-studied epigenetic modification [45]. Specifically CpG methylation, the addition of a Methyl group to the 5 carbon of a cytosine base (figure 1.7), in a

CpG dinucleotide context (figure 1.8). (hereafter ‘DNA Methylation’ or ‘DNAm’ refers to 5mC CpG methylation unless stated otherwise.) DNA methylation is stable and relatively easy to characterise from biological samples by comparison with other epigenetic marks. The methylation status of ancient DNA has even been characterised, for example, native American remains ranging in age from 230 to 4500 years were measured using bisulfite sequencing [46].

5mC is quite chemically stable such that for most tissue sample preservation techniques for which DNA can be extracted from the sample DNAm can also be characterised [47]. DNA extracted with normal laboratory methods can be assayed for DNAm. Histone modification status, on the other hand, is less robust against environmental stresses on samples and often requires considerably more laborious sample preparation to characterise [48]. DNA methylation also offers the appearance of relative simplicity in comparison to the complex picture of the ‘histone code’, the simple binary nature of DNA methylation state makes it easier to model and possibly to interpret. However, interaction between DNAm and histone modifications is well documented [49–51], this interplay between DNAm and histone modifications renders the interpretation of DNAm more complex, as its effects may be conditional on the chromatin environment.

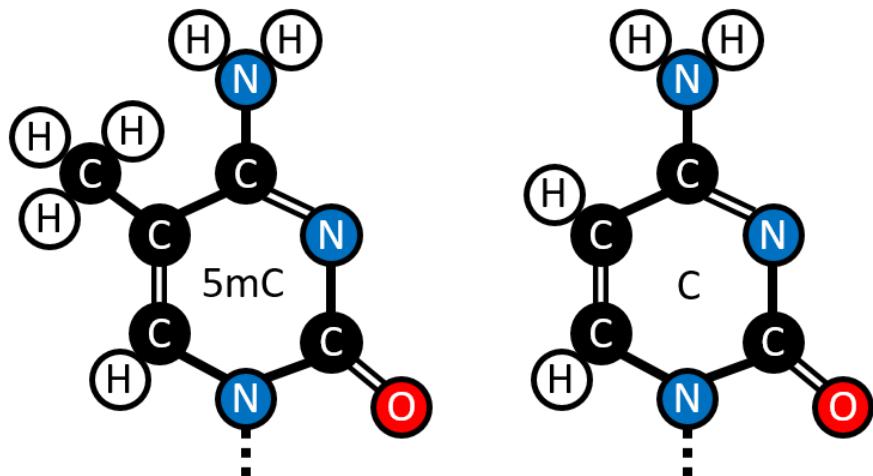


Figure 1.7: The structure of 5-methylcytosine contrasted with Cytosine. (Figure created by the Author.)

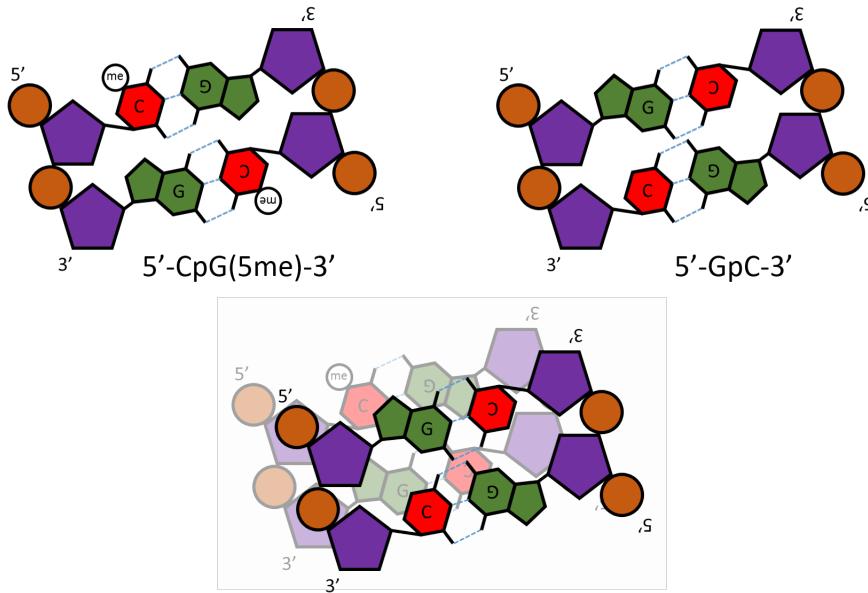


Figure 1.8: A schematic representation showing a Methylated **CpG** dinucleotide in its structural context. This is contrasted with a **GpC** dinucleotide, which cannot be superimposed on a **CpG**. upper left: methylated cytosine residues in a **CpG** dinucleotide, showing 5mC is 5' of G on both DNA strands. upper right: unmethylated cytosine residues in a **GpC** dinucleotide, showing G is 5' of C on both DNA strands. lower: superimposition of a **GpC** and **CpG** dinucleotide illustrating the mirrored nature of these structures and the impossibility of superimposing them through rotation. (Figure created by the Author.)

#### 1.4.2 Distribution and Global Trends

There are 28,299,634 CpGs [52] in the hg19 assembly of the human genome [53]. Given that the GC content of the human genome is 42% the prior probability of getting a CpG dinucleotide is:  $0.21 \times 0.21 = 0.0441$ , 4.41%. CpGs represent  $\sim 1.8\%$  of the dinucleotides in the human genome ( $28,299,634 \div (3.23 \times 10^9 \div 2) \approx 0.0175$ ). Making them  $\sim 2.5$  fold less frequent than would be expected *a priori*. Methylated cytosines are prone to deamination to thymines, resulting in mismatch lesions [54]. This increased mutagenic potential generally means they are selected against, accounting for at least some part of their under-representation in the genome. The three major classes of repeat elements SINEs, LINEs and LTRs contain some 46% of all CpG sites, with a further 5% in other repetitive elements [52].

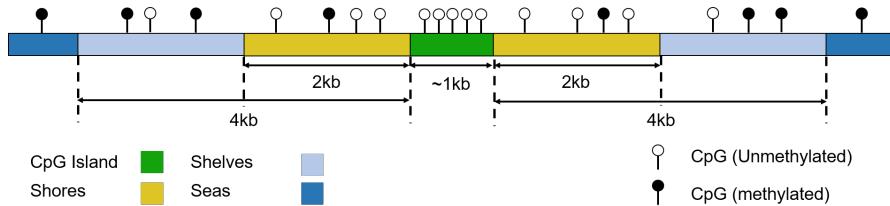


Figure 1.9: Schematic representation of CpG Islands, shores, shelves and seas. CpG density declines from shores to seas and methylation increases (CpG density and methylation proportion not to scale). (Figure created by the Author.)

CpGs are not uniformly distributed in the genome, they occur at higher frequency in some regions. “CpG islands” or CGIs are regions of high CpG density. The total number of CpGs in the UCSC repeat masked CGI annotation list of 28,691 CGIs is 1,990,729. Therefore, these CpGs comprise ~ 7.0% of the total number of CpGs in the genome. The mean percentage of the sequence of these CGIs that is comprised of CpG dinucleotides is ~ 18.5% and their mean length is 761bp. CGIs overlap the promoters of ~ 70% of genes [55], this, however, leaves roughly half of all CGIs as “orphans” not associated with a known transcription site. These orphan CGIs frequently constitute cell-type specific enhancers [56] and alternate promoters [57,58]. They are found in both intergenic and intragenic regions. CpGs are distributed quite sparsely through the genome occurring at low density in “seas” and at increasing density in CGIs and their flanking regions, see figure 1.9. It is a characteristic of CGIs that they are generally unmethylated.

Early work showed 70-80% of CpGs are constitutively methylated [59,60] Stadler et al. produced a more detailed picture of the distribution of CpG methylation in mouse embryonic stem cells (ESCs) [61] see figures 1.10 & 1.11. Stadler et al. categorised ...Mouse embryonic stem cells hidden markov model [61]

Irizarry et al. found that the 2kb regions flanking CpG islands which they termed “CpG island shores” exhibited greater tissue-specific differential methylation than the islands themselves [62]. This nomenclature has subsequently been expanded further with “CpG island shelves” which are 2kb - 4kb from the CGIs, and “seas” referring to the rest of the genome, see figure 1.9. Ziller et al. identified ~5.6 million CpGs dynamically regulated across diverse cell-types regions, these clustered into ~716,000 differentially methylated regions. More than three quarters of these regions were under 1kb in size and located away from transcription start sites and >70% of which had methylation levels >75% [63].

Irizarry et al., Stadler et al. and Ziller et al. data indicate that the most dynamic DNA methylation changes are occurring at the level of small distal regulatory features.

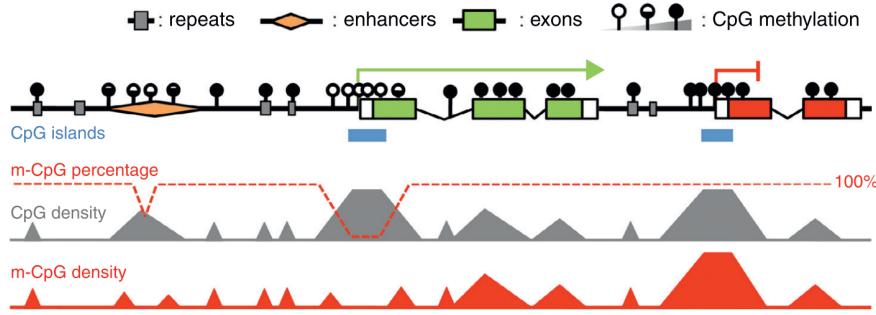


Figure 1.10: Stylised representation of the relationship between CpG Density, CpG methylation and functional DNA elements. (Figure reproduced from Baubec & Schubeler [64] figure 1)

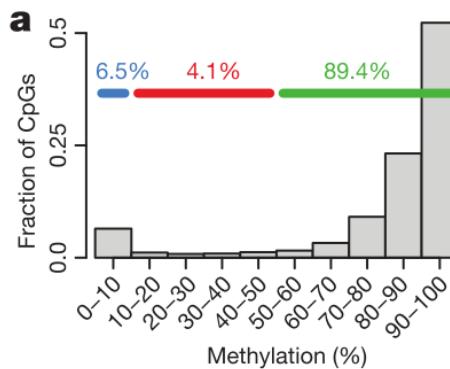


Figure 1.11: The fraction of CpGs with a given methylation level. Data from Whole-Genome Bisulfite Sequencing (WGBS) of Mouse embryonic stem cells. Colours denote: unmethylated regions (blue), low-methylated regions (red), and fully methylated regions (green). These groups are three segments produced by a hidden markov model. low-methylated regions with 10-50% methylation are evolutionarily conserved distal regulatory elements with dynamic cell-type specific regulation unlike the consistently unmethylated regions at CGIs. (Figure reproduced from Stadler et al. [61] figure 1a)

The regulatory role played by DNA methylation varies with context and the scale at which it is examined. For example methylation at a single CpG in a transcription factor binding site can determine whether or not the factor can bind at that site. There are methylation-sensitive transcription factors which can have their affinity for DNA either increased or decreased by methylation [65,66]. In addition, larger-scale changes in DNAm can through interactions with chromatin modifiers alter DNA compaction and more general accessibility of DNA regions for binding, and/or impacting on their topological organisation. DNAm also varies with nucleosome occupancy with lower methylation in linker sequences than on nucleosome-associated DNA [67,68]. The oxidation products of 5mC can also affect chromatin organisation, 5-Formylcytosine can impact nucleosome positioning through covalently binding to histones [69]. This raises the

question if some DNA methylation occurs specifically to act as an intermediate step for regulatory mechanisms which utilise its oxidation products.

### 1.4.3 Pathways of DNA methylation and demethylation

CpG methylation is produced and maintained by DNA methyltransferase (DNMT) enzymes. All the DNMT enzymes use S-adenosylmethionine as the source of the methyl donor group. DNMTs form a covalent intermediate between a conserved cysteine residue and the target base, through a nucleophilic attack on the C6 position in the cytosine ring [70]. This is followed by the transfer of the S-adenosylmethionine methyl group to C5, and deprotonation of the C5 to reform the double bond between C5 and C6; which is mediated by a base provided by the enzyme (Figure 1.12).

DNMT1 is associated with the replication machinery [71] and reproduces the methylation state of the parent strand on the daughter strand during replication. The largest of the DNMT family DNMT1 contains a ‘replication foci targeting sequence’ (RFTS) domain required for its targeting to replication forks. DNMT1 specialises in recognising hemimethylated DNA and methylating the unmethylated C in a palindromic CpG dinucleotide site, the UHRF1 protein assists in the recognition of these sequences [72]. DNMT3a and DNMT3b are responsible for *de novo* DNA methylation along with DNMT3L a catalytically inactive, but DNA-binding subunit [73]. The location of *de novo* methylation by the DNMT3s is influenced by a number of factors including the Chromatin state and other DNA binding factors [70].

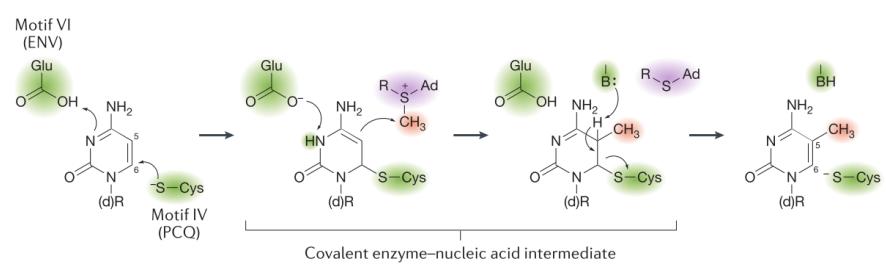


Figure 1.12: Motif VI (ENV/PCQ) refers to the conserved catalytic motif in DNMTs. ‘B:’ Represents the enzyme supplied base (Figure reproduced from the review of DNMT biology by Lyko [70] (fig. 1b))

DNMT2 is unlike the other members of the DNMT family in that it targets an RNA substrate. DNMT2 is a tRNA methyltransferase which methylates a specific subset of tRNA genes, mostly Asp isoacceptors, at a site adjacent to the anticodon which protects them from endonucleolytic cleavage under stress conditions. Unprotected tRNAs produce fragments which compete with small interfering RNAs interfering with their signalling pathways [74].

The mechanisms of the demethylation of DNA were more recently characterised than those

which govern its methylation [31,75]. They are more complex and less well understood. TET2 mutations are quite common in cancer, ranking 65th across 12 major cancer types [76] and more common in diverse myeloid malignancies, where loss of its catalytic activity favours tumourigenesis [77]. A distinction is drawn between active and passive demethylation, in passive demethylation, 5mC bases are diluted out in the process of DNA replication. In the leading model of active demethylation, they are oxidised one or more times by an enzyme from the TET (Ten-Eleven-translocase) family. They are then either passively removed by DNA replication or actively removed by a DNA glycosylase (thymine DNA glycosylase TGD) to create an apurinic site which is restored to a C by the Base Excision Repair (BER) pathway. This cycle of cytosine methylation and demethylation is illustrated in Figure 1.13 [78].

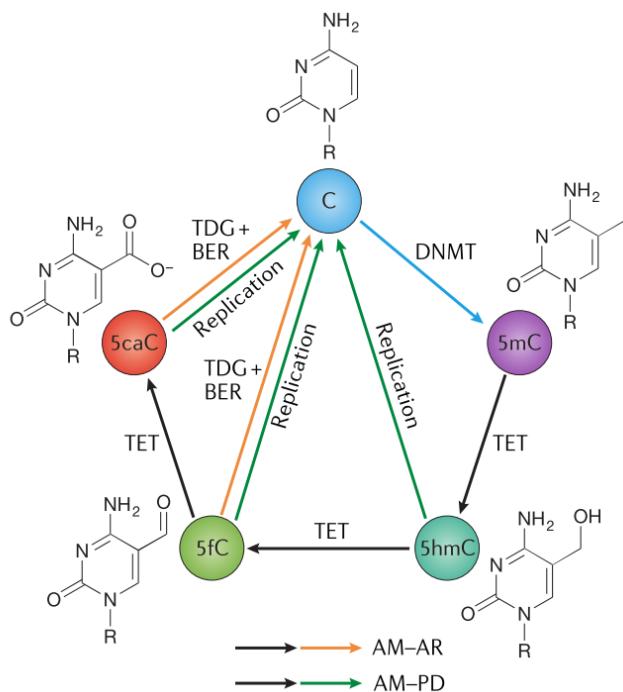


Figure 1.13: BER = Bases Excision Repair; AM = Active Modification; AR = Active Removal, PD = Passive Dilution (Figure reproduced from the review of TET mediated active demethylation by Wu & Zhang [78] (fig. 1a))

#### 1.4.4 Maintenance and Fidelity

DNA methylation is highly dependent on the underlying DNA sequence. Sequence features of particular importance to determining DNA methylation status are; transcription factor and other DNA binding protein recognition motifs, and CpG density [63,79]. Alterations in underlying DNA sequence such as SNPs and copy number variants (CNVs) can have a significant impact on methylation level and the susceptibility of the methylation level to change. In addition, CNVs can result in dosage effects on measures of DNAm, causing regions to appear, respectively,

substantially more or less methylated when fewer or greater copies than expected are present [80–85]. Despite the strong influence of sequence on methylation, global CpG methylation exhibits change over developmental time as illustrated in figure 1.14 and exhibits tissue-specific changes in distribution and amount [63].

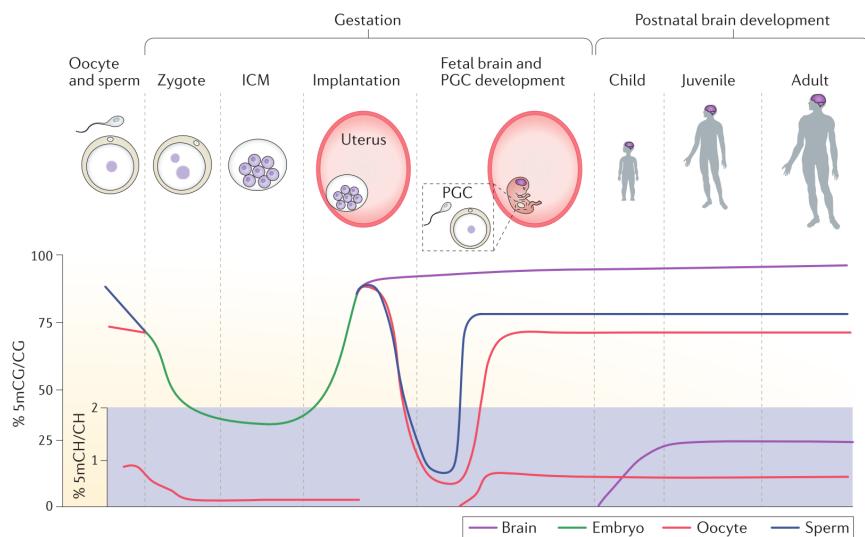


Figure 1.14: 5mC levels over developmental time. 5mCH or 5-methylcytosine-(A, T or C) levels are represented on a separate (purple) Axis from the main 5mCG axis. PGC = Primordial Germ Cell. (Figure reproduced from [86] figure 3)

In mitosis, DNA methylation is inherited by daughter cells with an error rate on the order of  $1 \times 10^{-3}$  per site per generation [87], several orders of magnitude less than that of DNA replication (error rate of  $1 \times 10^{-7}$  -  $1 \times 10^{-8}$  [88]). The fidelity of DNA methylation copying can be assayed by Hairpin-Bisulfite PCR (Polymerase Chain Reaction) [89]. Laird et al. looked at two alleles of a portion of the CpG island from the human *FMR1* gene in uncultured lymphocytes, one hypermethylated and one hypomethylated. In the hypermethylated allele, they found that 96% of sites methylated in the parent strand remained methylated in the daughter strand and 86% of unmethylated sites remained unmethylated. By contrast in the hypomethylated allele, there were no methylated sites to be retained and >99% of unmethylated sites remained so following replication. Laird et al.'s work suggests a higher degree of overall methylation fidelity for hypomethylated DNA and a propensity for unmethylated sites in hypermethylated DNA to become methylated.

Change in DNAm levels over time (divisions) can be modelled using the differential equations [90], which predict that a fully methylated site and a fully unmethylated site will converge on an equilibrium level. This level is determined by the probability of maintenance of the methylation state and of *de novo* methylation for a given locus. This stochastic model of DNAm is in agreement with experimental findings [89,91], Figure 1.15.

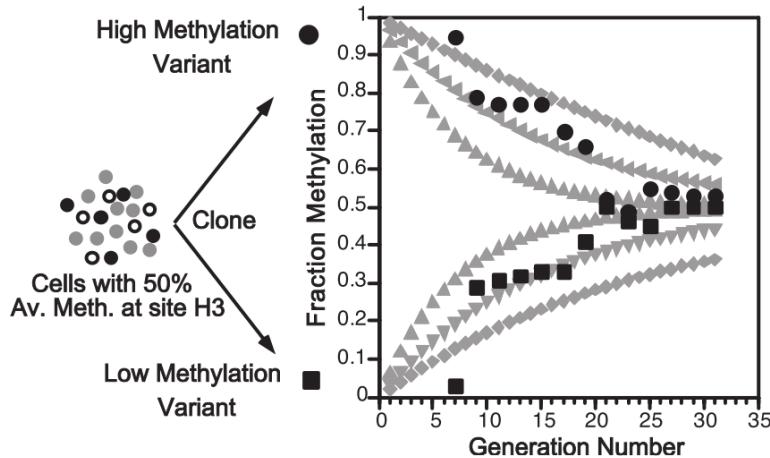


Figure 1.15:  $E_m$  = probability of methylation maintenance,  $E_d$  = probability of *de novo* methylation.  $M$  &  $U$  = the number of methylated and unmethylated molecules at specific CpG sites respectively. Modelling (**Grey points**)  $E_m, E_d$  values of 0.90, 0.10; 0.95, 0.05; 0.97, 0.03 (outermost to innermost, two curves for each  $E_m, E_d$  set, one starting at  $M = 1, U = 0$  and one at  $M = 0, U = 1$ ). Experimental data (**Black points**) from the *HpaII* locus (site H3) in 17 clones of mouse cell line BML-2 which has a known methylation level of 50%. (Reproduced from Riggs & Xiong [91])

Jenkinson et al. [92] used an information-theoretic approach, modelling DNA methylation as a binary communications channel using a 1-dimensional Ising model from the field of statistical physics. This permitted them to examine properties of DNA methylation not accessible to conventional means of analysis which typically capture the mean methylation level and perhaps the variability for a given locus. When considering methylation fidelity using this lens the maintenance of a given methylation state can be seen as an information processing task which requires the consumption of free energy in order to reduce the probability of error in the transmission of that methylation state. Thus Relative Dissipated Energy (RDE) can serve as a measure of the work expended by a cell at a given locus in order to preserve the current methylation state of that locus.

Zhao et al [93] applied Hairpin-Bisulfite PCR genome-wide in mouse embryonic stem cells (ESCs). They also found high degrees of methylation fidelity in hypomethylated regions such as CGIs and Promoters, as well as a high degree of fidelity in sites bound by transcription factors. This is in agreement with the findings of Jenkinson et al. [92] who noted that entropy (methylation stochasticity) was lower and more variable in CGIs and TSS (transcription start sites).

Methylation inheritance fidelity is reduced in cancer [87] and increases with differentiation [93]. The RDE (relative dissipated energy) at CGIs and TSSs is higher in differentiated tissues such as the brain, implying low entropy, and lower in embryonic stem cells, implying greater entropy [92] (see Figure 1.16). In addition, the correlation between CpG sites increases in cancer [92],

suggesting reduced higher-level regulatory control and tendency to fall back on lower level feedback mechanisms leading to a return to the baseline equilibrium described above as the cell is expending less energy to maintain an out of equilibrium methylation state.

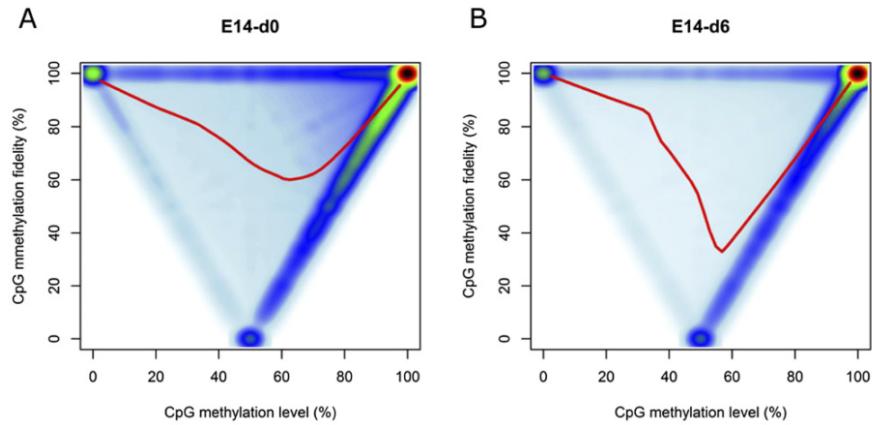


Figure 1.16: Bimodal Distribution of DNA methylation level and fidelity. Methylation fidelity exhibits a bimodal distribution with most and least methylated regions exhibiting the highest fidelity. Additionally, fidelity for methylation levels of 10-50% is considerably better than fidelity for levels of 50-90%. These data are from the mouse cell line ES-E14TG2a which is capable of self-renewal when cultured with Leukaemia Inhibitory Factor (LIF) and spontaneous differentiation upon removal of LIF (denoted as E14-d0 at day 0 and E14-d6 at day 6 after the withdrawal of LIF, respectively). Fidelity represents the percentage of symmetrically methylated or unmethylated CpG dyads for a given position as determined by Hairpin-Bisulfite PCR. (Reproduced from Zhao et al. 2014 [93] (figure 3 a and b).)

[94] late replicating domains

#### 1.4.5 DNA Methylation Assays

There are numerous DNA methylation assays. Some are designed to assay global DNAm levels producing a single global measurement of the amount of 5mC present in a sample, others assay DNAm in a manner traceable to a specific genomic locus. These methods can be further subdivided into targeted methods where the loci to be examined are known in advance, and those where the loci covered are quasi-randomly sampled. Both the targeted and untargeted methods vary in their granularity, from individual CpGs to large regions of the genome, on the order of megabases. These assays also vary in their coverage of the genome from locus-specific approaches looking at small numbers of individual loci to whole-genome methods providing methylation information on all ~28 million CpG sites. Between these two extremes are ‘epigenome-wide’ methods which are focused with varying degrees of specificity on genomic regions of interest.

These DNA methylation assays generally exploit one of four following principles:

1. **Methylation sensitive restriction digestion** A variety of methods exploit enzymes which differentially cut methylated and unmethylated DNA. An example of such enzymes is: HpaII and MspI which are isoschizomers for the sequence: 5'-C<sup>^</sup>CGG-3' but only HpaII is methylation-sensitive and unable to cut the sequence if the internal cytosine is methylated. The methods which use such enzymes include: Restriction landmark genomic scanning, which gives a roughly megabase resolution indicators of DNAm levels by 2D electrophoresis [95]. MRE-seq, which makes use of methylation-sensitive restriction digests to enrich for unmethylated DNA [58]. This is then sequenced to provide the genomic location of this unmethylated signal. Reduced Representation Bisulfite Sequencing (RRBS), which uses methylation-sensitive restriction digestion to enrich for unmethylated sequences which are then subject to bisulfite sequencing, (see point 3 below).
2. **DNA methylation-sensitive binding of DNA by antibodies or other proteins selectively which bind methylated DNA** Anti-5mC antibodies were first used to isolate methylated DNA in 1985 and were subsequently paired with array-based technologies to assay DNAm levels at specific loci in 2005 [96,97]. This was followed by MeDIP-seq (methylated DNA immunoprecipitation and sequencing) [98,99], which enriches methylated DNA that is then sequenced and the number of reads mapping to a locus is indicative of the relative methylation level. MBD-seq uses a recombinant Methyl Binding Domain (MBD) protein to enrich for methylated double-stranded DNA, prior to sequencing [100]. Inferring the absolute methylation levels from the number of reads mapping to a locus enriched by one of these pull-down methods is complicated by variation in CpG density in the genome and the fact that CpG dense regions tend to have low methylation levels, and thus tend not to be pulled down by antibodies binding 5mC or MBD proteins. The best results for estimating absolute methylation levels with these methods come from pairing them with a complementary method to enrich for unmethylated sites such as MRE-seq (described above) [101]. Though relative methylation levels remain effective in identifying differentially methylated regions when using these methods alone [58,102].
3. **Sodium bisulfite conversion of unmethylated cytosines to uracil** Conversion of cytosine to uracil changes the base complementary to this site from G to A [103]. The basic chemistry of the conversion process is illustrated in Figure 1.17. This conversion can be detected with a variety of technologies such as: The Illumina bead chip methylation arrays; WGBS; and Targeted bisulfite PCR or sequencing to examine a small number of selected loci. It can also be performed on a larger number of loci with microfluidic multiplexing such as that provided by the Fluidigm access array [104]. It is worth noting that bisulfite conversion cannot differentiate between methylated and hydroxymethylated cytosines. Because 5hmC represents a small fraction of modified bases compared to 5mC many analyses have made the working assumption that unconverted bases are methylated,

but it is beginning to be recognised that distinguishing between the two may be biologically important especially in tissues where 5hmC is more prevalent such as neurons. A variant of bisulfite conversion, oxidative bisulfite conversion, exists which can permit 5hmC to be distinguished from 5mC [105]. Additionally, New England Biolabs has recently developed an enzymatic alternative to chemical bisulfite conversion [106]

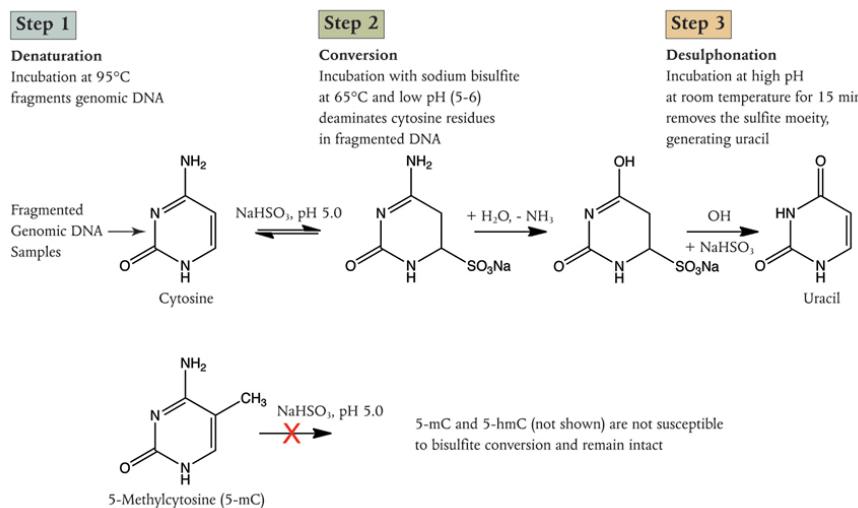


Figure 1.17: Sodium Bisulfite conversion of C to U, 5mC and 5hmC are unaffected. (Image Credit New England Biolabs).

**4. Physical differences in the methylated base** These are exploited by the not yet widely used ‘3rd generation’ sequencing technologies such as nanopore and single-molecule real-time sequencing. These methods respectively exploit the effects of modified bases on ion flow through the nanopore and impact on incorporation time of new bases whilst in the polymerase [107,108].

## 1.5 DNA Methylation and Ageing

Returning now to the relationship between DNA methylation and ageing we will review a little of the history of that field, the roles thought to be played by DNAm in ageing and the utility of DNAm clocks as biomarkers of ageing.

Early work on DNA methylation and ageing used biochemical, chromatographic and radiolabelling techniques to assay global changes in the amount of 5mC. These studies found differences in 5mC composition with cell type [59] but did not see changes with age. Other work around that time however found differences in 5mC levels with the age of cells in culture [109]. Wilson et al. noted that DNA methylation decreased across several tens of generations of cell lines in culture, but not in immortal cell lines. Immortal cell lines had lower absolute levels of 5mC to start with but remained constant over time. Wilson et al. also cite earlier work by Romanov at

al. 1981 [110] and Berdyshev et al. 1967 [111] which documented decreases in 5mC with age in cows and salmon, respectively.

Later work by Wilson et al. 1987 [112] in mice *P. leucopus* and *M. musculus* found decreases in 5mC with age and that the rate of decline was less in the longer-lived *P. leucopus* than in *M. musculus*. Interestingly, a recent study using different methods by Cole et al. found no global differences between young and old mice in short and long-lived strains, however, their other observations would seem to corroborate these initial trends. For the assayed sites long-lived mice exhibited 10x more hypermethylation than wild-type (WT), and WT mice had 3x more Differentially Methylated Regions significantly associated with age (DMRs/aDMRs) than did long-lived mice. Notably, the WT and long-lived mice shared many of the same aDMRs which differed in their degree of methylation as opposed to affecting different sites in the genome [113]. The lack of apparent global changes may be due to the biases of the reduced representation bisulfite sequencing (RRBS) method used. RRBS uses a restriction enzyme-based approach to enrich for regions with high GC content such as CGIs which tend to have low levels of methylation, and consequently may not be sensitive to loss of methylation in generally hypermethylated regions which could contribute to a global trend [114].

Wilson et al. also noted that the mitotic index of tissues did not relate to the loss of DNA methylation with age in tension with their earlier observations *in vitro*. The persistent loss of 5mC over time and the dramatic changes in methylation seen in cancer cells led Wilson et al. to suggest that dysregulation of DNA methylation may have a substantial role to play in the age-dependency of cancer risk and ageing more generally [115]. It has been suggested that “Epimutations” may be able to substitute for mutations in the multi-hit model of carcinogenesis [116]. For example; hypermethylation of the *BRCA1* promoter [117], Or, the development of Wilms’ tumour due to Beckwith–Wiedemann syndrome, a disorder arising from loss of imprinting of the gene encoding insulin-like growth factor 2 (IGF-2) leading to a double dose of IGF-2 protein [118]. Furthermore, the loss of methylation and increased entropy of closed chromatin regions may lead to increased susceptibility to structural mutations commonly found in cancer cells through the exposure of homologous sequences [92,115].

Methods which allowed the examination of changes in DNA methylation at known loci in the genome permitted a more nuanced picture of changes in DNA methylation over time to develop. Fraga et al. [119] introduced the concept of “epigenetic drift” being the divergence of DNA methylation and other epigenetic modification patterns with time. Fraga et al. looked at the divergence in DNA methylation along with global Histone H3 and H4 acetylation patterns between monozygotic (MZ) twins over a wide range of ages. They observed that older twins had greater epigenetic differences with time, in DNA methylation as well as H3, H4 Acetylation. Fraga et al. also noted that the divergence in epigenetic state was greater in twins who had lived longer apart and had different medical histories. Figure 1.18 is a useful visual encapsulation of Fraga

et al.'s DNA methylation results. Slieker et al. [120] identified 6366 CpGs whose methylation variability increased with age using the Illumina 450k array platform on whole-blood from 3295 individuals. Both the increase in variability with age and the increasing divergence of twins support a narrative of epigenetic dysregulation and increasing entropy with age.

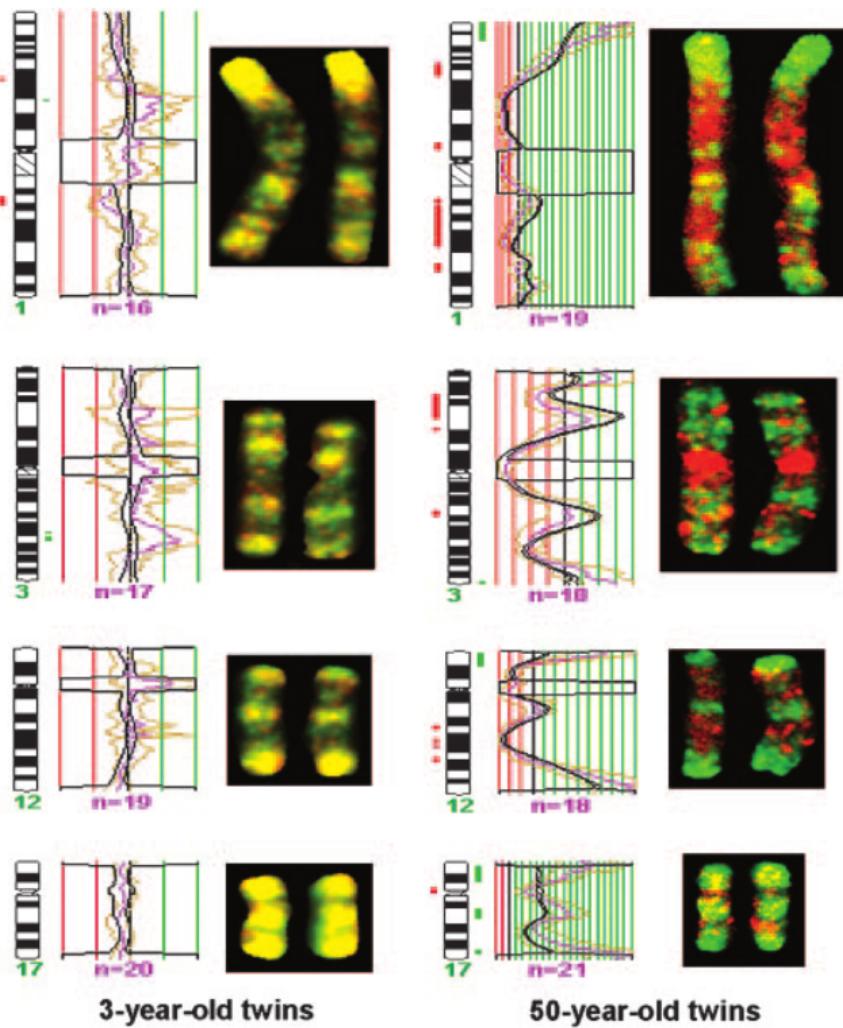


Figure 1.18: Mapping Chromosomal Regions of Differential DNA Methylation. Examples from Chromosomes 1, 3, 12 & 17 are shown for a 3 and 50-year-old twin pair. Methylation status is shown by competitive hybridization of AIMS (amplification of inter-methylated sites) products to metaphase chromosomes. Green and Red signals indicate hypermethylation and hypomethylation events between twins, Yellow indicates little difference. Red and Green blocks adjacent to ideograms indicate areas of significant DNA methylation change. (Reproduced from Fraga et al. [119] (figure 3).)

### 1.5.1 Age-Related Differential Methylation and Epigenetic Clocks

With the advent of array-based techniques which permitted the measurement of the methylation status of specific CpGs placed strategically throughout the genome, researchers were able to undertake a more fine-grained analysis of the changes in DNA methylation with age. This has allowed the prediction of chronological age from a metric of “DNA methylation age” (DNAm age), and interpretation of the differences in predicted and chronological age in terms of the pathophysiology of ageing.

The Illumina Golden gate array [121] with 1505 probes targeted to cancer-related genes was the first of these. Bjornsson et al. observed DNAm changes with age as well as noting substantial intra-individual differences. In contrast earlier work by Eckhardt et al. [122] and Ehrlich et al. [59] found no age-related changes using averages across individuals. Bjornsson et al. [123] also found that the intra-individual differences in differential methylation with ageing were highly heritable by use of familial clustering. Work by Boks et al. [80] also using the golden gate array corroborated Bjornsson et al.’s findings in monozygotic twins.

The Bjornsson and Boks studies used the golden gate array on peripheral blood samples, whereas Christensen et al. [124] used the golden gate array to look for age-related changes in DNAm in several tissue types. CpG sites that are differentially methylated with age when searching across tissues were frequently specific to a small number of tissue types. In addition, CpGs in CGIs tended to be hypermethylated with age and *vice versa*. Christensen et al. also found that the methylation profiles for different tissues were highly predictive of tissue type.

Rakyan et al. [125] looked at age-related changes in DNAm using the Illumina 27k array [126] with whole blood samples in a discovery set. They replicated their findings in sorted cell fractions to see if they could capture changes in DNAm that could be attributed to changes in blood cell-type composition over time, which they did not find. Teschendorff et al. [127] found that promoters of targets of the polycomb group proteins (PCGTs) were more likely than non-PCGTs to become methylated with age [128]. Repression of PCGTs is required for Stem cells to differentiate. PCGTs whose methylation status was associated with age were associated with pre-neoplastic conditions in a large cohort using the 27k array with blood and epithelial cell samples.

Bocklandt et al. [129] also using the 27k array but with saliva samples created the first “epigenetic clock” used to predict the chronological age of donors based on the DNAm landscape of their cells, the mean error was 5.2 years. Bocklandt et al. were able to achieve similar predictive accuracy with as few as 3 CpG sites in their model. Koch et al. [130] used publicly available 27k array datasets from a variety of different tissues to train a model using a different statistical method but only achieved an error of 11 years. They did, however, identify a number of CpGs also identified by Bocklandt et al. Bell et al. [131] used 27k array data to look for correlations between differential methylation and age-related phenotypes as well as chronological age. They

found that few age-related phenotypes were correlated with differential methylation, only 5 CpG sites were identified for the 16 age-related phenotypes examined. Whereas 490 significant CpGs were associated with chronological age. Many of the CpGs they identified persisted across tissue types and replicated in a second cohort, several had been identified previously by Rakyan et al. [125] and Bocklandt et al. [129].

Heyn et al. [132] examined the DNA methylation status of a newborn and a centenarian in CD4+ T cells using whole-genome bisulfite sequencing (WGBS), along with a group of newborns and nonagenarians using the Illumina 450k array [133]. Heyn et al. observed a global decrease in methylation from newborn to centenarian, as well as an intermediate level of methylation on a sample of intermediate age (see figure 1.19). This observation was replicated in 450k data. In addition, Heyn et al. noted that adjacent CpGs normally exhibit a substantial degree of correlation in methylation status and that they were less well correlated with their neighbours with increasing age. Interestingly Jenkinson et al. [92] found that correlation among nearby CpGs increased in cancer tissues. Modelling work done by Affinito et al. [134] agrees with this correlation between neighbouring CpGs under physiological conditions, which is particularly pronounced in CpG dense regions like CGIs. The greater physical proximity of CpGs in CpG dense regions means that greater correlation is expected. Due to the fact that these CpGs are more likely to be affected by the same proteins and regulatory features than CpGs with greater distance separating them [135]. Garagnani et al. used 450k array data from a cohort of 64 subjects aged 9-83 to identify those CpGs most well correlated with age as had been done with previous array technologies [136].

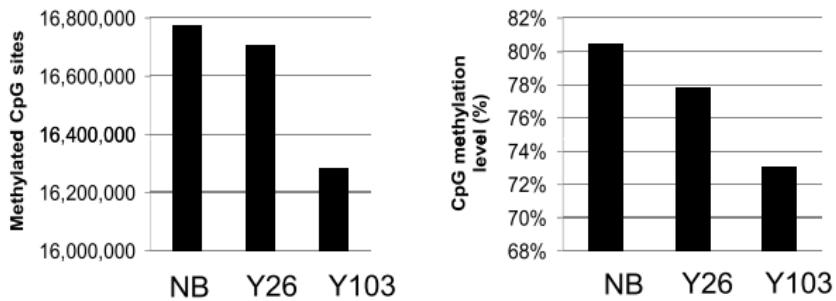


Figure 1.19: Global Hypomethylation with Age. NB = Newborn, Y26 = 26 year old, Y103 = 103 year old. (Reproduced from Heyn et al. [132] figure 1 b.)

In January of 2013 Hannum et al. [137] published a study using the Illumina 450k array and whole blood samples for 656 participants aged 19-101. Their optimised model was able to predict age with an error of 3.9 years. In December of the same year, Horvath [138] published a study using multiple datasets, including Hannum's data, some produced on the 27k and some the 450k array. Both used 'elastic net' penalised multivariate regression models to identify CpGs which

collectively provide the best predictor of DNA methylation age rather than the CpGs which are individually most highly correlated with age [139].

Florath et al. [140] identified some additional age-associated CpGs from two large cohorts totalling N = 898 and an N = 67 8yr follow-up longitudinal cohort. More than 3/4 of the CpG sites they identified began as hypomethylated and increased in methylation with age, a proportion likely skewed by the large number of 450k probes in CGIs and other typically low methylation regions. Bacalini et al. [141] performed a meta-analysis of existing DNA methylation datasets and employed a “region-centric” approach to try to identify loci larger than single CpGs which they anticipate will be more biologically meaningful than lone CpGs. They found that their approach increased the number of common features identified using the Hannum et al. and Heyn et al. datasets. Zaghlool et al. [142] performed 450k array DNA methylation study in Qatari population where they identified 12/88, 23/490 and 102/162 of the CpGs found by Bocklandt et al., Bell et al. and Florath et al. respectively. Zaghlool et al. also found that Horvath’s age predictor had an error of 3.7 years in their dataset. This is in agreement with existing findings and indicates that ethnicity has no major effects on the DNA methylation age signature. Benton et al. [143] looked at changes in DNA methylation with age in a genetically isolated population on Norfolk Island with similar results to previous studies but identifying some novel age-associated CpGs.

Most DNA methylation relationships with age noted in previous studies have been linear, Johnson et al. [144] used the 450k array in peripheral blood samples. They identified 21 CpGs whose DNA methylation changes at a rate that changes with age from an initial pool of 27,723 CpGs which were differentially methylated with age. Two sites exhibited an increasing rate of increase in DNA methylation with age, and 18 sites a decreasing rate of increase.

### 1.5.2 Genetic Influences on DNA Methylation

Epigenetic variation falls on a continuum of genetic influence that can be summarised by three categories:

1. **pure** - DNA sequence has no predictive value for epigenetic state.
2. **facilitated** - DNA sequence biases epigenetic state.
3. **obligatory** - DNA sequence permits exact prediction of epigenetic state.

Regional methylation state is strongly influenced by genotype, by single nucleotide polymorphisms (SNPs) [145] and by structural variants [84]. An example of an obligatory effect on methylation is a point mutation at a CpG site, a C to T transition precludes methylation at that site in future. The effect of structural variants on methylation can be hard to determine as changes in sequence dosage often lead to measurement artefacts. Efforts have been made to correct for the influence of genetic factors in EWAS, when searching for purely epigenetic effects but the potential interaction of the somatic mutations known to accumulate with age (Figure

1.20) and the changes in DNAm with age remain largely unexplored. This is of particular relevance to DNAm as the profile of the types of mutation which accumulate with age (Figure 1.21) distinctly favours C to T transition mutations which can disrupt CpG dinucleotides [146].

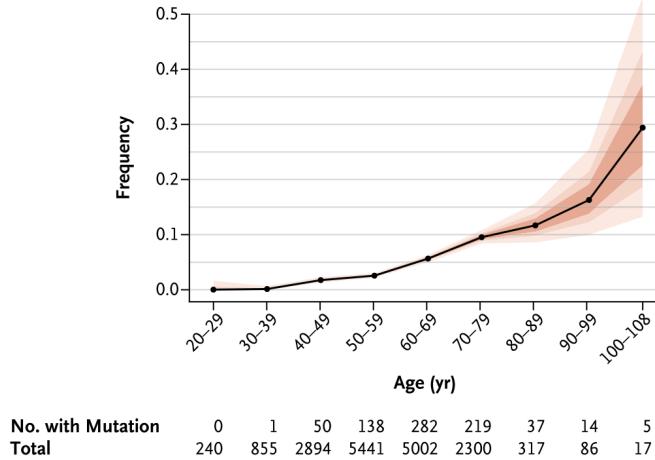


Figure 1.20: Somatic mutations increase with Age. (Reproduced from Jaiswal et al. [146], figure 1)

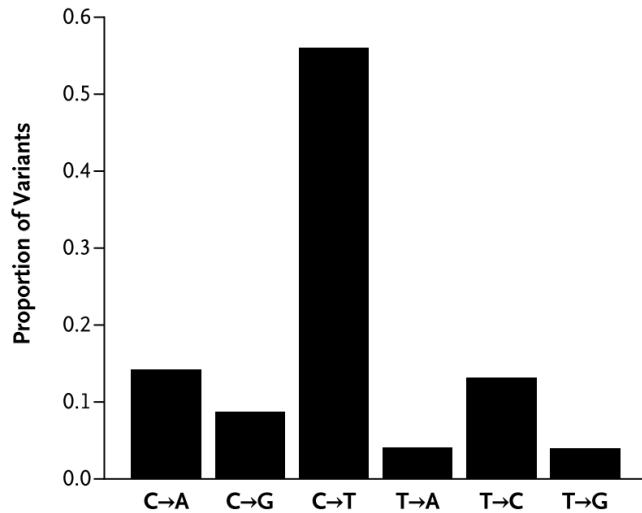


Figure 1.21: C to T transitions are the most common substitutions that occur with age. (Reproduced from Jaiswal et al. [146], figure2 c)

### 1.5.3 DNAm Age as a Biomarker of Ageing

We have seen that DNAm can be a fairly accurate predictor of age, however, there is more to being a good biomarker of ageing than predictive accuracy. Weidner et al. [147] conducted a study with the specific intent of identifying CpG sites that would serve as the best biomarkers of ageing. Below is a list of specific criteria for a high-quality biomarker of ageing laid out by

Johnson et al. [148].

1. It must predict the rate of ageing. In other words, it would tell exactly where a person is in their total lifespan. It must be a better predictor of lifespan than chronological age.
2. It must monitor a basic process that underlies the ageing process, not the effects of disease.
3. It must be able to be tested repeatedly without harming the person. For example, a blood test or an imaging technique.
4. It must be something that works in humans and laboratory animals, such as mice. This is so that it can be tested in lab animals before being validated in humans.

We have seen that DNA methylation clocks can be good predictors of chronological age, we will revisit the question of whether or not they are superior to chronological age in their ability to predict lifespan below. Bell et al. [131] found DNAm did not generally correlate well with ageing phenotypes but did correlate well with chronological age and Teschendorff et al. [127] found that their age-associated DNAm signature remained fairly constant across several disease states including ovarian cancer and type 1 diabetes. Suggesting that DNAm age prediction does indeed: “monitor a basic process that underlies the ageing process, not the effects of disease”, meeting criterion number 2. Several of the DNAm clocks discussed so far have been based on whole peripheral blood samples meeting criterion number 3. DNAm age clocks have been shown to work in Chimpanzees [138], Dogs/Wolves [149], Mice [150], Naked Mole Rats [151], Rhesus Macaques, Humpback Whales [152] and are likely to work in other mammalian model organisms. Furthermore, a cross-species clock based on conserved ribosomal RNA genes has been developed [153]. It should be noted that Horvath’s Human clock uses array data and that the Dog/Wolf and Mouse clocks use RRBS. Model organisms such as *C. elegans* which essentially lack DNA methylation would not, however, be likely to be tractable for DNAm based age prediction. Consequently DNAm age prediction at least partially meets criterion number 4.

Is DNA methylation a better predictor of lifespan than chronological age? Both Hannum and Horvath identified the difference between DNAm age and chronological age or  $\Delta_{Age}$  as a potential indicator of disproportionate biological ageing. Marioni et al. [154] explicitly set forth to test the value of DNAm age as a predictor of mortality. Adjusting for age, sex, childhood IQ, education, social class, hypertension, diabetes, cardiovascular disease, and *APOE* e4 status a  $\Delta_{Age}$  of +5yr is associated with a 16% increased mortality risk (See figure 1.22). A longitudinal study of twins found 3.2 fold increase in the risk of dying first per 5yr  $\Delta_{Age}$  within twin pairs [155]. A meta-analysis of DNAm age predictors by Chen et al. [156] also found that measures of age acceleration based on Hannum and Horvath clocks were superior predictors of mortality than chronological age before and after correction for various potentially confounding factors. Chen et al. also noted that correction for blood cell composition improved predictive power. Indicating that DNAm age is indeed superior to chronological age as a predictor of lifespan meeting criterion number 1 for an ageing biomarker.

A new DNAm age clock explicitly designed to capture phenotypic age and outperform age acceleration as a biomarker of ageing has recently been produced by Levine et al. [157]. The ‘PhenoAge’ metric was created by selecting nine biomarkers of ageing from 42 possible metrics using a proportional hazards penalised regression model and combining these with chronological age. (The nine biomarkers are: Albumin, Creatinine, Serum glucose, C-reactive protein, Lymphocyte per cent, Mean red cell volume, Red cell distribution width, Alkaline phosphatase, White blood cell count.) DNAm data was then regressed against PhenoAge using elastic-net regression to produce a DNAm based PhenoAge predictor which made use of 513 CpGs. The PhenoAge predictor outperformed the Horvath and Hannum clocks at predicting all-cause mortality, comorbidities, coronary heart disease risk, and measures of physical functioning. 41 of the 513 CpGs in the PhenoAge clock were present in the original 353 CpG sites used in the Horvath clock.

The ‘GrimAge’ clock [158] takes a different approach to previous epigenetic clocks in that it has a 2-tiered model predicting first seven surrogates for a variety of other biomarkers of ageing then fitting a model to predict time to death using these surrogates. ‘GrimAge’ outperforms age accelerations from the Horvath, Hannum and PhenoAge clocks at predicting time to death.

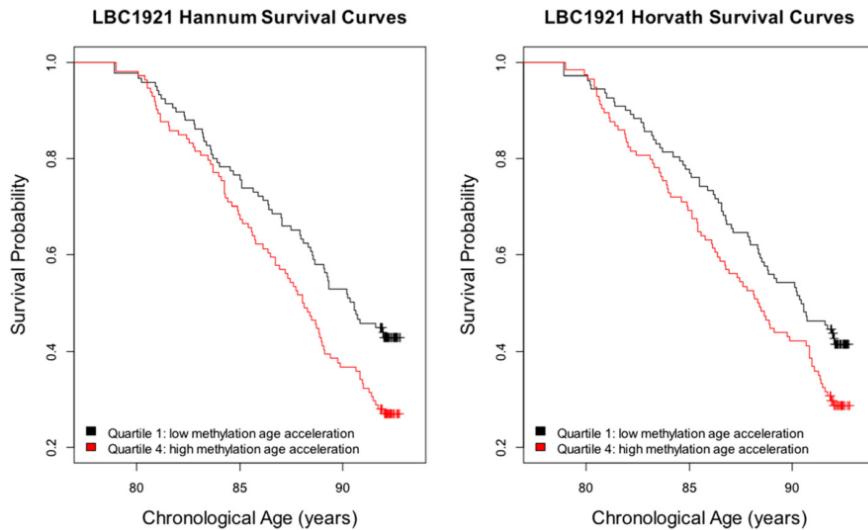


Figure 1.22: Survival probability by quartiles of  $\Delta_{Age}$  in LBC 1921 adjusted for sex, and chronological age. LBC = Lothian Birth Cohort. Using the Hannum and Horvath predictors’ values for  $\Delta_{Age}$ . (Reproduced from [154] figure 3.)

#### 1.5.4 Discussion and Analysis of the Findings of DNAm Age Studies

Interestingly CpGs whose methylation status has been identified as correlating with age in various studies show relatively little overlap. Despite this, changes in DNAm with age are probably the most robust and reproducible large-scale epigenetic change yet captured. This lack of agreement in absolute terms, however, presents a challenge about how to interpret DNAm changes with age.

The result of the various DNA methylation studies display some seemingly contradictory trends. For example, there is the global loss of methylation with age [112,119,132] and yet the majority of the highly significant age predicting CpGs in clocks are undergoing hypermethylation by a ratio of very roughly 2:1 [125,129,130,136,140] (see figure 1.23). Interestingly when expanded to look at all age-associated CpGs and not just the most highly correlated CpGs a slightly different picture emerges. Teschendorff et al. [127] noted in their original study and Zhang et al. [159] noted in a subsequent analysis if more loosely age correlated sites are included hypomethylation with age is more prevalent. Zhang et al. state roughly 60% of sites are hypomethylated and 40% hypermethylated, they also observed that all of the highly age-predictive CpGs which overlapped between the studies they examined were hypomethylated.

It is worth noting that the technologies used to measure the DNA methylation introduce their own biases into global changes in DNA methylation levels. The more Global measures such as WGBS and biochemical measures seem to favour hypomethylation with age [112,119,132]. The arrays seem to have agreed more with these global measures as they have increased in size. RRBS also seems to skew in favour of hypermethylation with age [113]. These more narrow technologies are biased towards promoters and CGIs with their higher CpG density and lower baseline methylation levels making the prior probability of observing a methylation event greater than demethylation by dint of having a dearth of methylated sites to start with.

Sites which are hypomethylated as their “ground state” are more likely to be in CGIs than sites which are hypermethylated as their “ground state”, and thus hypomethylated sites are more likely to be variable in their methylation state and subject to tight regulatory control than are hypermethylated sites. Expanding the pool to look at CpG sites more poorly correlated with age is likely to introduce more sites whose hypomethylation with age is due to increased stochasticity with age. However, from the background information, one would expect that hypermethylated sites in CpG dense regions like CGIs whose methylation level declines with age would make good quality indicators. Given that CpGs in CGIs have a high prior probability of being hypomethylated, it follows that one would be more likely to observe a hypermethylation event by chance. Consequently, a CpG in a CGI that exhibits hypomethylation with age is less likely to be a result of noise than its converse, this may explain the higher reproducibility of these hypomethylated sites.

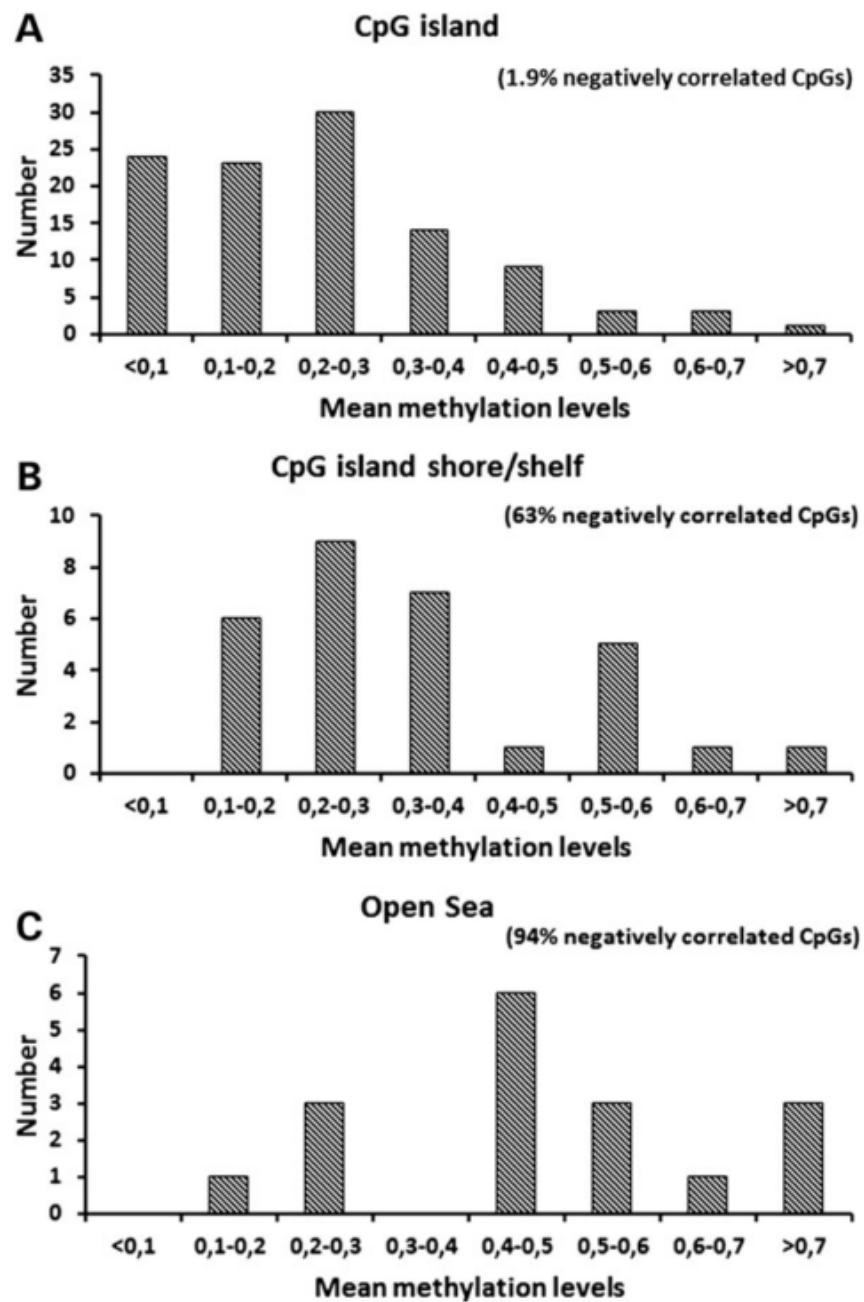


Figure 1.23: Mean methylation level by genomic region showing proportions negatively correlated with age, i.e. become hypomethylated as age increases. (Reproduced from [140] figure4 a, b, c.)

In a review of DNAm and ageing Jones et al. [160] drew a distinction between epigenetic drift changes and epigenetic clock changes with age wherein they defined drift as “the collection of DNA methylation changes that are associated with age within an individual but are not common across individuals.”, which is slightly different from the definition used by Fraga et al. [119]. Clock changes are defined as “those sites that are associated with age across individuals and can thus in some cases be used to predict chronological age”. Drift constitutes stochastic change over time

not driven by some process which strongly biases the resulting changes in DNAm in direction or to specific sites. Clock changes are the converse, some underlying process biases DNAm changes in direction and/or position.

This picture is complicated by the tissue specificity of DNAm changes. Whilst it is evidently the case that multi-tissue DNAm age prediction is possible, it also seems to be the case that an age predictor trained on samples from a specific tissue will outperform one trained on multiple or other tissues when predicting the age of a sample of matching tissue type. This is supported by the out-performance of the Horvath predictor by the Hannum predictor in blood samples [154] and differentially predicted ages of tissues such as breast noted by Horvath [138]. DNAm was highly correlated with age for a variety of tissue-specific sites in many studies, suggesting that the definition of clock sites is frequently tissue-dependent. This may also be true of drift changes, stochastic changes may occur at sites in one tissue that would be tightly controlled in another. If one was interested in universal signatures of ageing common across cell-types, one could restrict the definition of clock changes further to only those which are present in all/most tissue types to make a multi-tissue predictor.

Apparent changes in the DNAm state of a mixed cell population could be due not to global changes in that population but rather its composition, blood exhibits changes in the cellular composition with age [161,162]. As hematopoietic stem cells age their numbers increase and their regenerative potential declines. Their output becomes skewed towards the myeloid lineage but produces poorer quality cells. The number of naive T-cells decreases with age, memory and effector T-cells accumulate as to memory B cells limiting the diversity in the cell-types and contributing to immunosenescence [163]. Cellular composition of tissues can have a substantial impact on epigenome-wide association studies (EWAS), especially when dealing with small effect sizes [164]. Thus shifts in cell composition such as those in blood during ageing can be substantial sources of confounding when searching for cell-intrinsic changes in DNA. Some attempts to assay and/or correct for cell-type composition have been made, for example, Rakyan et al [125] isolated  $CD4^+$  T-cells and found that they had a 60% overlap with age-related CpGs in whole blood. Chen et al. [156] used the method developed by Houseman et al. [165] to correct for cell-type composition and found that it did improve the performance of the Horvath and Hannum age predictors. The changes were significant but the effect of this correction was relatively modest. Whilst cell composition has an effect, DNAm Age definitely seems to be capturing something beyond shifts in cell-type composition. However, the age acceleration computed with Hannum's method, which is trained on blood data only, was more strongly correlated with cell-type composition than Horvath's method, indicating that extrinsic changes in cell-type composition likely have some impact the model.

Jenkinson et al. [92] took an information-theoretic approach to DNA methylation maintenance, modelling the binary methylation state of CpGs as a noisy communications channel. Permitting

them to compute the information capacity (CAP), methylation entropy (ENT), and the relative dissipated energy (RDE) of these channels. High capacity channels represent reliable transmission of methylation state, consistent with high RDE as more energy is needed to ensure reliability and low entropy as methylation is more ordered at sites under tighter control than at less highly controlled sites. They found that the Transcription Start Sites (TSS) exhibited high levels of RDE and high information capacity as well as low levels of entropy. These values trended in the opposite direction as they moved outwards to CGIs, shores and seas. This information-theoretic framing provides a useful lens through which to view drift and clock loci.

It could be hypothesised that as a cell differentiates those sites the methylation status of which are most important for successful differentiation will be the focus of the expenditure of cellular resources to set and maintain their methylation state, and will thus exhibit high capacity and low entropy. Whatever mechanism is focusing the resources of the cell on the maintenance of the methylation at these sites may be responsible for the tissue-specific biases in DNAm changes over time i.e. their targets will be clock sites. Those sites not involved in the maintenance of the specific differentiated state would be a lower priority for the expenditure of cellular resources and thus display lower information capacity and higher entropy i.e. be drift sites. This would also lead one to predict that multi-tissue age predicting loci would be enriched for factors common to the maintenance of a differentiated state. The polycomb and Trithorax group proteins are very broadly evolutionarily conserved and play a key role in developmental cell fate choice and particularly in long-term stable maintenance of epigenetic state [166]. Consistent with this hypothesis Teschendorff et al. [127] observed that PCGTs are enriched among CpGs associated with age and Levine et al. [157] also noted that their 513 clock CpGs were enriched for PCGTs. This is also complementary to the idea that DNAm dysregulation with age promotes carcinogenesis by stabilising stem cell character but with reduced capacity for differentiation as proposed by Teschendorff et al. [127] and Rakyan et al. [125].

In his paper [138] Horvath proposed that DNAm age is a measure of the work done by an “Epigenetic Maintenance System” (EMS) which maintains epigenetic stability over time. Considering work in its physical sense of the integral of power over time, Horvath considered power, the rate of change in the use of energy by the EMS as the ‘tick-rate’ of the epigenetic clock. He proposed that during development power would be high due to the high level of stress on the system and that power would drop to a consistent level once maturity was reached. Horvath laid out four predictions of this EMS model and argued that there was evidence in support of each of them. The predictions were:

1. Cancer tissue should show signs of accelerated age, reflecting the protective actions of the EMS
2. Many mitogens, genomic aberrations, and oncogenes, which trigger the response of the EMS, should be associated with accelerated DNAm age.

3. High age acceleration of cancer tissue should be associated with fewer somatic mutations given the protective role of the EMS
4. Mutations in TP53 should be associated with a lower age acceleration of cancer tissue if one further assumes that p53 signalling helps trigger the EMS

A fifth prediction of this model would be that DNA methylation fidelity should be elevated during development and differentiation as a reflection of the greater amount of work done by the EMS during this period. This is borne out by the results of Zhao et al. [93] who noted a global increase in DNA methylation fidelity during the early differentiation of mouse embryonic stem cells (see figure 1.16).

If Horvath's concept of an EMS is accurate and Jenkinson et al.'s RDE captures the energy expended by the cell on the maintenance of the methylation state of a locus, one would expect Horvath's concept of an EMS to tie in quite nicely with Jenkinson et al.'s methylation channel model. With the Relative dissipated energy (RDE) for a CpG channel corresponding to the work done by Horvath's EMS at a given locus. Jenkinson et al. did indeed observe high levels of RDE in stem cells and a decline of RDE with age. Interestingly, however, some cancers exhibited decreased RDE, CAP and increased entropy which is consistent with an accelerated ageing profile but seemingly less so with elevated activity of an EMS. It is possible that whilst the distribution of RDE at various loci in the genome shifts giving a lower mean RDE that the total amount of energy expended is increasing but becoming less well targeted such that whilst there is increased energy flux through the EMS it is not being properly directed. There may also be a plateau effect wherein the energy flux through the EMS is maxed out and cannot keep up with demand from many loci. One would also predict more rapid equilibration of DNAm levels at loci with high RDEs, which could be examined by looking at the second derivatives for the stochastic model of DNA methylation described by Pfeifer et al. [90] using the experimental systems Laird et al. [89] and Riggs et al. [91] used which supported the original model.

Jenkinson et al. [92] found a global increase in entropy with age, but not with cell passage in culture suggesting an increase in entropy is associated with epigenetic age independent of mitotic age. This observation is in agreement with the finding that epigenetic state becomes more stochastic and diverges with age [119,120], indeed genes whose methylation were most divergent with age were enriched for ageing associations. They also observed a loss of entropic sensitivity with age. Entropic sensitivity is an indicator of how plastic DNAm state is to extrinsic effects. Jenkinson et al. noted that there is a general loss of phenotypic plasticity with age, but cited no specific instance of loss of DNAm plasticity to environmental effects with age. Hahn et al. [167] provided a possible example of this effect when they reported that dietary restriction caused fewer differences in the methylation state of older than younger mice.

In summary, DNAm age is a high-quality biomarker of ageing which can be tailored to the particular requirements of the context, be it estimating the chronological age of individuals, measuring

age acceleration/deceleration of individuals relative to their chronological age or estimating age disparity of tissues [168,169]. There is predictably a cost/accuracy trade-off. DNAm age is to some degree capturing “biological age” and likely not being driven by changes due to specific ageing-related diseases as illustrated by lack of relationships with specific pathologies and phenotypes of ageing noted by Bell et al. [131]. The PhenoAge clock is a better predictor of morbidity and mortality than age acceleration but may not be a better predictor of underlying biological age if the risk it is capturing is driven more by disease than that captured by age acceleration. PhenoAge is strongly associated with comorbidity count, more so than either the Hannun or Horvath clocks, interestingly zero comorbidities are associated with a negative PhenoAge score suggesting some of the PhenoAge signal is derived from disease status. Disease may, of course, interact with biological age causing an acceleration thereof, making them difficult to disentangle. Neither Age acceleration nor PhenoAge perfectly captures biological age but both are highly informative on the subject

Global DNA methylation decreases with age, DNA methylation becomes more disordered and less plastic to environmental influences with age. A set of loci, not necessarily individual CpGs but rather some functional unit of methylation exhibit consistent predictable changes with age. Which loci are in this set varies with tissue type; there is, however, a subset which seems to be common to most tissues. The preponderance of these loci are in areas of high to moderate CpG density and become hypermethylated with age. The high frequency of CpGs and the hypomethylated “ground state” of CGIs would lead one to expect a bias towards being located in CGIs and becoming hypermethylated. Thus it is unclear if changes in DNA methylation over time are biased towards regions of high CpG density and hypermethylation at a level that is greater or less than one would predict with the prior information.

## 1.6 Aims

Going forward this work covers several aspects of change in DNA methylation related to ageing phenotypes. Starting at the beginning of the human life-course with EWAS for *in utero* markers and interventions to impact on long-term bone health through the lens of the developmental origins of health and disease hypothesis. Continuing with an examination of the age-related changes in the methylation state of a key functional region of the genome, whose methylation state has not previously been characterised in detail, namely the tRNA genes. Lastly, the generation of new DNAm clocks based on a family of repetitive elements, specifically the Alu repeat family. Making use of another region of the genome the DNAm state of which has not previously been characterised in detail, but which represents a substantial proportion of all CpG sites and the derepression of which has long been considered to play a role in ageing.

1. Identify Age-related changes in DNA methylation in regions of the genome characterised

in the TwinsUK MeDIP-seq dataset and poorly covered by or covered in small samples by other technologies.

- a. The tRNA genes, following up on the previous finding of age-related DNA hypermethylation at tRNA-iMet-CAT-1-4 [170]. tRNA genes have a core role in cellular metabolism, many emerging regulatory functions both structural and as signalling molecules. tRNA genes also interface with many systems the modulation of which impact longevity, making any age-related changes in their epigenetic state potentially very consequential.
  - b. Alu repeat elements, a primate-specific family of SINEs present in over 1 million copies in the human genome. We aim to construct a DNA methylation clock using only these elements. We reason that age acceleration based on this clock will yield different information about biological ageing than from previous clocks as the effects of DNA methylation changes differ between repressed repetitive elements and regulatory sequences.
2. Identify epigenetic associations with bone health outcomes and of vitamin D supplementation during pregnancy as an intervention to improve bone outcomes [171,172]. Provide some groundwork for mechanistic studies to extend understanding of what systems influence bone health by identifying what genomic features change epigenetically with bone traits and interventions intended to affect these traits.

# Chapter 2

## Methods

### 2.1 Illumina DNA Methylation arrays

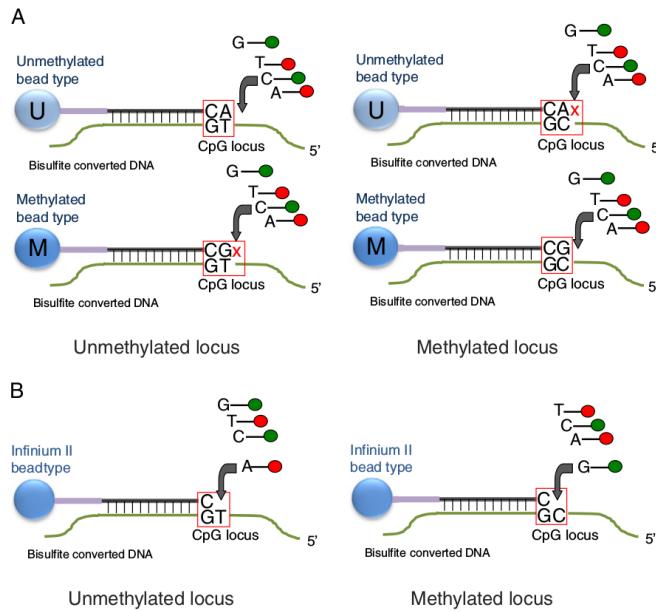
The Illumina methylation bead chip arrays make use of 50 bp long probes with sequences corresponding to the loci of interest in the genome, with an additional 23 base sequence for positional addressing on the array [121,126,133,173]. The CpG of interest is located at the 3' end of the probes. These probes are attached by their 5' ends to  $3\mu m$  diameter silica beads which are embedded in indentations acid etched into the surface of planar silica slides. Each locus is represented by an average of  $\sim 30$  beads and a minimum of  $\sim 5$  beads which are distributed randomly across the slide. The number of beads for a given locus is approximated well by the poisson distribution, so this can be used in conjunction with the number of loci to ensure that the probability of having fewer than 5 beads per locus is kept negligibly small by adjusting the ratio of loci to possible bead sites. Thus each BeadChip array has a unique random distribution of probes which must be decoded to yield the mapping between positions on the array a target loci being assayed. The positions of the beads are decoded through the sequential hybridisation of fluorescent oligonucleotide probes within the 23 base addressing sequence, with probes for each locus having a unique combination of hybridisation events to identify them [174,175]. An optical scanner such as Illumina's iScan system with a charge coupled device (CCD) to detect the light and a laser to excite the fluorescent dyes is used the read the arrays. There are two fluorophore channels the respective fluorescence intensities of which are recorded at each spot. Initially the probe decoding oligos are used to map the probe locations and then fluorescently tagged DNA bases are cycled to assay the DNA methylation as described below.

Samples are prepared for the array assay by fragmentation, bisulfite conversion and amplification. bisulfite conversion creates a 'pseudo-SNP' by selectively deaminating unmethylated cytosines to uracil and thus changing their base pairing properties and permitting this technology, originally developed to detect SNPs, to assay DNA methylation [121]. These sample fragments anneal

to complementary probes on the array and unbound DNA is washed away. A single DNA polymerase-mediated elongation step is carried out to permit a fluorescently labelled base to be incorporated at the end of each probe. C and G are tagged with one fluorophore and A and T with another. There are two probe designs:

1. The type I probes, which make use of two distinct bead types for each locus.
2. the type II probes, which make use of a single type of bead at each locus.

The ‘colour’ of the fluorophore incorporated signifies the methylation state and the intensity of the proportion of sites which are in that state. The logic of which ‘colour’ corresponds to methylated or unmethylated varies with the probe type and Figure 2.1 illustrates and explains this in-depth.



**Figure 2.1: Illumina Methylation Bead Chip Array Probe Designs** The probes on the array are 50bp in length. In type I probes, unlike type II probes the base that is incorporated (or not) in the elongation step is the base that follows the CpG site. The probes on the array are bound by the amplification products derived from the bisulfite converted DNA, not direct conversion products. Thus, Cs in the sample DNA represent methylated positions and Ts unmethylated positions, as Ts take the place of Us created by bisulfite conversion during the amplification process. **Type I:** Panel A shows the type I probe design. Each locus is represented on two beads, a methylated and an unmethylated bead. *Methylated Sample:* When sample DNA binds a methylated probe the G at the second position in the CpG on the methylated probe will be complementary to the C marking the location of the 5mC. Extension will then be able to occur incorporating a fluorescently labelled base complementary to the base just 5' of the 5mC position in the probe. On the unmethylated probe, the C marking the position of the 5mC will not be complementary to the A base in the probe and thus extension cannot occur. No fluorescent base will be incorporated at the unmethylated probe. *Unmethylated Sample:* When sample DNA binds the methylated probe the T marking the position of the unmethylated C is not complementary to the G at the second position in the CpG site on the probe. Thus, no fluorescent base is incorporated. Whereas, on the unmethylated probe the T marking the position of the unmethylated C will be complementary to the A at the second position in the CpG site on the probe. Thus, a fluorescently labelled base can be incorporated. **Type II:** Panel B shows the Type II probe design. For the Type II probes only one probe represents each locus. In the Type II probes the incorporated base is at the second position in the CpG site. A fluorescently labelled A is always incorporated opposite a T in the sample DNA marking the position of an unmethylated C, and a fluorescently labelled G is always incorporated opposite a C marking the position of an unconverted 5mC in the sample DNA. For the Type II probes, in contrast to the Type I probes, methylation is always signalled in the green channel and unmethylation in the red channel. (figure reproduced from Bibikova et al. [133].)

These arrays have now undergone several iterations [126,133,173]. The current array probe design was preceded by the ‘golden gate’ array which used a variant on the SNP probe design with methylation-specific PCR. This array contained ~1.5k probes covering 371 genes with a strong focus of cancer-related genes [121]. The first array using the probe design outlined above was the ‘Infinium’ array which contained 27k type I probes focused on the promoter regions for 14,475 consensus coding genes with 110 miRNA promoters [126]. The 450k array which had ~480k mixed type I and II probes covering 21,231 RefSeq genes, 26,658 CpG islands, 80,538 predicted enhancer regions and a variety of other features, including the MHC region [133]. The 450k array systematically underestimates the methylation level in highly methylated regions [176], this may well be because it samples only single sites, and highly methylated regions are often not fully methylated. A Markov chain model of DNA methylation state developed by Affinito et al. [134] suggests that once a threshold level of DNAm is reached on a molecule further methylation becomes less likely. They reason that DNMTs recruited by the high methylation state have greater difficulty accessing the remaining unmethylated sites. Thus sampling a single site is likely to produce a systematic underestimate as once neighboring sites are methylated methylation at a given CpG is less likely to increase further. The ‘EPIC’ array has ~850k probes including >90% of those on the 450k with greatly expanded coverage of loci with more dynamic methylation states than promoters and CGIs [173]. Promoters and CGIs tend to have relatively low methylation variability and the EPIC array aims to capture more functional methylation variation at regions such as enhancers [63]. The ‘EPIC’ array also contains both Type I and Type II probes, with many of the new sites being type II. The type I probes have greater dynamic range than the type II but take up twice as much space on the array so there is some trade-off between maximising the number of sites covered and the quality of data at those sites.

Methylation at each site on the array is commonly reported as a  $\beta$  value, which corresponds to the proportion of the sample DNA which was methylated at that site. This is computed from the intensity values extracted from Illumina’s IDAT format files which in the case of the methylation arrays are a binary format which had to be reverse engineered to permit analysis outside of Illumina’s ‘GenomeStudio’ [177].

$$\beta = \frac{\text{intensity}[M]}{\text{intensity}[U] + \text{intensity}[M]}$$

Where  $M$  = methylated and  $U$  = unmethylated.

The type I and type II probes have slight systematic differences not accounted for by their genomic context [178]. Type I probes have a wider range and are more reproducible than the type II probes. This is likely due to the dual complementary probes with fluorescence in different colour channels providing additional information for methylation level calling leading to more robust estimates at extreme values near 1 and 0. This difference means that to operate uniformly

on this mixed data a normalisation procedure is needed to correct for these differences.

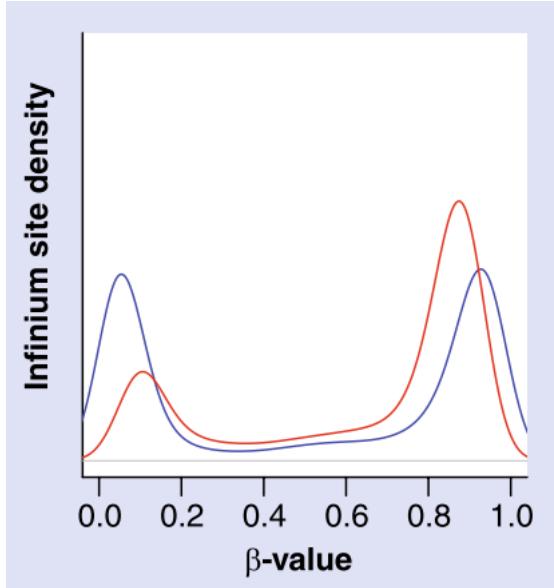


Figure 2.2: **Type I probes have a wider range and less variability than type II.** Type I probes shown in blue, Type II in red. Plot shows kernel density of beta values with gaussian smoothing. (Reproduced from Dedeurwaerder et al. [178])

There exist a number of normalisation methods to account for probe-type bias [179] such as beta-mixture quantile normalization (BMMQ) [180]. This method was applied as a part of the pre-processing of the Twins UK 450k array datasets used in (Chapter 4) in this thesis. In addition, there are other techniques capable of correcting for probe bias as well as other sources of technical variation such as batch effects [181]. Functional normalisation [182] which makes use of the technical control probe data on the arrays to inform its transformation of the DNA methylation values was used to normalise EPIC array data from the MAVIDOS and Southampton Women's Survey (SWS) data in (Chapter 3) of this work.

Bisulfite conversion efficiency can be estimated by looking at the rates of conversion of non-CpG cytosines as these are very rarely methylated they should be bisulfite converted at a very high rate. If this is not the case it can lead to the overestimation of DNA methylation levels, bisulfite conversion also fragments DNA so there is a trade-off between achieving sufficient conversion efficiency and not producing overly fragmented DNA [183]. On the Illumina methylation arrays control probes ([184] p226) are located at non-CpG cytosines to assay the conversion efficiency for use in quality control assessment of array data.

In addition to probe-type bias, there are other sources of potential problems for individual probes on these arrays [185]. Zhou et al. divide these into three categories [186]:

1. Probes with SNPs internal to the probes especially if located near the 3' end. These can

cause issues by:

- a) preventing extension through mismatches near the 3' end.
  - b) Altering the ability of the target cytosine to be methylated through effects on the sequence context, including the case where a C/T polymorphism in place of the C/T bisulfite induced conversion.
  - c) Inducing a colour change in a type I probe by altering the extension base.
2. Probes which map to multiple possible locations in the bisulfite-converted genome. Methylation values from multi-mapping probes are an amalgam of the methylation levels at the different loci to which the probes can map and cannot be disentangled to yield reliable results for the different loci.
  3. Probes with partial off-target matches to other loci. Probes with partial matches at the 3' end can lead to extension and an erroneous readout of methylation levels for that locus or partial matches elsewhere can lead to competition for binding with the probe intended for that locus.

To account for these potential problems probe ‘masks’ have been produced which identify probes which commonly exhibit one or more of these problematic properties so that findings at these loci can be excluded or followed up for validation with other methods [185,186]. These masks have been used here to identify potentially problematic probes.

When looking for a purely epigenetic effect the presence of a genetic variant which influences the DNA methylation state of a CpG is known as genetic confounding. This is not necessarily a problem if genetic effects on a phenotype of interest are a useful outcome but it is an issue if looking for effects mediated purely through epigenetic mechanisms. Genetic variants which bias the DNA methylation in a particular direction, or modulate the plasticity of DNA methylation state to environmental influences may be important to mechanistic understanding the biological system under examination [187,188]. Simply excluding all probes with a SNP above a certain minor allele frequency and with a certain distance of a target CpG can both miss rarer genetic confounders and throw away actual non-genetic differences in DNA methylation. Andrews et al. developed a method called ‘gap-hunting’ to identify probes which commonly have multi-modal DNA methylation distributions which may be associated with different genotypes influencing a probe locus [189]. ‘gap-hunting’ was also employed in this work to search for DNA methylation changes which may be under genetic influence (Chapter 3). In addition to SNPs causing sequence context changes that affect the DNAm of probes, structural mutations can also have this effect on DNA methylation [84]. Thus probes in the vicinity of non-SNP mutations should also be subject to additional scrutiny.

## 2.2 MeDIP-seq

MeDIP-seq uses a monoclonal anti-5mC antibody to bind denatured fragmented genomic DNA at methylated CpG sites. This antibody-bound fraction of DNA is isolated and sequenced [98], Figure 2.4 outlines the workflow. Unlike bisulfite conversion approaches this method permits 5mC to be differentiated from 5hmC, as the antibody binds specifically to 5mC. The methylation level across the CpG sites in the region of the genome to which the resultant sequencing reads map can subsequently be estimated by counting the number of reads and accounting for the CpG density with software tools such as the MEDIPS R package [99].

MeDIP-seq can also cover much more of the methylome than do the 450k or EPIC arrays. The 450k array captures ~1.6% of the 28 million CpGs in the genome. At near saturation coverage of ~40 million reads per sample, MeDIP-seq can cover ~60% of all CpGs with at least 1 read, this is almost all methylated CpGs. At half that total read count MeDIP-seq still covers ~50% of CpGs at  $\geq 1\times$  and ~20% at  $\geq 5\times$  [190]. Other estimates place the threshold for saturation coverage by MeDIP-seq lower at between 20 and 30 million reads and suggest that the maximum number of CpGs covered at least once could extend up to 90% [191]. MeDIP-seq provide a particular advantage by comparison with the Illumina beadchip array technologies with respect to the coverage it provides for repetitive regions of the genome. This is illustrated by Clark et al. who show that MeDip-seq covers an estimated 91.4% of repetitive sequences and the 450k array only 3.4% [176] (Figure 2.3). This makes it particularly suited to the aims of this project (Chapters 4 & 5).

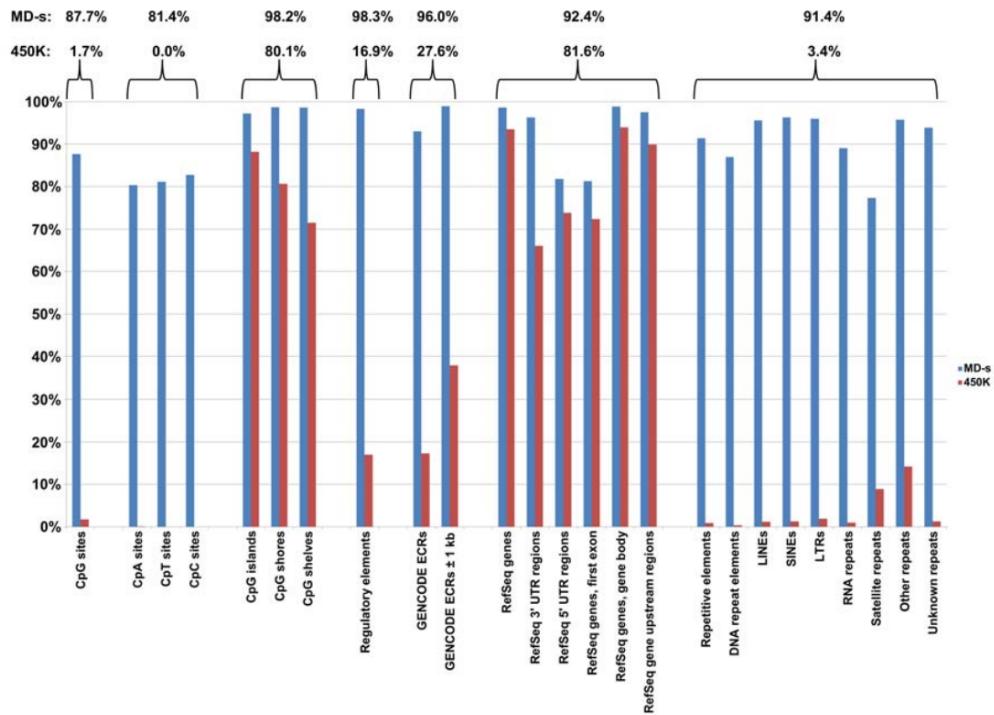
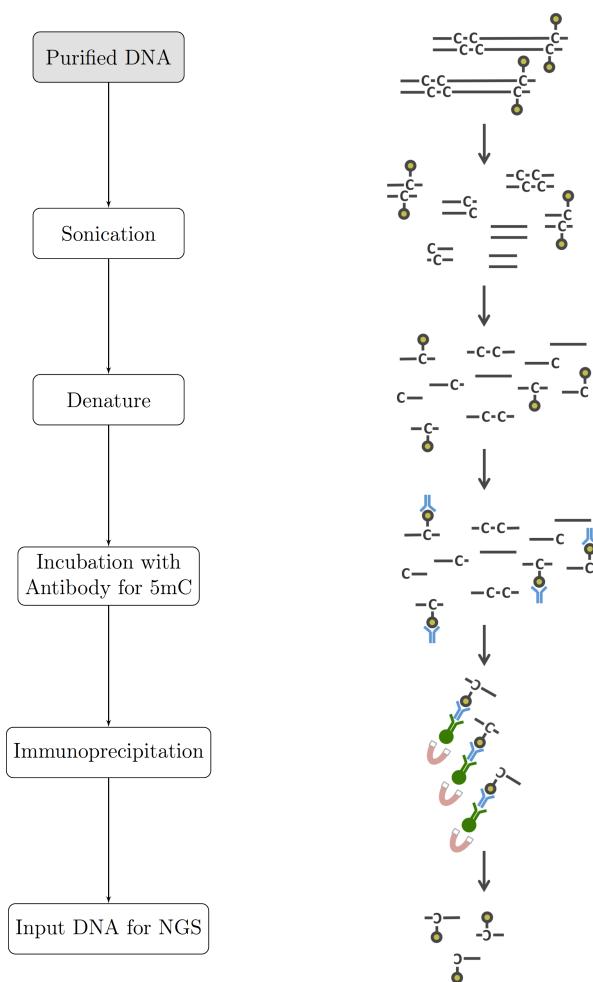


Figure 2.3: MeDIP-seq provides superior coverage of repetitive regions of the genome than the Illumina array platform. This presumes saturation coverage for the MeDIP-seq assay. MD-s = MeDIP-seq, 450k = Illumina 450k DNA methylation bead chip array (Reproduced from Clark et al. [176]) The CpG islands, shores, and shelves major RefSeq gene features are covered almost as well in the 450k array as they are by MeDIP-seq but repetitive regions and other regulatory elements are where MeDIP-seq provides substantially greater coverage.



**Figure 2.4: Graphical Summary of the MeDIP-seq process** Purified DNA is fragmented by sonication, denatured and incubated with anti-5mC antibodies. It is then immunoprecipitated resulting in fragments containing methylated CpGs which are subsequently sequenced.

DNA sample preparation, MeDIP reaction, and Illumina sequencing were performed by BGI-Shenzhen (Shenzhen, China). Fragmentation of genomic DNA from whole peripheral blood samples from the Twins UK cohort was performed by sonication using a Covaris system (Woburn, MA, USA).  $5\mu\text{g}$  of DNA was used in the Illumina Single-End DNA Sample Prep kit followed by end repair, base addition and adaptor ligation. The MeDIP reaction was performed, according to the protocol for the MagMeDIP kit (mc-magme-048), with a monoclonal antibody for 5-methylcytosine (5mC) (Cat. No.: CO2010021 mc-magme-048 from Diagenode (Liège, Belgium)). The antibody was incubated with Adaptor-ligated DNA combining  $0.5\mu\text{l}$  antibody +  $0.5\mu\text{l}$  water; then  $0.6\mu\text{l}$  MagBuffer A,  $1.4\mu\text{l}$  water and,  $2\mu\text{l}$  MagBuffer C; yielding a final volume of  $5\mu\text{l}$  for the immunoprecipitation reaction. The MeDIP reaction was validated with quantitative PCR. The product of the MeDIP reaction was purified with Zymo DNA Clean & Concentrator-5 (Zymo Research), and amplified with adaptor-mediated PCR. Size selection of fragments (200–500 bp)

was performed by gel excision and quality assessed by Agilent BioAnalyzer (Agilent Technologies, Santa Clara, CA, USA). The resultant libraries were subjected to highly parallel 50bp single-end sequencing on the Illumina HiSeq2000 platform.

Quality control (QC) and alignment were performed by Twins UK. Sequencing data were subject to initial QC for call quality and base composition using FASTQC (v0.10.0) (<https://github.com/s-andrews/FastQC>) [192]. Duplicates were removed with SAMTools (<https://github.com/samtools/samtools>) [193]. The alignment was performed using the Burrows-Wheeler Aligner (**bwa**) (<https://github.com/lh3/bwa>) [Li2009] with a minimum mapping quality score of Q10. MEDIPS (v1.0) was used to perform MeDIP-specific QC, as well as to generate reads per million (RPM) and Absolute methylation score (AMS) values (<https://bioconductor.org/packages/release/bioc/html/MEDIPS.html>) [99]. AMS and RPM values are binned into 500bp windows with a 250bp slide in the BED format, resulting in ~12.8 million windows on the genome (Build: hg19/GRCh37). Additional quality control checks were performed in the R language: correlation matrix, hierarchical clustering, dendrogram, heatmap, density plot, and batch effects inspection by principal component analysis. There was a mean of ~16.89 million (SD 3.04M) reads per sample. This number of reads per sample is consistent with covering ~50% of all CpGs at 1x, ~15% at 5x and ~10% at 10x (estimated using the results of Taiwo *et al.* 2012 [190].) MeDIP-seq data from regions of interest was extracted using **bedtools** v2.17.0 [194].

### 2.2.1 Participants

Peripheral blood samples for DNA extraction were provided by the adult volunteers from the UK Twin Register (TwinsUK Resource [195]). The participants are rigorously phenotyped at visits at St Thomas' Hospital, London. A twinning questionnaire determines twin zygosity and is confirmed by genotyping. Ethics were approved by Guy's & St Thomas' NHS Foundation Trust Ethics Committee (EC04/015—15-Mar-04) and written informed consent was obtained from all subjects in accordance with this. Samples are stored at -80 °C in EDTA tubes before extraction via the Nucleon Genomic DNA Extraction Kit. DNA is subsequently stored in TE Buffer at -20 °C.

The MeDIP-seq dataset used in this work consists of 4350 whole blood methylomes with age data. 4054 are female and 270 male. 3001 have full blood counts. There are 3652 individuals in this data set. These individuals originate from 1933 unique families. There are 1234 monozygotic (MZ) twin pairs (2468 individuals), and 458 dizygotic (DZ) twin pairs (916 individuals).

The above-described processing steps for Methylated DNA Immunoprecipitation sequencing (MeDIP-seq) data used here have been previously described in Bell *et al.* [84,170]. These processed data are available, subject to approval, from the European Genome-phenome Archive (EGA) (<https://www.ebi.ac.uk/ega>) under study number EGAS00001001910 and dataset (<https://www.ebi.ac.uk/ega/datasets/EGAD00010000983>).

## 2.3 Targeted Bisulfite sequencing

Targeted Bisulfite sequencing [196] is very similar in principle to ordinary targeted DNA sequencing but preceded by a bisulfite conversion step as illustrated in Figure 1.17 [102]. However, bisulfite conversion has several implications for downstream processing. Bisulfite library prep can be either directional, sequencing reads correspond to a bisulfite-converted version of either the forward or reverse strand, or non-directional in which sequencing reads correspond to bisulfite-converted versions of either strand giving a total of four bisulfite-converted sequences with the strand of origin unknown. If a library is directional there are only two bisulfite-converted sequences. Library preparation in this work was non-directional.

Mapping can be performed with customised software wrapper to implementations of existing alignment tools. **Bismark** is a popular example of such as wrapper which uses the bowtie alignment tool and calls site methylation levels [197]. Alignment is performed with *in silico* bisulfite converted versions of the genome and reads, with G to A conversions for reverse strand reads. This requires four parallel instances of bowtie for forward and reverse G to A and C to T conversions, which are combined to determine the unique best alignment. (Figure 2.5 panel A). Bisulfite conversion results in loss of sequence complexity as many Cs effectively become Ts meaning short reads can be challenging to map. For targeted sequencing approaches it is advisable to align to both the targeted sites and their flanking regions as well as to the whole genome, especially if any of the targeted regions contain repetitive sequences [198]. Substantial mapping to areas other than those targeted could be indicative of off-target amplification, this step forms a part of the quality control process. **Bismark** produces methylation calls either combined or by strand and which can be filtered by methylation context CpG, CHG or CHH.

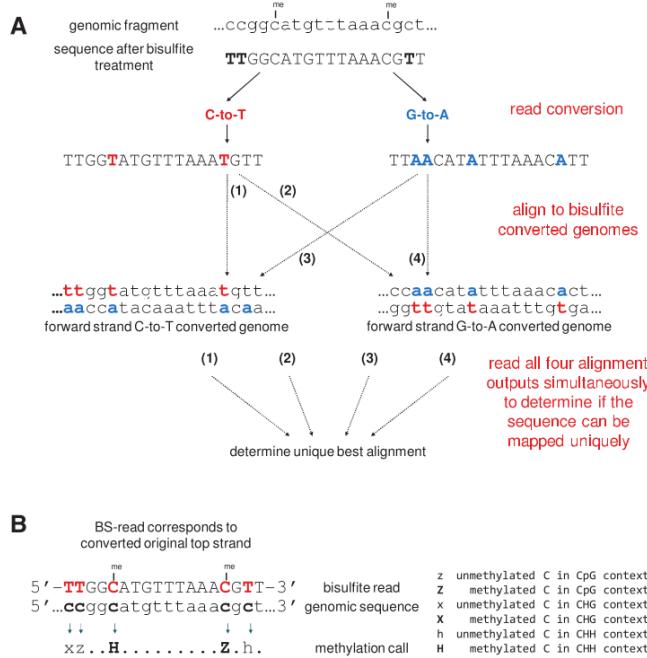


Figure 2.5: Bisulfite sequencing alignment process (Reproduced from Krueger et al. [199])

Targeted sequencing methods where the aim is to sequence a number of amplicons ( $N_a$ ) from a number of samples ( $N_s$ ) means the total number of PCR reactions in the pairwise matrix is ( $N_a \times N_s$ ). This can rapidly become difficult to manage with conventional bench-top PCR methods especially where consistency between reactions is important for quantitation. Microfluidic systems for conducting multiple PCR reactions in parallel such as the Fluidigm 48.48 access array employed here improve the ease with which a large number of reactions can be performed.

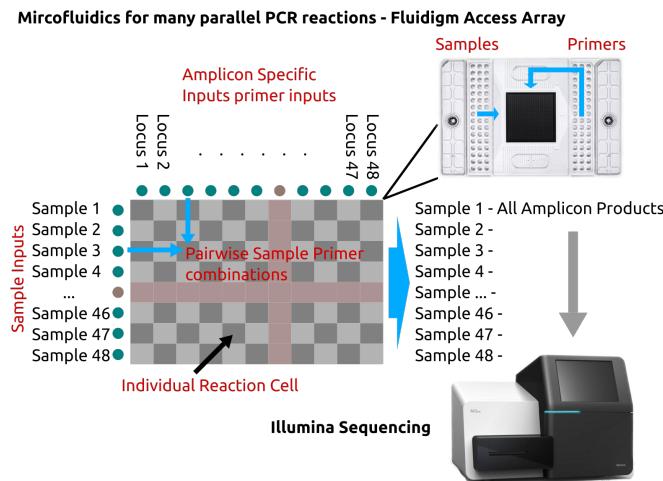


Figure 2.6: Microfluidic system for multiple parallel pairwise PCR reactions. Fluidigm 48.48 access array.

A four primer step-out PCR process is performed on the bisulfite-converted genome to generate the final PCR product used in the Illumina sequencing reaction (Figure 2.7). After the loci of interest are selected locus-specific primers are designed to be complementary to the bisulfite converted regions flanking the area of interest and avoiding CpG sites in those flanking regions. A number of tools for designing such primers are available, the tool used in this work was: ‘methPrimer’ [200].

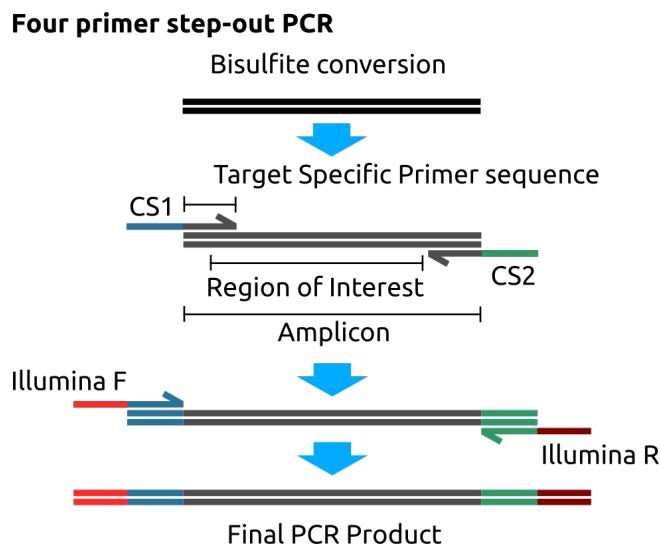


Figure 2.7: **Four primer step-out PCR for targeted bisulfite sequencing on the Fluidigm access array.**

After designing the target-specific primer sequences the CS1 and CS2 adaptor sequences are added to the forward and reverse primers respectively to permit the addition of Illumina sequencing primers in a second PCR step. The subsequent analysis is similar to that for whole-genome bisulfite sequencing. The PCR products are sequenced and the resulting reads aligned and methylation levels called with **Bismark**.

The targeted BiS-sequencing data is available at: <https://www.ncbi.nlm.nih.gov/bioproject/PRJNA635108>



## Chapter 3

# Epigenome Wide Association Studies for Bone Health Outcomes in Umbilical Cord Blood and Tissue

### 3.1 Abstract

Long term bone health and fracture risk is strongly influenced by the peak bone mass that is attained in early adulthood, this has its origins in development from late pregnancy through early childhood. The availability of vitamin D to facilitate the uptake of calcium and phosphate is essential to healthy bone development, in addition to the role of vitamin D as a hormone mediating this process. Greater mechanistic understanding of this process is desirable to permit the design of more effective interventions to reduce the substantial health burden of osteoporotic fractures in the elderly. Maternal vitamin D levels during pregnancy have previously been shown to impact on the methylation status of specific genes associated with vitamin D signalling and metabolism, such as *RXRA*. Long-term bone health outcomes have also been associated with both maternal vitamin D during pregnancy and methylation of the *CDKN2A* gene. To attempt to identify other loci, the epigenetic regulation of which is pertinent to the impact of maternal vitamin D levels on long-term bone health, Epigenome-Wide Association Studies (EWAS) were performed with DNA methylation data from the Illumina EPIC & 450k methylation bead chip arrays. We tentatively identify two CpGs whose DNA methylation state is associated with bone mineral content at 6 years of age and periosteal circumference at 6 years of age respectively.

## 3.2 Introduction

In England and Wales a survey of general practice data indicates that 53% of women and 21% of men will experience a bone fracture in their remaining lifetime from the age of 50 years [201] (Figure 3.1). Hip fractures in particular are associated with a reduction in survival of an estimated 10-20% with most deaths occurring in the first 6 months, vertebral fractures though approximately an order of magnitude less common carry a similar mortality risk [201,202]. The national osteoporosis foundation estimated the societal cost of osteoporosis at 58 Billion USD for 2018 in the USA alone [203] with serious negative implications for the quality of life of those with this condition. Thus interventions to increase bone health and reduce the rate of osteoporosis as well as developing a mechanistic understanding of the development of this condition are of considerable medical interest.

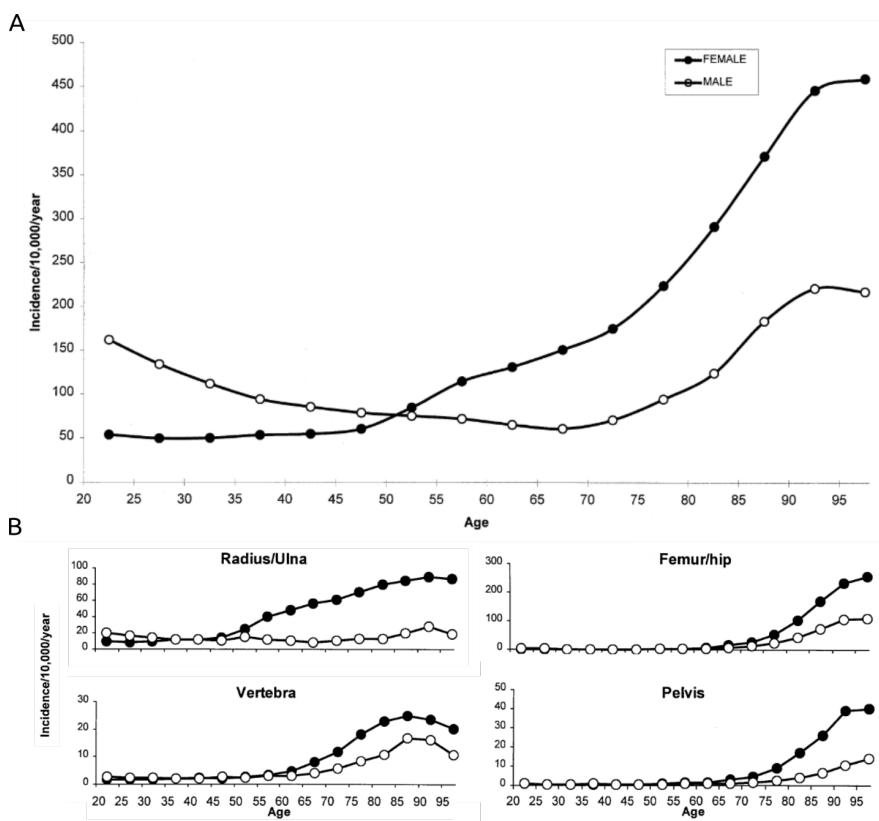
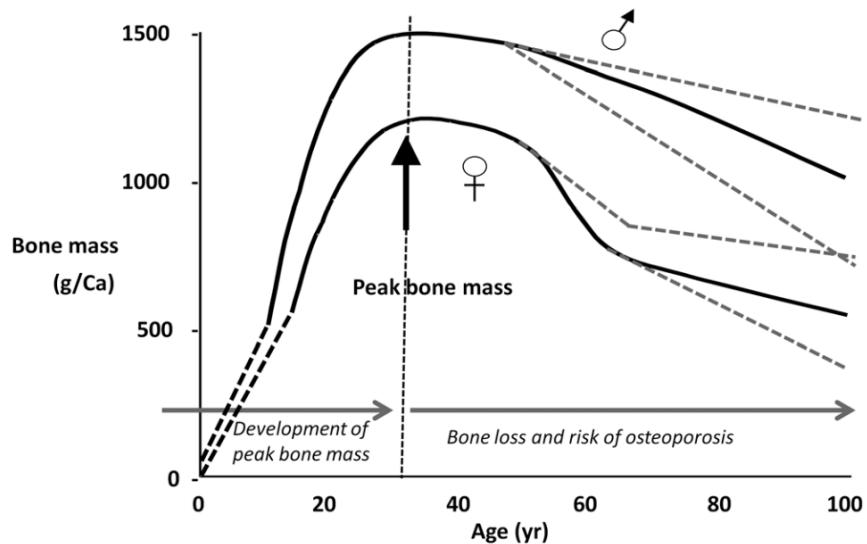


Figure 3.1: **Overall Fracture Incidence Increases With Age** Data from 5 million adults in the General Practice Research Database 1988-1998. **A** Incidence of all fractures at any site by sex and age. **B** Fractures by sex and age at select sites, pelvis, Femur/Hip & Vertebra confer the greatest increased mortality rates and Radius/Ulna has the greatest sexual dimorphism. (Adapted from van Staa et al. [201].)

Modelling work indicates that the largest determinant of osteoporosis risk is peak bone mineral density (BMD) with a 10% change in peak BMD leading to a predicted 13 year delay in the

development of osteoporosis [204]. Though the baseline BMD is a better predictor of fracture risk, the rate of loss also independently predicts risk [205–207]. Thus factors altering the peak bone mass will be important for modulating osteoporosis risk. Bone mass increases from intrauterine development through to a peak in early adulthood (Figure 3.2).



**Figure 3.2: Bone Mass Over The Lifecourse** Bone mass increases from intrauterine development through to a peak in early adulthood. Intervention in early life to modulate growth trajectory and increase peak mass may provide a higher starting point and delay bone loss in later life. (Reproduced from [208].)

Calcium, phosphate and vitamin D are needed for bone development [209] and maintenance. Vitamin D is needed for sufficient intestinal absorption of calcium [210] and phosphate [211] in addition to the hormonal role it plays in the regulation of bone development. Calcium and vitamin D intake are modifiable factors and thus constitute a reasonable starting point for interventions to improve bone health. This raises the question of how and when they can best be deployed to maximise their impact on bone health.

Birth weight predicts bone mass in adulthood [212]. A longitudinal study with long term follow up at Helsinki University Central Hospital linked low childhood growth rate with greater fracture risk in later life [213]. Data from the Southampton Women's Survey (SWS)[214] suggest that intervention in the interuterine period and infancy may afford the greatest opportunity to impact on long term bone health. Rate of growth in late pregnancy (19-34 weeks) and early post natal growth (up to ~2 years) predict bone mass at 4 years of age [215], and estimates of hip strength at 6 years [216]. The proportion of individuals crossing between tertiles of the distribution of growth in length decreases with age, indicating that the variability and thus potential malleability of growth is greatest during the period from 11 weeks gestation to ~2 years.

Observational data indicated a relationship between maternal vitamin D status and bone outcomes in children [217,218]. Prospective studies showed the effects of maternal vitamin D status on bone growth in newborns [219] and at 20 years of age [220]. MAVIDOS (Maternal Vitamin D Osteoporosis Study) is a randomised, double-blind, placebo-controlled trial of oral vitamin D supplementation (1000 IU cholecalciferol/day or placebo) from 14 weeks gestation in women with initial circulating 25(OH)-vitamin D levels of 25-100 nmol/l [171,172]. The primary outcome of the study was whole body bone mineral content, assessed by Dual Energy X-ray absorptiometry (DXA) within 14 days of birth, and follow up DXA at 4 years of age. MAVIDOS did not find significant differences in neonatal bone outcomes overall but there were significant differences between intervention and placebo for births taking place during the winter months.

$1,25(OH)_2$ -Vitamin D has been implicated in the control of Plasma membrane  $Ca^{2+}$ -ATPases (PMCA<sub>s</sub>). The expression of the *PMCA3* gene in the placenta has been associated with umbilical cord calcium concentration and intrauterine accrual of bone mineral content [221]. This provides part of a candidate mechanism linking intrauterine vitamin D availability to bone outcomes, but much remains to be elucidated about how maternal vitamin D status may impact bone development. Work in rat models has demonstrated the ability of maternal nutritional exposure during pregnancy to influence the epigenetic state and expression of genes in their offspring, specifically a change in DNA methylation was identified [222–224]. Consequently, the DNA methylation status of target genes of interest for vitamin D and bone development related pathways have been examined [225].

Retinoid-X-receptor-alpha (RXRA), which forms a heterodimer with the vitamin D receptor is essential for the nuclear action of  $1,25(OH)_2$ -Vitamin D. DNA methylation at four of ten CpG sites in the *RXRA* gene were significantly ( $p \leq 0.05$ ) lower in the umbilical cord of offspring from cholecalciferol supplemented mothers compared to placebo with a mean difference in methylation of -1.98% (n=447, 95% CI -3.65 to -0.32, p=0.01). One CpG site in the gene is related to estimated free 25(OH)D levels [226].

Methylation of cyclin-dependent kinase inhibitor 2A (CDKN2A) in umbilical cord tissue has been implicated in bone cell activity mediating skeletal development and homeostasis. The *CDKN2A* locus has quite complex biology encoding two cell cycle inhibitors *p14<sup>ARF</sup>* and *p16<sup>INK4a</sup>* as well as the long non-coding RNA *ANRIL* which inhibits *p16<sup>INK4a</sup>*. Targeted analysis of nine CpGs in a 300bp stretch of the *ANRIL* promoter region of *CDKN2A* was carried out. Methylation at several of these sites showed an inverse correlation with bone size, mineral content, and mineral density at age four years [227].

Epigenetic states are a product of both genetics and environment, they exist on a continuum from complete or obligatory genetic determination through facilitative variation influenced by genetics, to the purely epigenetic subject entirely to environmental factors. This picture is complicated by the observation that genetics can not merely influence epigenetic state but also

the extent to which that state will vary depending on environmental factors [228]. Epigenome-wide association studies have been performed on DNA methylation data for many traits including age (Section 1.5.1), smoking status [229], exposure to atmospheric particulate matter [230], and obesity [231]. Some robust and reproducible DNA methylation changes have been identified with EWAS, particularly for traits such as age and smoking. Unlike genome-wide association studies (GWAS) associations found by EWAS can be due to reverse causation. That is to say, a genetic variant associated with a disease is very unlikely to have been caused by the disease state whereas an epigenetic association can be the consequence of a disease state. Whilst this can complicate the dissection of disease aetiology it can also be a boon in unpicking the biology that follows on from a particular disease or environmental exposure by revealing what biological networks are affected.

The previous work identifying DNA methylation changes at the *RXRA* [226] and *CDKN2A/ANRIL* [227] loci associated with childhood bone outcomes lead to an interest in establishing if DNA methylation changes were occurring at other loci. EWAS using the Illumina 450k and EPIC array technologies were performed for maternal Vitamin D status and childhood bone outcomes in data from the MAVIDOS and SWS cohorts.

### 3.2.1 Outline of EWAS

Three sets of EWAS analyses were performed: Phase I of the MAVIDOS EPIC arrays and comparison with 450k array results. Phase II of the MAVIDOS EPIC analysis with additional samples. Analysis of the SWS cord blood EPIC array data.

- MAVIDOS phase I results (EPIC n = 140, 450k n = 60), EWAS outcomes:
  - Intervention Vs. Placebo
  - Bone Mineral Content at Birth (g), measured by DXA
  - Maternal circulating 25(OH)-vitamin D levels (nmol/l) at 34 weeks
  - Change in Maternal Vitamin D from 11 to 34 weeks
- MAVIDOS phase II results (n = 237), EWAS outcomes:
  - 4 year Total BMC (g), without heads, measured by DXA
- SWS cord blood, EWAS outcomes:
  - 8 year Total BMC (kg), without heads, adjusted for sex and age, measured by DXA (n = 408)
  - 6 year Total BMC (kg), without heads, adjusted for sex and age, measured by DXA (n = 402)
  - 6 year Periosteal circumference of the Tibia at 38% from the distal end (mm), measured by Peripheral Quantitative Computed Tomography (PQCT) (n = 141)
  - 6 year Cortical density of the Tibia at 38% from the distal end (mg/cm<sup>3</sup>), measured by Peripheral PQCT (n = 141)

### 3.3 Methods

Analysis of the Illumina EPIC and 450k methylation arrays was carried out in the R statistical programming language (v3.5.2) using the `meffil` package [232], which was chosen as it is capable of performing functional normalisation in a more memory efficient fashion than alternatives such as `minfi`.

#### 3.3.1 Functional Normalisation

Functional normalisation is an approach to removing unwanted variation associated with ‘batch’ effects such as the date on which a sample was analysed or which slide a sample is on which was developed by Fortin et al. [182]. This noise in the data masks the signal associated with the underlying biological effect of interest.

Functional normalisation makes use of control probes on the arrays which are designed to capture only technical variation as surrogates for the sources of unwanted variation. The control probes are processed into 42 summary measures, and principal components analysis is performed on the control probe summary matrices for all samples. The top  $m$  principal components (PCs) are used as the surrogates for technical variation going forward. The number,  $m$ , of PCs from the control probe summary matrices used is informed by the amount of residual variation remaining after normalisation. Picking the number of PCs which correspond to the last steep drop in residual variation is the approach recommended by the implementers of `meffil`, and Fortin et al. [182] recommend an  $m$  of 2 as performing well across a variety of analyses. (See supplementary material from [182] for details of the control probe summary process.)

The process used by functional normalisation is a variant on quantile normalisation. Instead of forcing the empirical marginal distributions of the samples to be identical at each site across arrays. To produce the normalised values it constructs a quantile function which only removes variation arising from surrogates for batch variation. This approach is effective even when comparing samples with large global differences in methylation levels such as between normal and cancer samples, but cannot overcome high degrees of confounding [182].

#### 3.3.2 Genetically Confounded and Multi-mapping Probes

Probes from different locations in the genome with similar sequences, especially following the reduction in sequence complexity associated with bisulfite conversion, can be cross-reactive on the arrays leading to erroneous signals and are thus commonly excluded from analysis. DNAm is often strongly influenced by genetic factors, the effect is especially pronounced when variants alter the sequence at CpG sites themselves as the site can then no longer be methylated. Thus, sites that have common genetic variants at the probe site are also excluded from the analysis, as well as some sites whose methylation is known to be under strong genetic influence by common

genetic variation.

43,254 probes on the EPIC array have been identified as multi-mapping, DNA binding to the probes may be derived from other locations in the genome invalidating these probes as a measure of methylation at their intended loci [185]. 12,510 probes were identified as having genetic variants at the CpG locus they are intending to assay, this can produce misleading results as mutant bases can resemble the products of bisulfite conversion [185]. 1,812 probes were found to overlap regions exhibiting haplotype-specific methylation associated with common non-SNP genetic variants (CNVs, Indels, STRs) and regional in-phase clustering of CpG-SNPs [84]. Zhou et al. [186] provided a list of probes which they recommend ‘masking’ from the 450k array due to multi-mapping issues, genetic variants overlapping the CpG sites, and other factors which may render results from these probes problematic.

In order to identify any potential additional sources of genetic confounding in the phase I MAVIDOS analysis, Probes with methylation values which cluster into distinct groups were identified using the ‘gap hunting’ method developed by Andrews et al. [189]. Such distinct clusters of methylation can arise from genetic variants which influence methylation levels being present in homozygous and heterozygous forms in the study population, see Figure 3.15. As the sensitivity and specificity of ‘gap hunting’ is limited, it is the advice of the authors not to exclude probes flagged by `gaphunter()` prior to performing EWAS. It is instead advised to check if any of the results appear in this list after the fact and examine the possibility of a genetic effect if they are. In phase II of the MAVIDOS analysis and in the SWS analyses only probes with specific technical QC issues were excluded prior to normalisation and EWAS. Flagging of problematic probes occurred after EWAS on any significant results.

### 3.3.3 EWAS Models

#### 3.3.3.1 Cell-Type Correction

Cell-type composition is a known confounder in epigenetic studies [164,165]. Observed variation in DNA methylation can be cell-type intrinsic, where changes in DNA methylation are not driven by changes in the cell-type composition of the tissue; or cell-type extrinsic, where changes in DNA methylation are due to changes in the cell-type composition of the population sampled (Figure 3.3). An established approach for addressing this potential source of confounding is to add terms to the regression model which reflect the cell-type composition of each sample. Cell-type composition can be ascertained through three main approaches: Direct cell count data; Estimating cell counts using the experimental data and models fitted on DNA methylation data from reference panels of known cell-type composition [165]; or reference free approaches which make use of mathematical methods to identify sources of confounding variation such as cell-type heterogeneity. Whilst there are several reference samples available for cord blood [233–236], none were available for cord tissue at the time of the initial analysis, and consequently a

reference free technique was used in phase I of the MAVIDOS analysis. Models were fitted using Surrogate Variable Analysis (SVA) [237] and independent Surrogate Variable Analysis (iSVA) [238]. The focus of these results is on SVA as minimal differences between the two methods were observed and SVA was recommended based on comparisons of the performance of various cell-type heterogeneity correction methods [239,240].

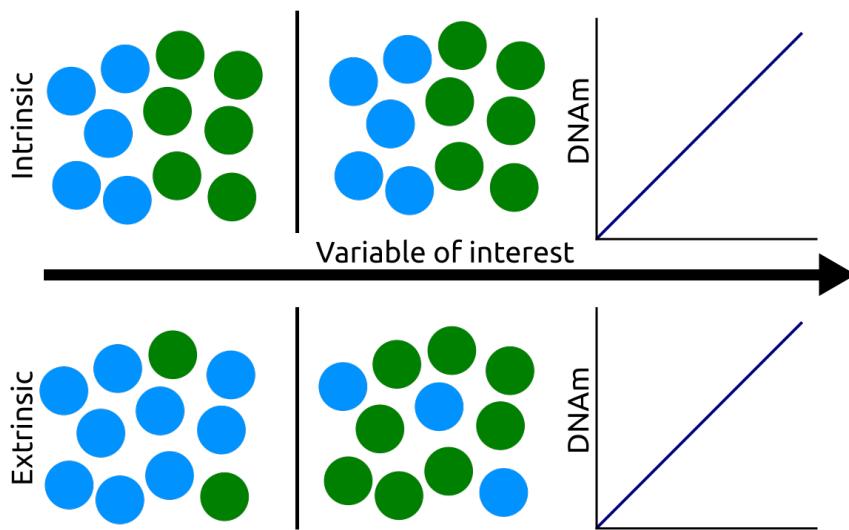


Figure 3.3: Diagrammatic representation of DNAm change arising from extrinsic or intrinsic changes in DNA methylation. Extrinsic changes are due to shifts in cell-type composition. Intrinsic changes in DNA methylation occur without changes in the proportions of cell-types. These two modes of change are of course not mutually exclusive and both can be occurring.

A reference panel for cord tissue samples was recently published [241]. This cell-type reference had not yet been integrated into the `meffil` R package used to perform the EWAS analyses so I created a fork of `meffil` including this reference panel in the required format, this can be installed directly from github. The code I used to add this reference to the package data before building the R package can be found in this [gist](#). Cell-type correction based on this reference was used in phase II of the MAVIDOS analysis.

### 3.3.3.2 Structure of models fitted for each EWAS

By default EWAS in `meffil` are run with four different models:

1. No covariates, attempting to predict methylation with the variable of interest.
2. All covariates, attempting to predict methylation with the variable of interest plus a user-supplied list of covariates.
3. Surrogate Variables + all covariates, attempting to predict methylation with user-supplied covariates and surrogate variables generated from SVA.

4. Independent Surrogate Variables + all covariates, attempting to predict methylation with user-supplied covariates and independent surrogate variables generated from iSVA.

Running EWAS with multiple models permits the effects of adding the various covariates on the results of the analysis to be seen.

### 3.3.4 Concordance of EPIC and 450k EWAS Results

In order to ascertain if the results from the EWAS in the 450k and EPIC arrays were producing similar sets of probes in the top ranking positions when ordered by p-value, the concordance (% overlap) between the top  $k$  probes, where  $k = 1..100,000$  was calculated. Only probes in common between the two arrays were used and  $k$  was incremented in steps of 50.

## 3.4 Results

### 3.4.1 MAVIDOS phase I

DNA methylation at none of the probes were significantly associated with any of the four variables of interest for each EWAS performed in either the EPIC or 450k datasets. The concordance between the probes with the top-ranked p-values in common between the EPIC and 450k data was at the level of chance.

#### 3.4.1.1 Whole array QC

**3.4.1.1.1 EPIC arrays** The predicted sex of the samples generated using sex chromosome probe intensities was checked against that in the sample annotation and two mismatches were found. These were MAVIDOS IDs 206 and 63 and the associated arrays were excluded from further analysis, (Figure 3.4).

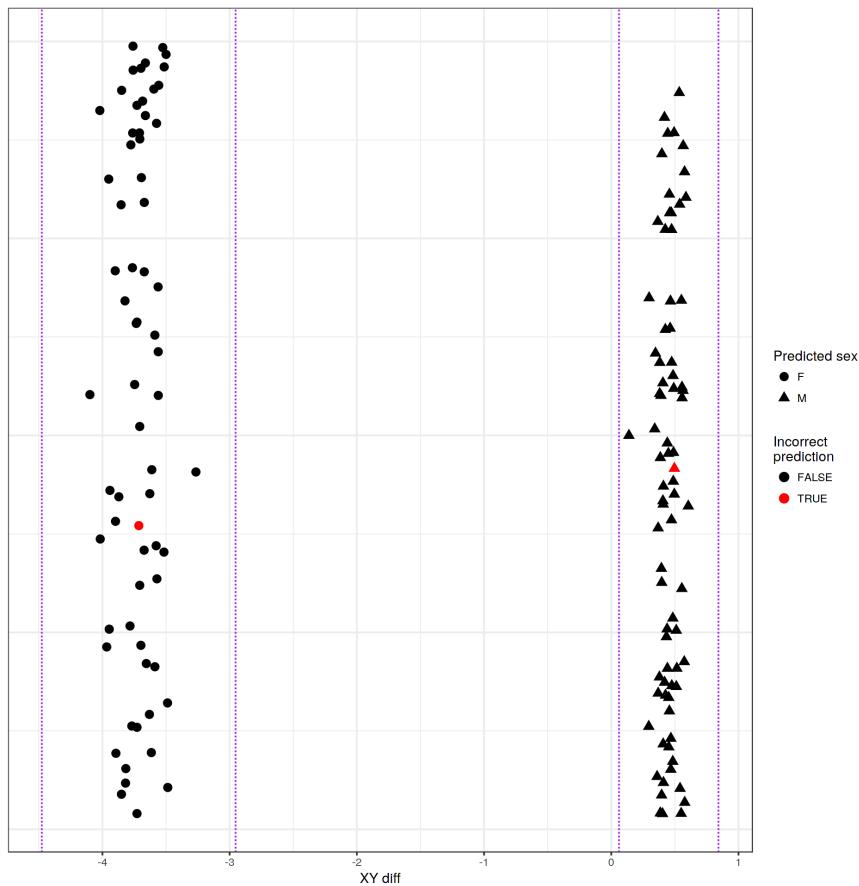


Figure 3.4: Predicted sex of each sample based on the sex chromosome copy numbers inferred from probe intensities for the EPIC array data. Mismatches between the predicted sex and that asserted in the sample annotation metadata are shown in red. Two predicted sex values differ from their annotations. Plot generated by `meffil` QC report.

The dataset also contained four samples for which there were two technical replicates, only the first replicate from each was used (144 arrays run for 140 individuals). Array 201516310023 (MAVIDOS ID 95) was excluded as its median methylated signal was more than  $3\sigma$  from the expected value, (Figure 3.5). No samples were excluded for having a higher than expected proportion of undetected probes (proportion of probes with detection p-value  $> 0.01$  is  $> 0.1$ ) (Figure 3.6). No samples were excluded for having a high proportion of probes with low bead counts (proportion of probes with bead number  $< 3$  is  $> 0.1$ ), (Figure 3.7). In total 3 of the EPIC arrays were excluded from the analysis for failing quality control leaving an n=137.

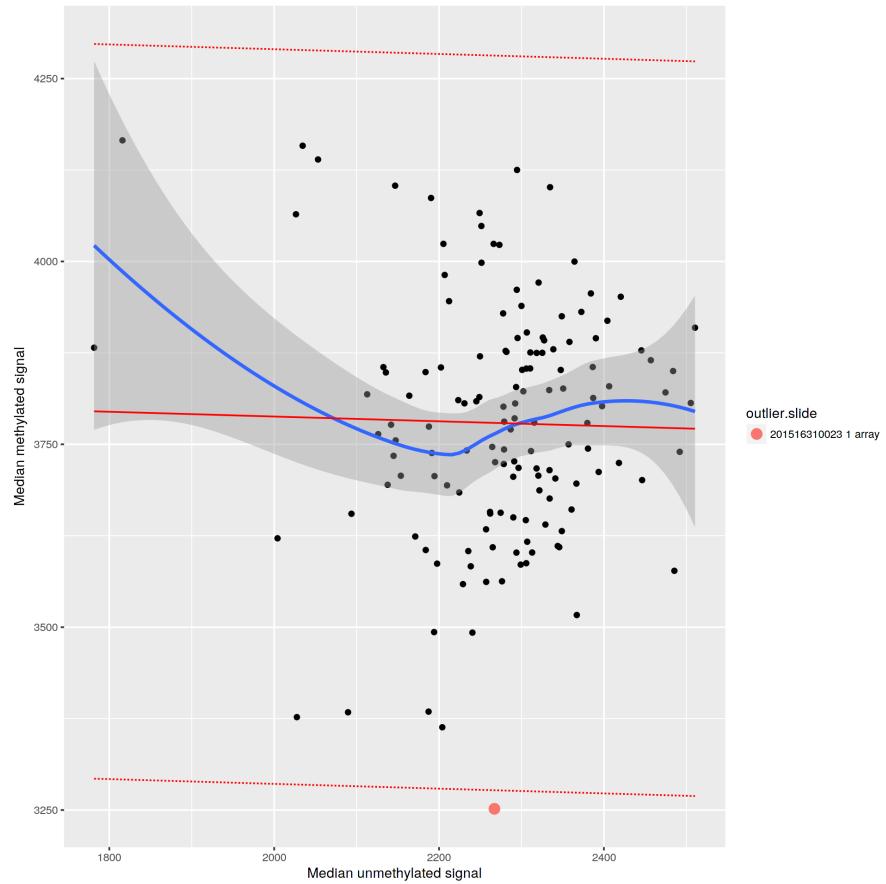


Figure 3.5: Median methylated signal vs unmethylated signal per sample for the EPIC array data, solid red line indicates linear regression of median methylated signal vs median unmethylated signal with dotted red lines representing  $3\sigma$  from the expected mean. Samples outside the expected range are indicated in the legend. Plot generated by `meffil` QC report.

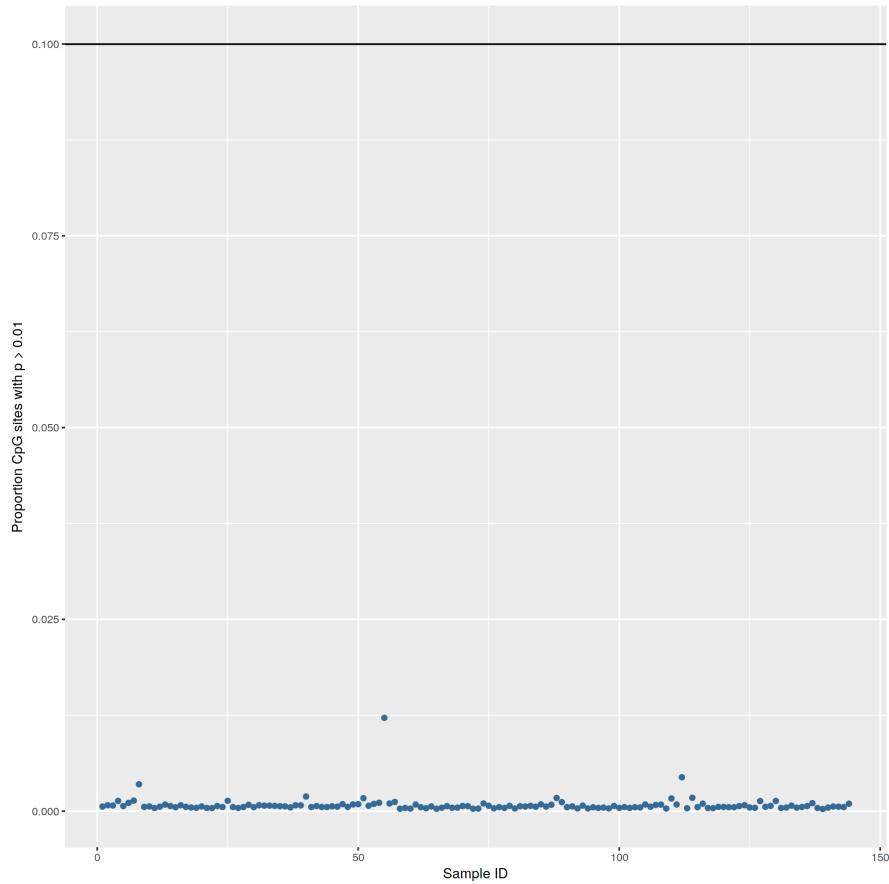


Figure 3.6: Proportion of probes with detection p-values  $>0.01$  by sample for the EPIC array data. Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil` QC report.

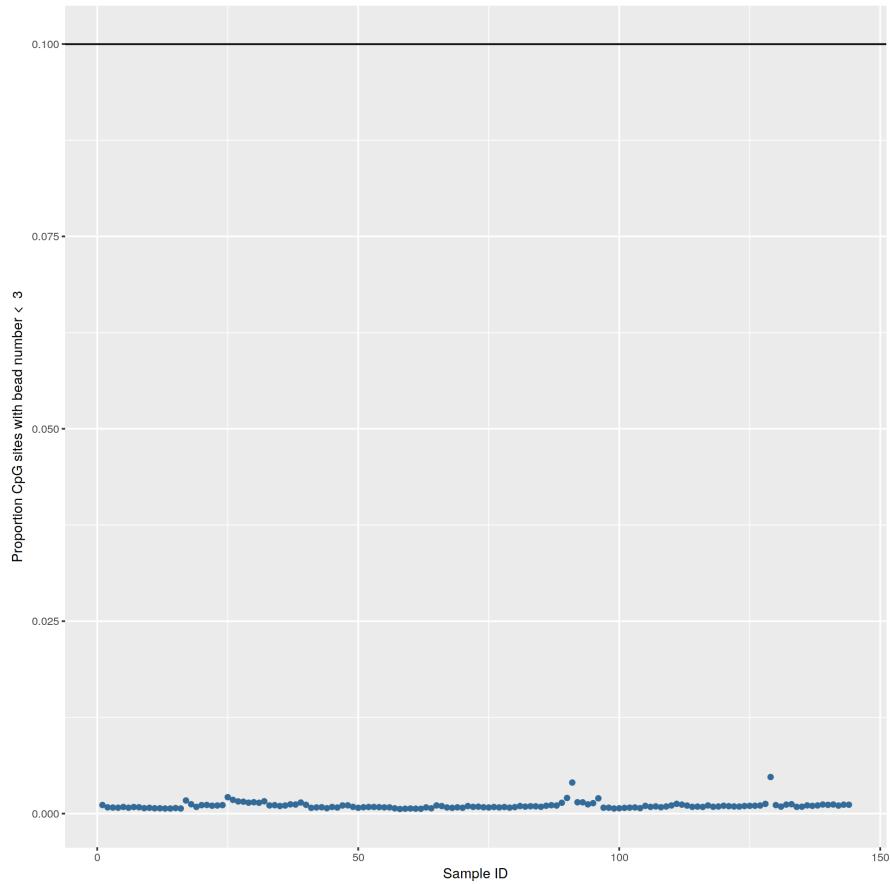


Figure 3.7: Proportion of probes with a bead count of  $< 3$  by sample for the EPIC array data. Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil` QC report.

**3.4.1.1.2 450k Arrays** There were no mismatches between predicted and annotated sex (Figure 3.8). There were not any samples with outliers in their methylated / unmethylated probe proportions (Figure 3.9). No samples were excluded for having a higher than expected proportion of undetected probes (proportion of probes with detection p-value  $> 0.01$  is  $> 0.1$ ), (Figure 3.10). No samples were excluded for having a high proportion of probes with low bead counts (proportion of probes with bead number  $< 3$  is  $> 0.1$ ), (Figure 3.11).

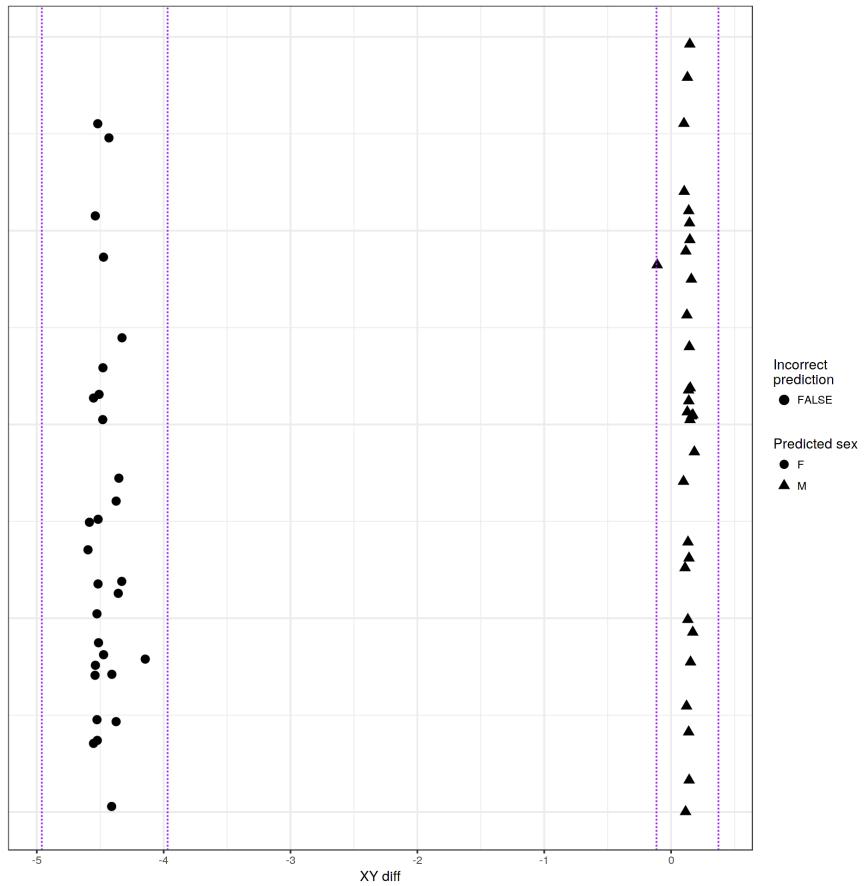


Figure 3.8: Predicted sex of each sample based on the sex chromosome copy numbers inferred from probe intensities for the 450k array data. No predicted sex values differ from their annotations. Plot generated by `meffil` QC report.

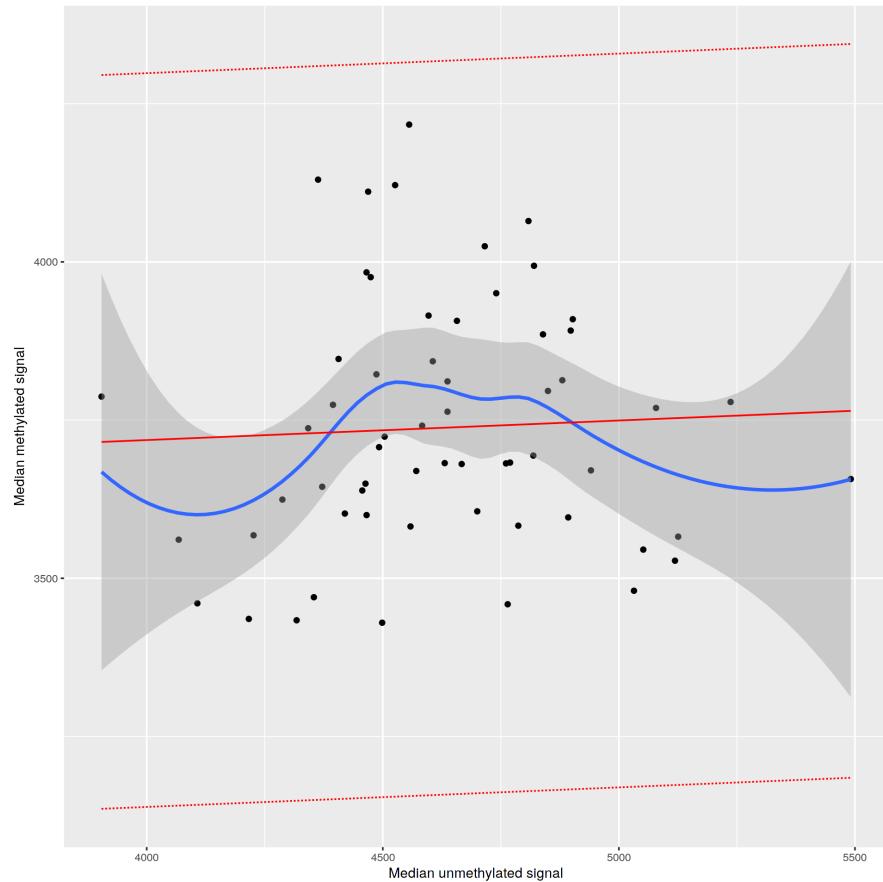


Figure 3.9: Median methylated signal vs unmethylated signal per sample for the 450k array data, solid red line indicates linear regression of median methylated signal vs median unmethylated signal with dotted red lines representing  $3\sigma$  from the expected mean. Samples outside the expected range would be indicated in the legend. Plot generated by `meffil` QC report.

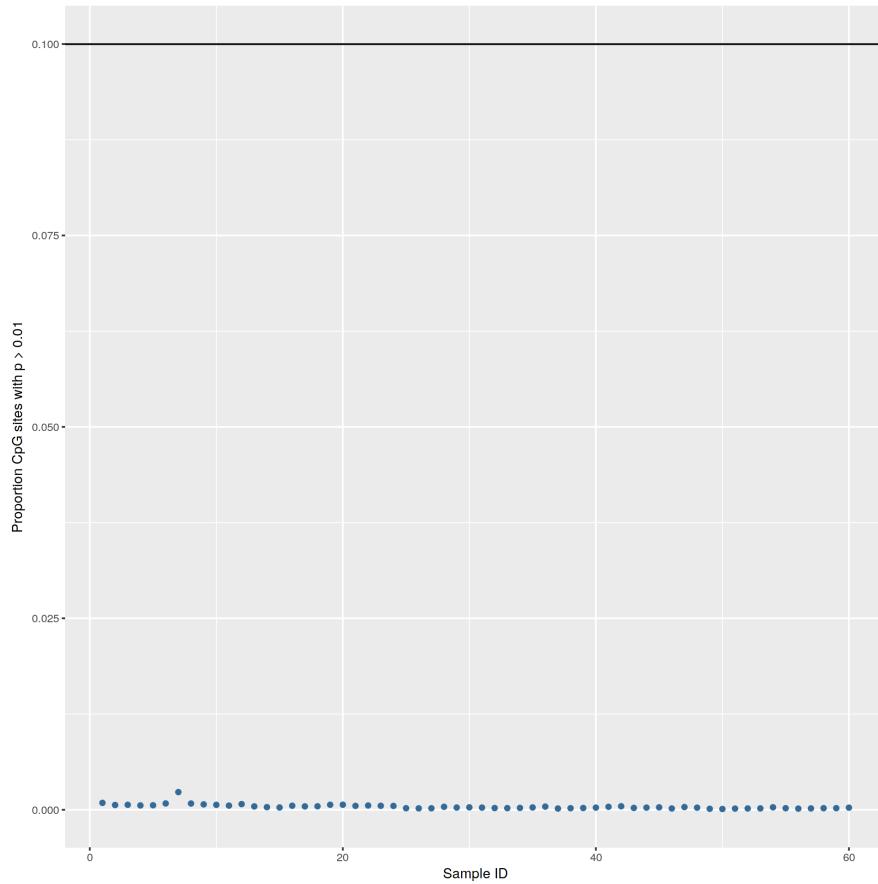


Figure 3.10: Proportion of probes with detection p-values  $>0.01$  by sample for the 450k array data. Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil` QC report.

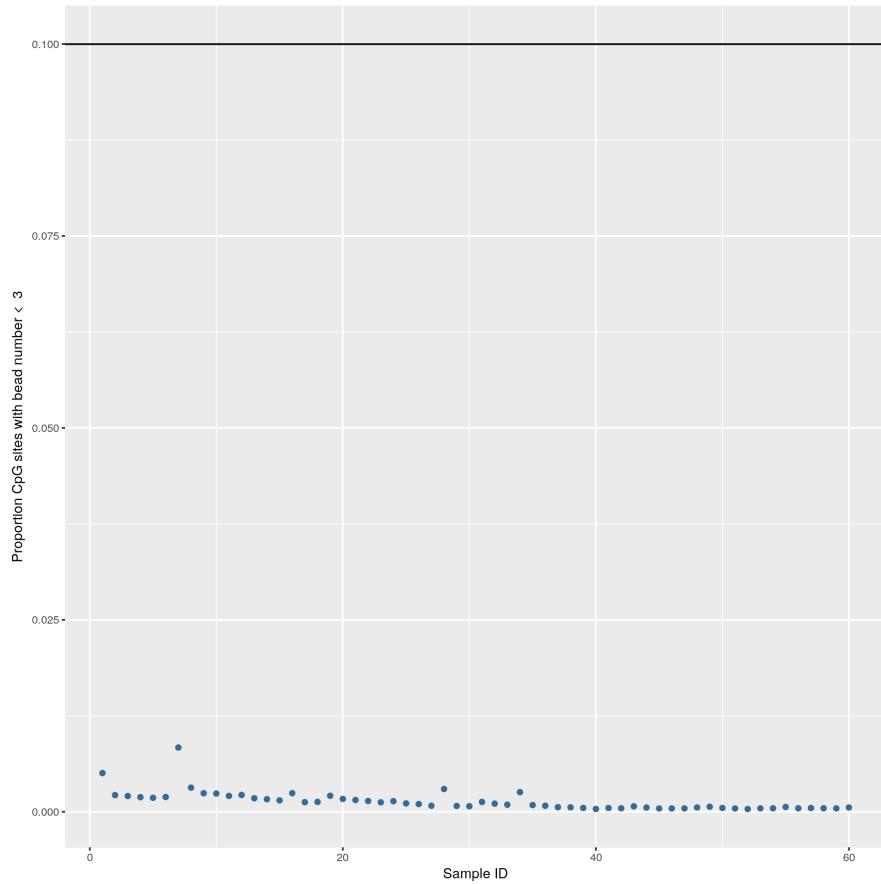


Figure 3.11: Proportion of probes with a bead count of  $< 3$  by sample for the 450k array data. Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil` QC report.

### 3.4.1.2 Probe QC

**3.4.1.2.1 Probe QC - EPIC Arrays** There were no outliers within the control probes Figure 3.12. 1,626 probes were excluded for having high background signal in a high proportion of samples (proportion of samples with detection p-value  $> 0.01$  is  $> 0.1$ ), (Figure 3.13). 162 probes were excluded for having low bead count in a high proportion of samples (proportion of samples with bead number  $< 3$  is  $> 0.1$ ), (Figure 3.14). Probes with poor technical quality were excluded from the analysis prior to functional normalisation.

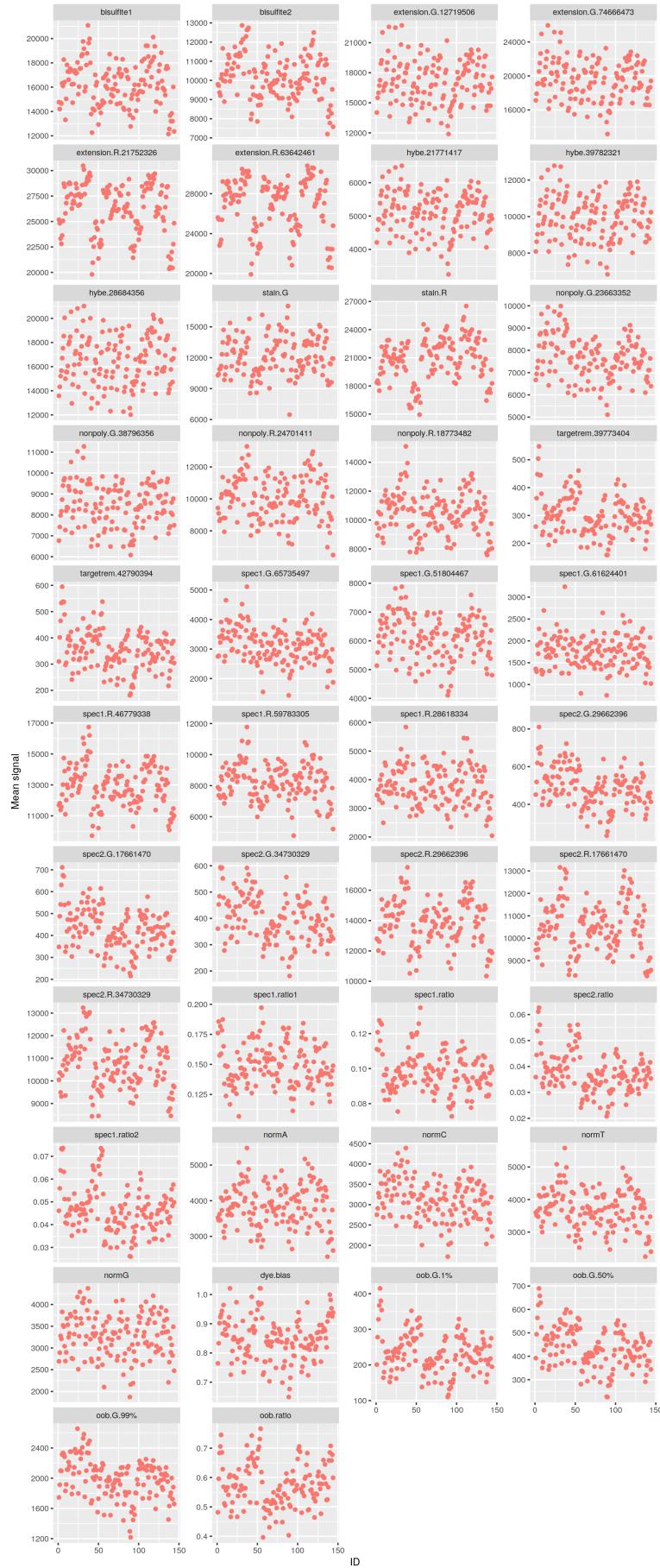


Figure 3.12: Control probe signal by sample for each summary group for the EPIC data. Outliers would be circled in black. Plot generated by `meffil` QC report.

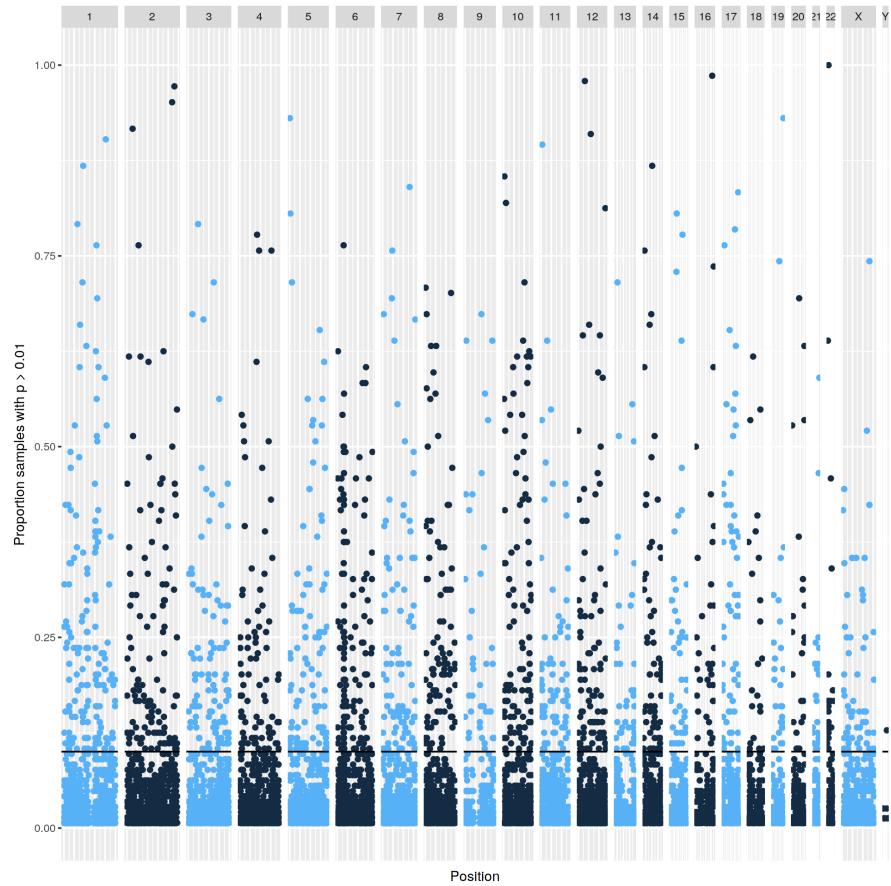


Figure 3.13: Undetectable probes across samples for EPIC data. Manhattan plot showing proportion of samples (y) in which a given probe (x) is not distinguishable from background noise, i.e. a detection p-value of  $> 0.01$ . Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil` QC report.

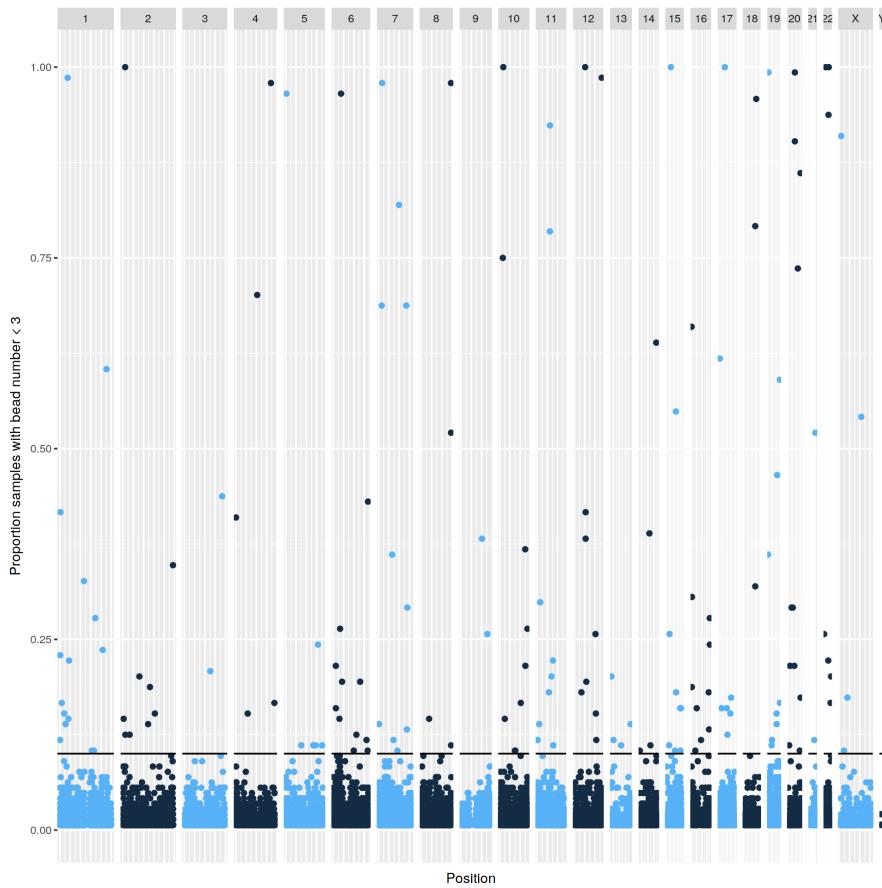


Figure 3.14: Low bead count probes across samples for EPIC data. Manhattan plot showing the proportion of samples (y) in which a given probe (x) has a bead count of  $< 3$ . Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil` QC report.

Problematic probes identified by Pidsley et al. [185] and those overlapping the regions identified by Bell et al. [84] were excluded from subsequent analysis after functional normalisation. This, including the poor quality probes, is a total of 57,396 unique probes excluded from the analysis (~6.62% of the total number of probes). Gap hunter identified a further 77,398 probes (8.9% of all probes) which might be subject to genetic confounding (Figure 3.15).

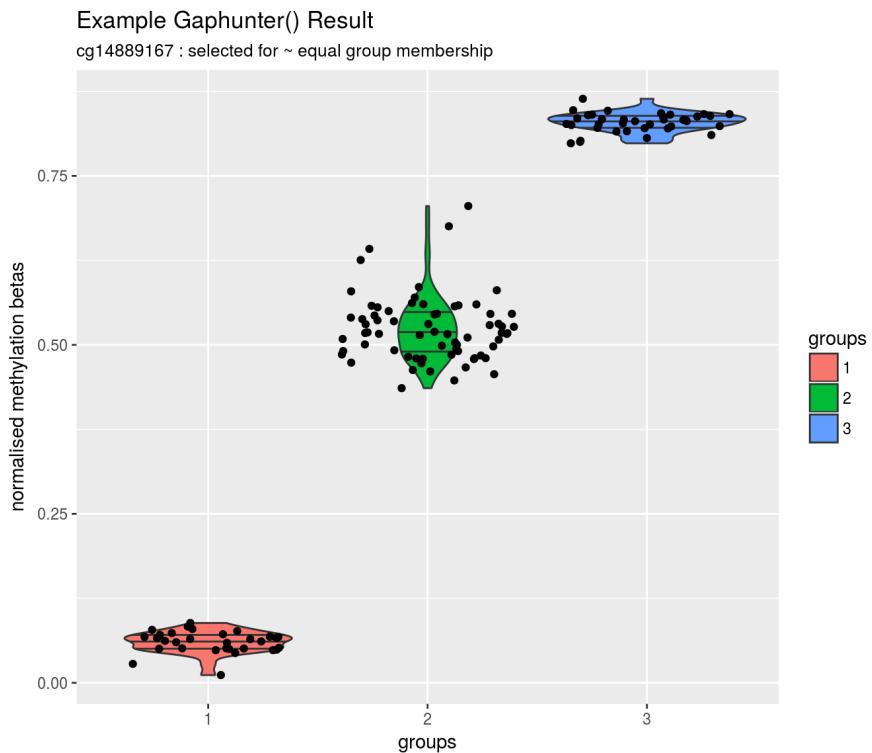


Figure 3.15: An example of the DNAm distribution for a result from the `gaphunter()` function. This is an example chosen to best exemplify the sort of result which is strongly suggestive of a genetic variant with an impact on methylation status acting on this site. It is unrepresentative of typical results from `gaphunter()` in that the groups have a relatively even membership, many results have a small number of individuals in one or more groups making it hard to distinguish methylation outliers caused by rarer genetic variants from those with other causes.

**3.4.1.2.2 Probe QC - 450k Arrays** 509 probes were excluded for having high background signal in a high proportion of samples (proportion of samples with detection p-value  $> 0.01$  is  $> 0.1$ ), (Figure 3.17). 1037 probes were excluded for having low bead count in a high proportion of samples (proportion of samples with bead number  $< 3$  is  $> 0.1$ ) (Figure 3.18). There was one sample (MAVIDOS ID 2183) with an outlier within the control probes, a dinitrophenyl labelled staining control probe, thus it was not excluded as only outliers in dye bias and bisulfite conversion control probes were deemed sufficient grounds for excluding a sample (Figure 3.16, Detailed contol probe description [184] p222).

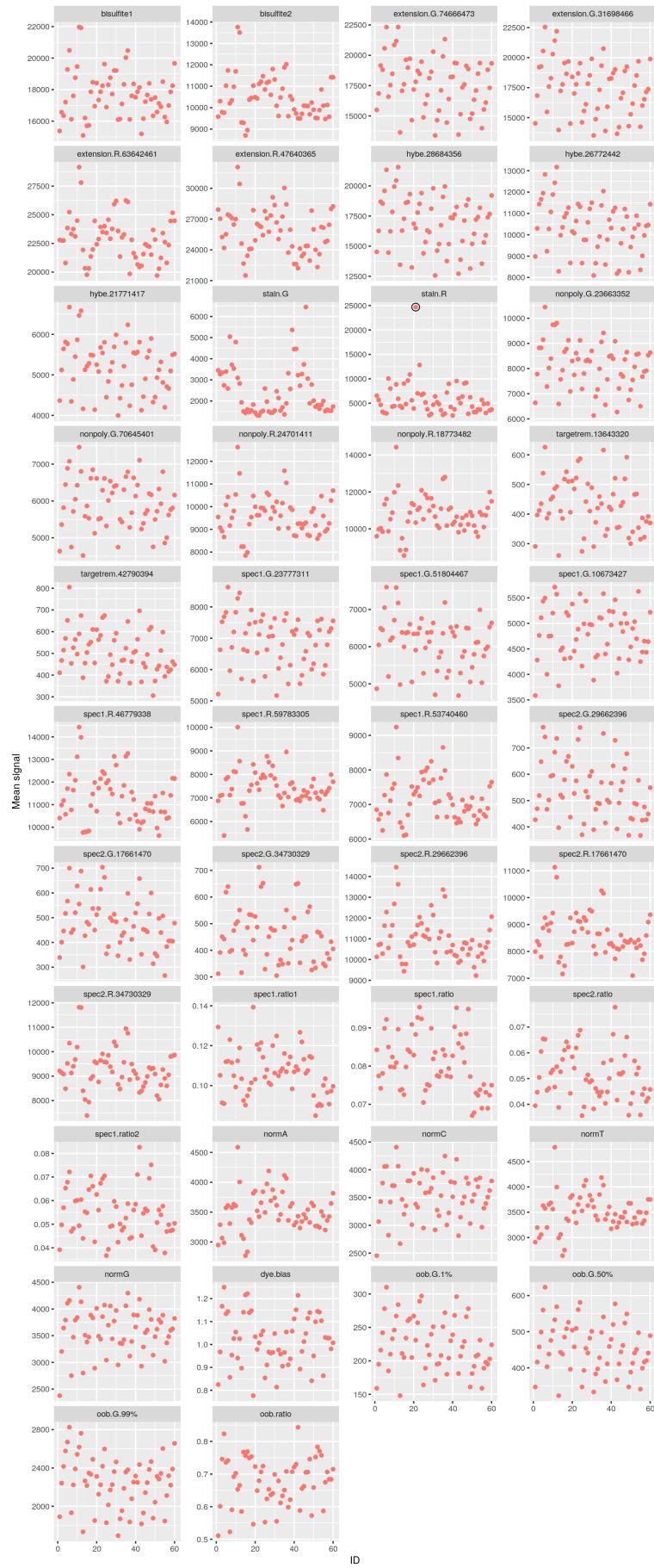


Figure 3.16: Control probe signal by sample for each summary group for the 450k data. Outliers would be circled in black. Plot generated by `meffil` QC report.

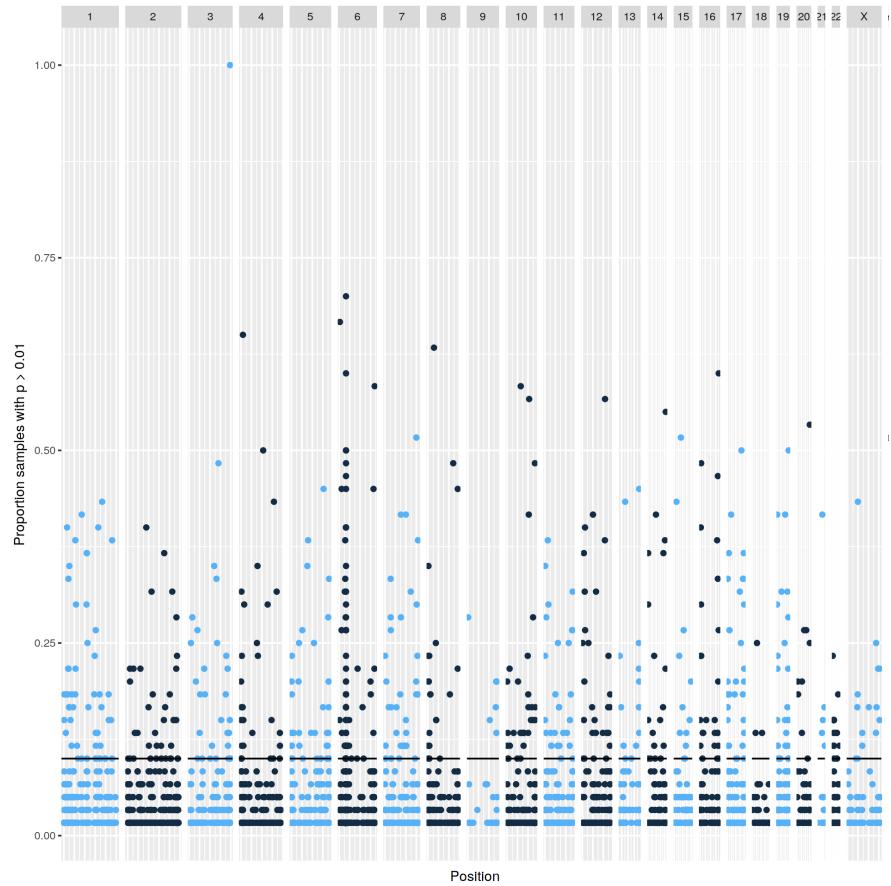


Figure 3.17: Undetectable probes across samples for 450k data. Manhattan plot showing proportion of samples (y) in which a given probe (x) is not distinguishable from background noise, i.e. a detection p-value of  $> 0.01$ . Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil` QC report.

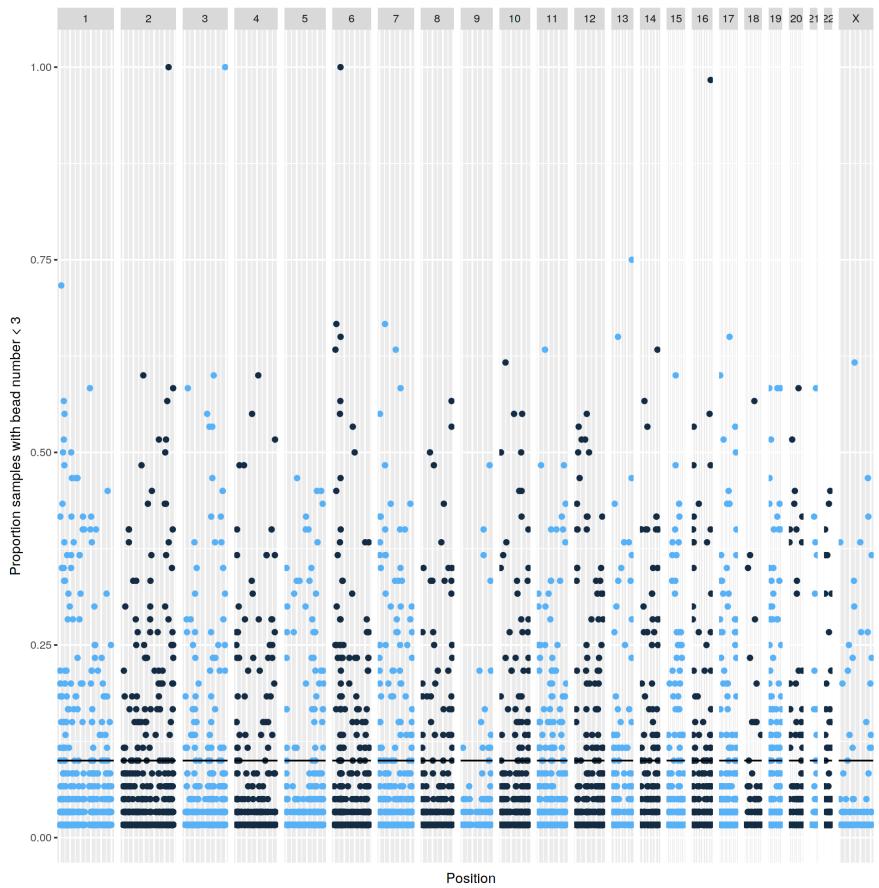


Figure 3.18: Low bead count probes across samples for 450k data. Manhattan plot showing the proportion of samples (y) in which a given probe (x) has a bead count of  $< 3$ . Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil` QC report.

All probes on the ‘general mask’ list from Zhou et al. [186] were excluded from the analysis following functional normalisation, leaving a total of 418,632 probes for subsequent analysis.

### 3.4.1.3 Functional Normalisation

An  $m$  of 6 was chosen for the EPIC arrays as this value produced the last steep drop in residual variation, see Figure 3.19. An  $m$  of 6 was chosen for the 450k arrays as this value produced the last steep drop in residual variation, see Figure 3.20.

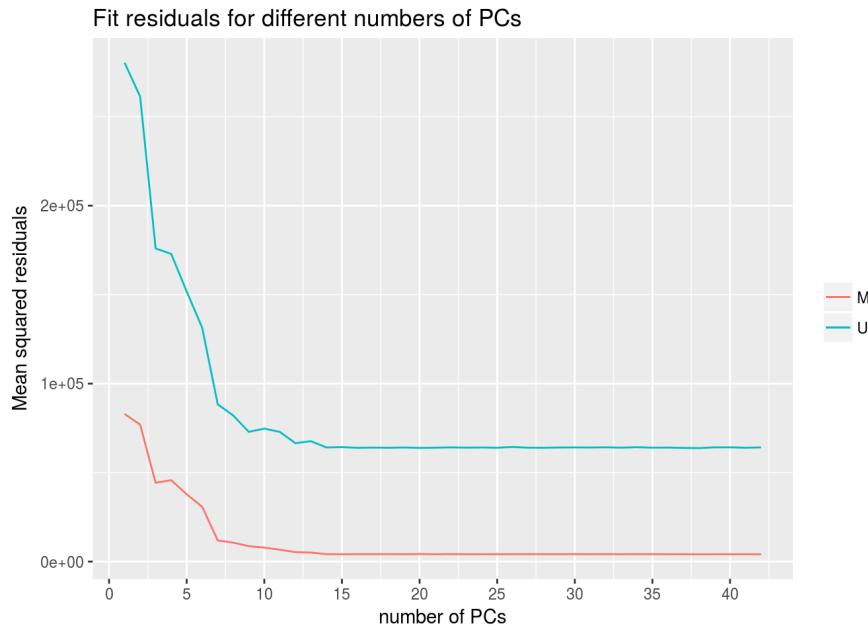


Figure 3.19: Residual variation remaining after functional normalisation of the top 20,000 most variable probes with  $m$  PCs from the control probe summary matrices for the EPIC array samples ( $n=137$ ), for M = methylated and U = unmethylated probes.

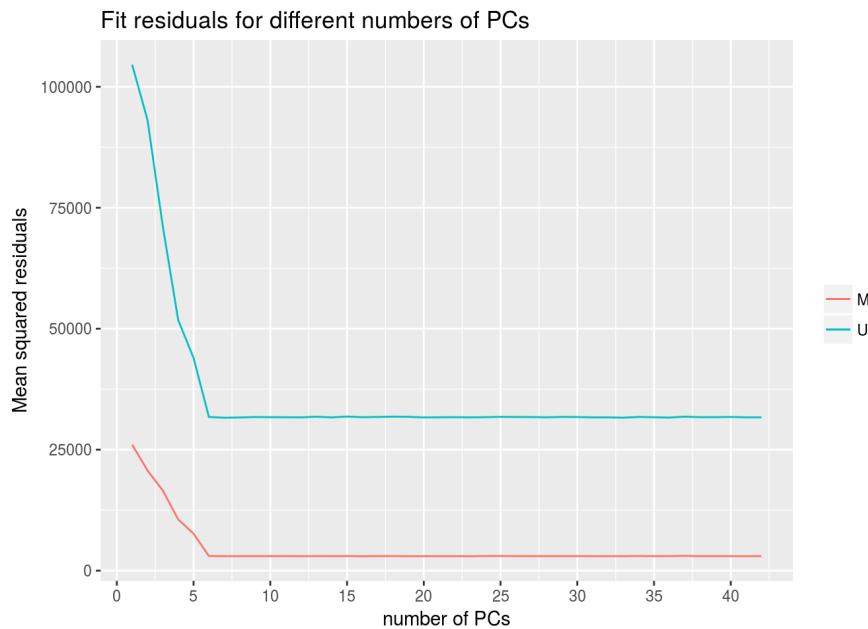


Figure 3.20: Residual variation remaining after functional normalisation of the top 20,000 most variable probes with  $m$  PCs from the control probe summary matrices for the 450k array samples ( $n=60$ ), for M = methylated and U = unmethylated probes.

### 3.4.1.4 EWASs

All EWAS performed below were also performed in exactly the same fashion for the 60 450k arrays, none of these results were significant and they are not included here.

**3.4.1.4.1 Neonatal Bone Mineral Content** Figure 3.21 illustrates the distribution of Neonatal Bone Mineral Content, the outcome on which this EWAS was performed.

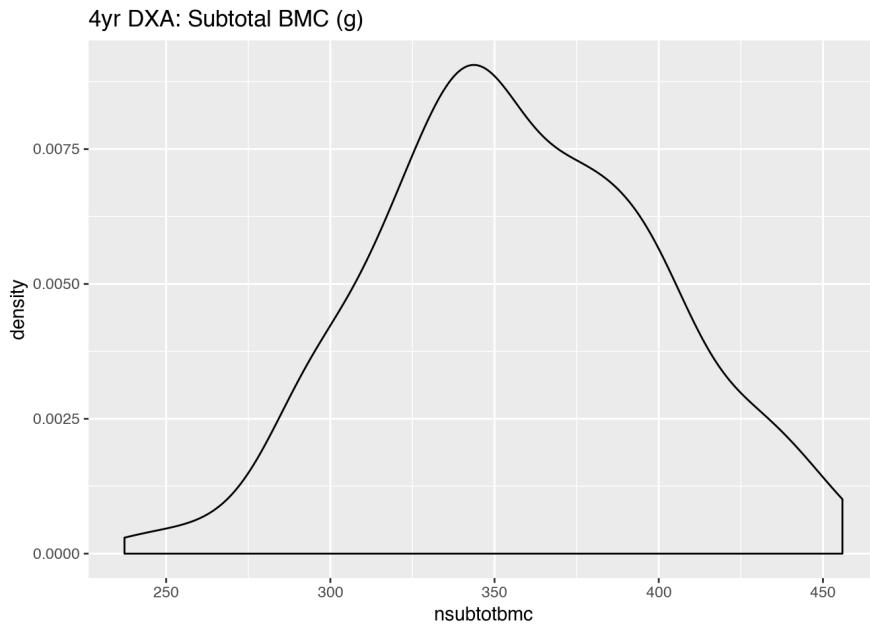


Figure 3.21: Distribution of Neonatal Bone Mineral Content (g) for individuals in the EWAS.

No probes fell below the Bonferroni corrected significance threshold for an association between DNA methylation at that locus and neonatal bone mineral content, Figure 3.22. Sex and sample age at DXA were included as covariates in the ‘all’ model. SVA generated 5 significant surrogate variables which were additionally used in the SVA model.

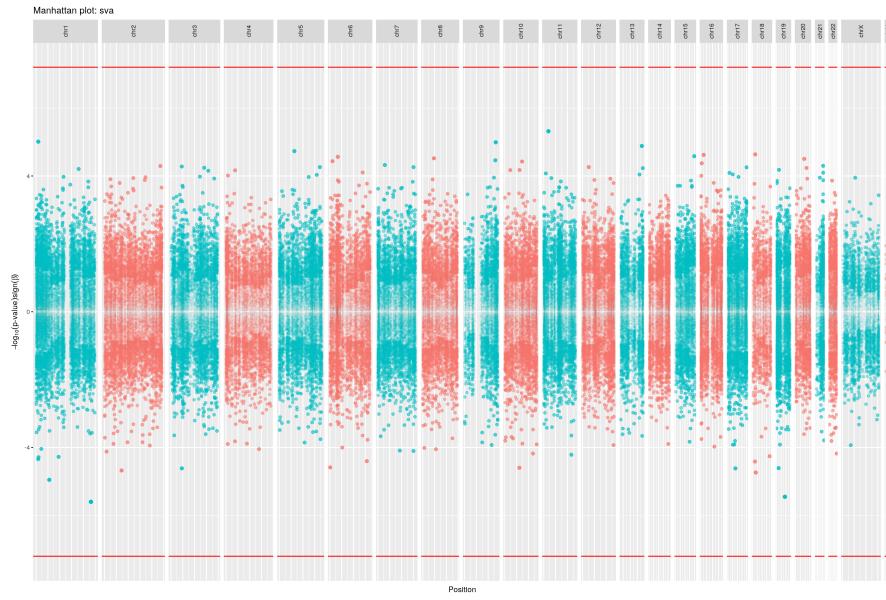


Figure 3.22: Results of EWAS for neonatal bone mineral content with SVA model. Bidirectional Manhattan plot on which  $-\log_{10}(p\text{-value})$  is plotted on the y axis and the sign of this value represents the direction of change. Size and transparency of points increases with  $-\log_{10}(p\text{-value})$  such that the most significant probes are represented by the largest and least translucent points. x axis represents chromosomes and position thereupon. Red line indicates the significance threshold of  $6.18 \times 10^{-8}$  ( $0.05 \div 808,585$ ).

**3.4.1.4.2 Intervention / Placebo** No probes fell below the Bonferroni corrected significance threshold for an association between DNA methylation at that locus and Intervention/placebo group status, Figure 3.23. Sex and sample age at DXA were included as covariates in the ‘all’ model. SVA generated 5 significant surrogate variables which were additionally used in the SVA model.

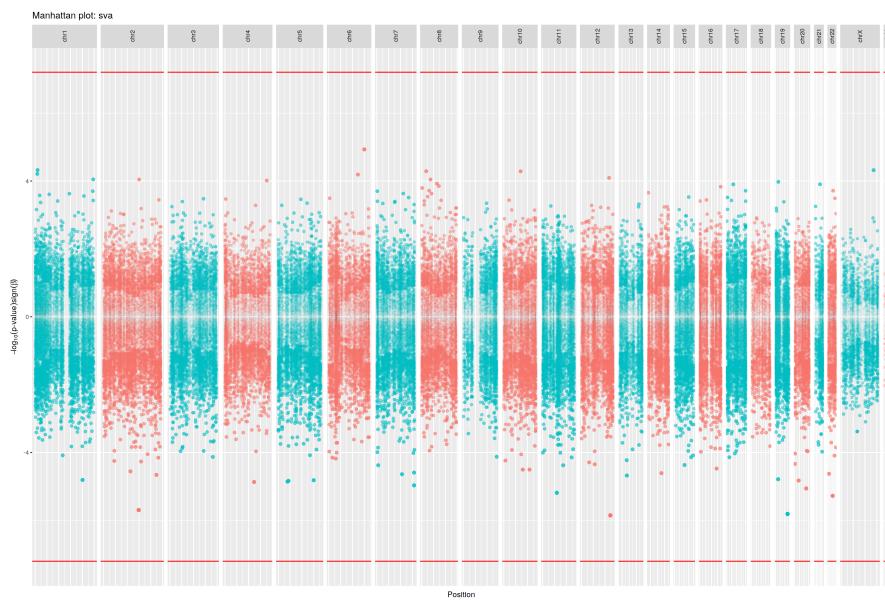


Figure 3.23: Results of EWAS for intervention/placebo group status with SVA model. Bidirectional Manhattan plot on which  $-\log_{10}(p\text{-value})$  is plotted on the y axis and the sign of this value represents the direction of change. Size and transparency of points increases with  $-\log_{10}(p\text{-value})$  such that the most significant probes are represented by the largest and least translucent points. x axis represents chromosomes and position thereupon. Red line indicates the significance threshold of  $6.18 \times 10^{-8}$  ( $0.05 \div 808,585$ ).

**3.4.1.4.3 Maternal Vitamin D (34wks)** Maternal Vitamin D levels remain substantially overlapping between intervention and placebo groups at 34wks, see Figure 3.24. Thus maternal vitamin D at 34wks may prove a more useful variable to model than intervention/placebo status.

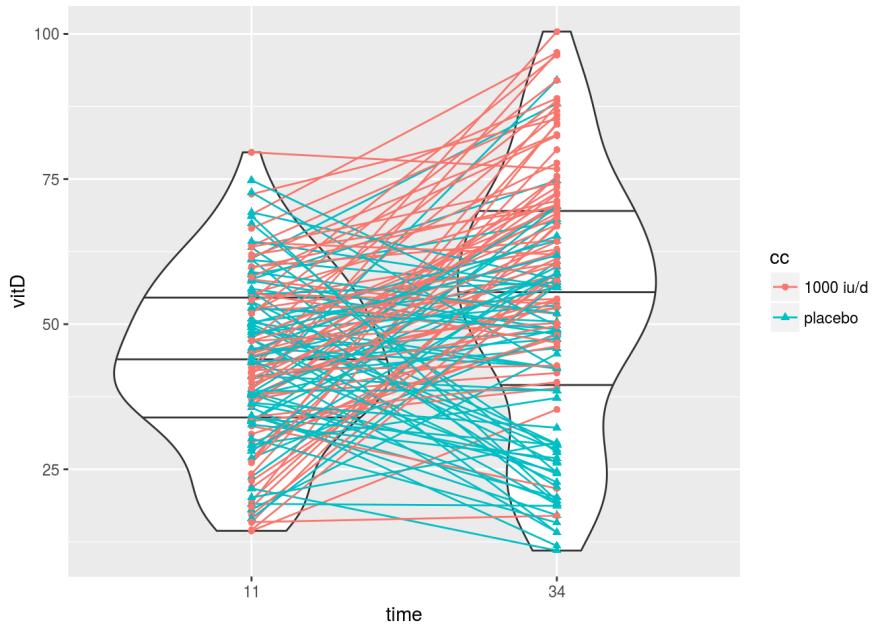


Figure 3.24: Maternal circulating 25(OH)-vitamin D levels (nmol/l) at 11 and 34wks gestation, supplementation with 1000 IU/day cholecalciferol began at week 14. Each participant is shown at both time points linked by a line to indicate the direction of change. The violin plots indicate the density of the distribution of vitamin D values at each time point with the 25<sup>th</sup>, 50<sup>th</sup>, 75<sup>th</sup> quantiles indicated with horizontal black lines. The colour indicates Intervention (Red) / Placebo (Blue) group

No probes fell below the Bonferroni corrected significance threshold for an association between DNAm at that locus and maternal vitamin D levels at 34wks gestation, (Figure 3.25). Sex and sample age at DXA were included as covariates in the ‘all’ model. SVA generated 5 significant surrogate variables which were additionally used in the SVA model.

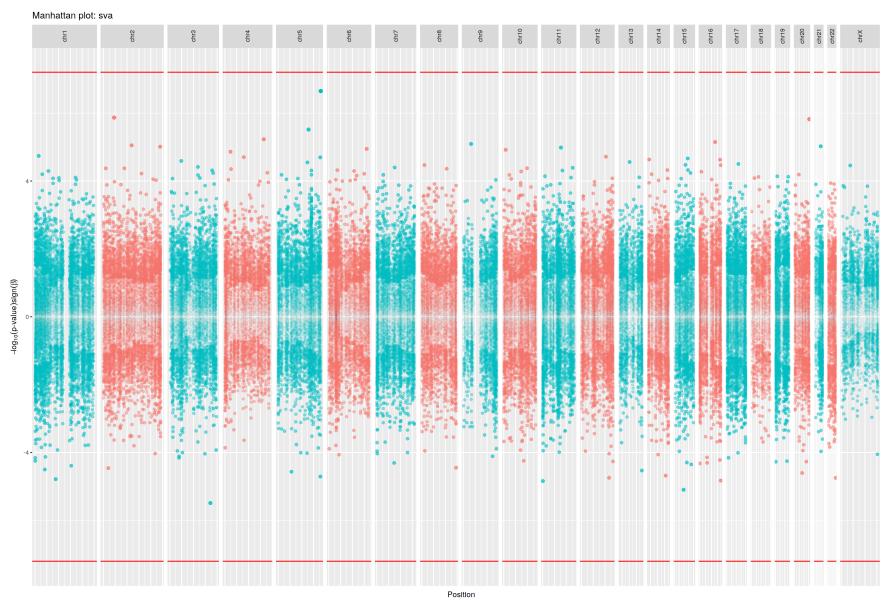


Figure 3.25: Results of EWAS for maternal circulating 25(OH)-vitamin D levels (nmol/l) levels at 34wks gestation with SVA model. Bidirectional Manhattan plot on which  $-\log_{10}(p\text{-value})$  is plotted on the y axis and the sign of this value represents the direction of change. Size and transparency of points increases with  $-\log_{10}(p\text{-value})$  such that the most significant probes are represented by the largest and least translucent points. x axis represents chromosomes and position thereupon. Red line indicates the significance threshold of  $6.18 \times 10^{-8}$  ( $0.05 \div 808,585$ ).

**3.4.1.4.4 Change in Maternal Vitamin D** Figure 3.26 illustrates the change in maternal vitamin D from 11 to 34 weeks gestation. No probes fell below the Bonferroni corrected significance threshold for an association between DNAm at that locus and change in maternal vitamin D from 11 to 34wks gestation, Figure 3.27. SVA generated 5 significant surrogate variables which were additionally used in the SVA model.

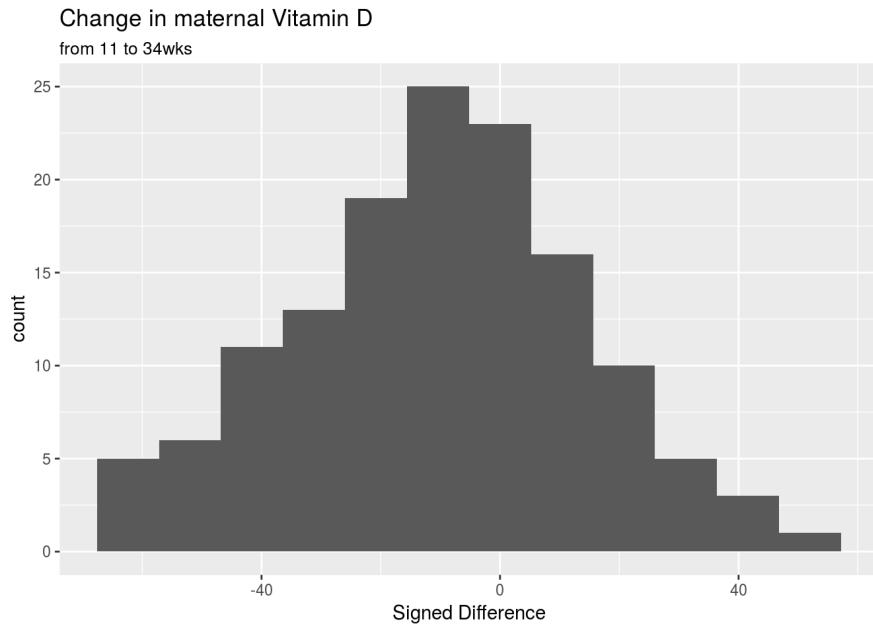


Figure 3.26: Distribution of the changes in maternal circulating 25(OH)-vitamin D levels (nmol/l) levels from 11 to 34 weeks gestation.

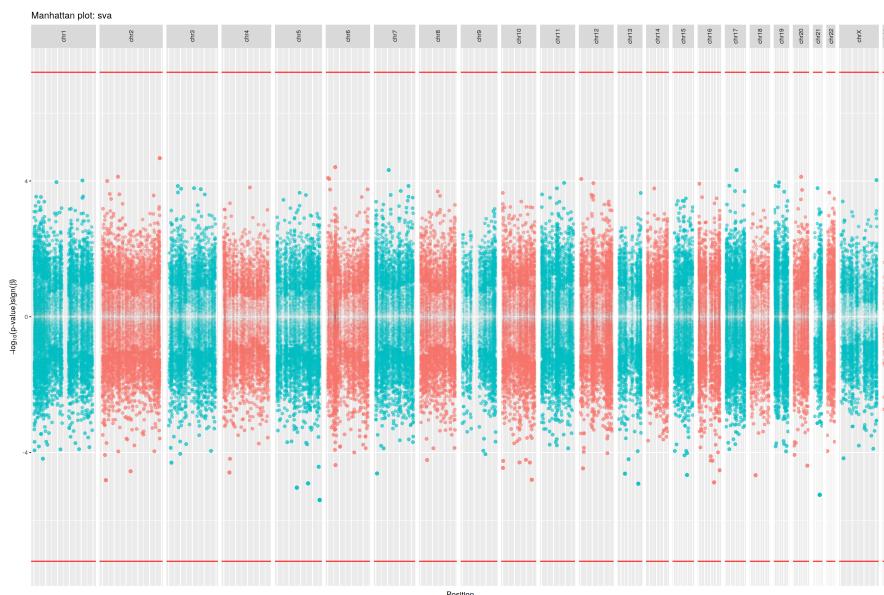


Figure 3.27: Results of EWAS for change in maternal circulating 25(OH)-vitamin D levels (nmol/l) levels from 11 to 34wks gestation with SVA model. Bidirectional Manhattan plot on which  $-\log_{10}(p\text{-value})$  is plotted on the y axis and the sign of this value represents the direction of change. Size and transparency of points increases with  $-\log_{10}(p\text{-value})$  such that the most significant probes are represented by the largest and least translucent points. x axis represents chromosomes and position thereupon. Red line indicates the significance threshold of  $6.18 \times 10^{-8}$  ( $0.05 \div 808,585$ ).

### 3.4.1.5 Concordance of EPIC and 450k EWAS results

Concordance between the rankings of probes would suggest that the EWAS may be capturing a ‘real’ signal that is simply below the significance threshold with the statistical sensitivity/power that is available in this dataset. The concordance between the EPIC and 450k datasets (Figures 3.28 & 3.29) appears to be at roughly the level expected by chance. This does not lend support to the possibility that there are associations between the variables of interest and DNAm that are beneath the current sensitivity of the study, it does not, however, rule this out. In the absence of any probes above the significance threshold and with poor concordance between the 450k and EPIC array p-value rankings no further analyses such as gene set enrichment and differentially methylated region (DMR) calling have been carried out at this time.

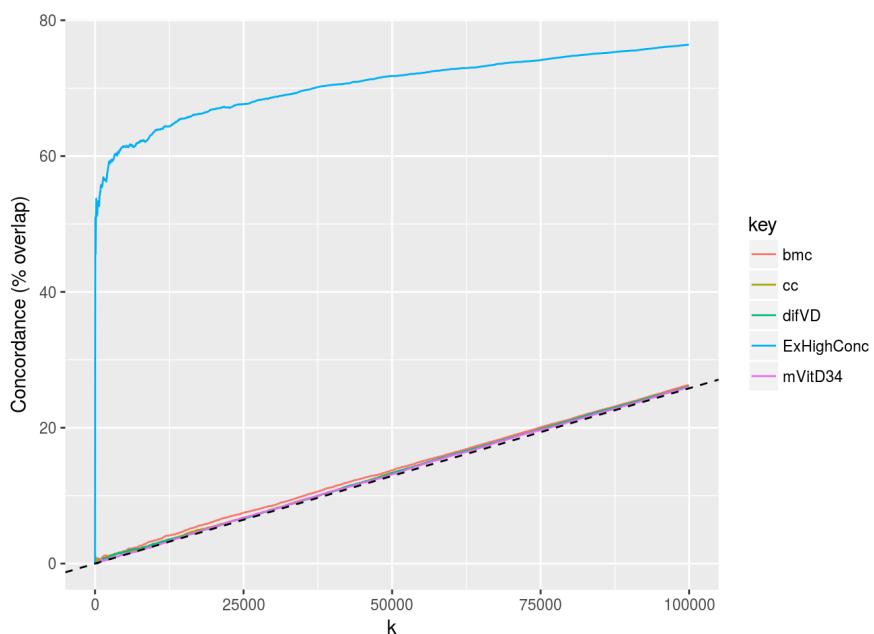


Figure 3.28: Concordance between the top 100,000 probes in common between the EWASs run on the EPIC ( $n=137$ ) and 450k ( $n=60$ ) data sets. bmc = bone mineral content, cc = Intervention / Placebo, difVD = Change in Vitamin D from 11 to 34wks, ExHighConc = Example of High Concordance generated using SVA vs iSVA results for the 450k intervention/placebo EWAS. Dotted line denotes concordance expected by chance (intersects 50% at 387,511, the number of shared probes).

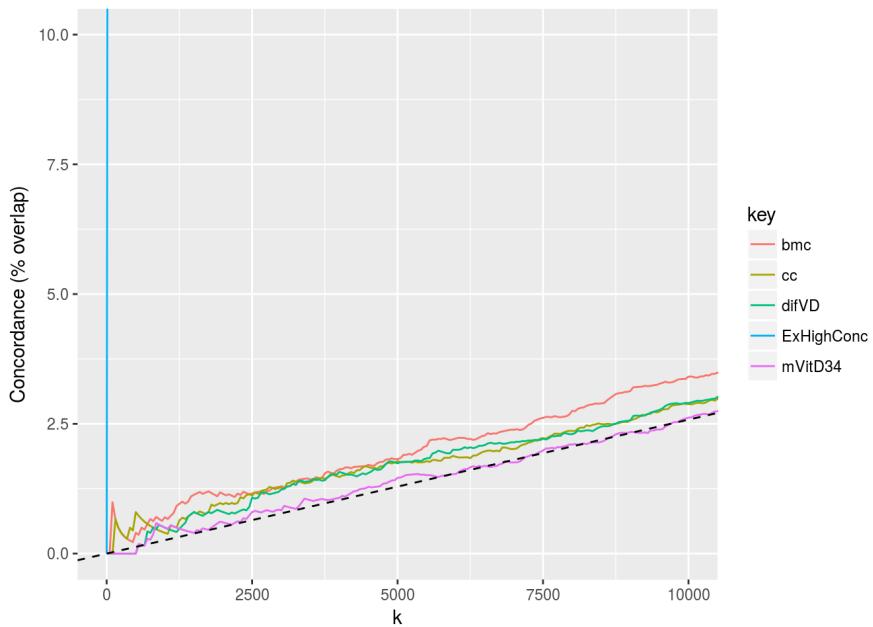


Figure 3.29: Concordance between the top 10,000 probes in common between the EWASs run on the EPIC ( $n=137$ ) and 405k ( $n=60$ ) data sets. bmc = bone mineral content, cc = Intervention / Placebo, difVD = Chance in Vitamin D from 11 to 34wks, ExHighConc = Example of High Concordance generated using SVA vs iSVA results for the 450k intervention/placebo EWAS. Dotted line denotes concordance expected by chance (intersects 50% at 387,511, the number of shared probes).

### 3.4.2 MAVIDOS phase II

DNA methylation at none of the probes was significantly associated with bone mineral content at 4 years.

#### 3.4.2.1 Whole Array QC

The predicted sex of the samples generated using sex chromosome probe intensities was checked against that in the sample annotation and 5 mismatches were found (Figure 3.30). Two samples did not have sex information, 2 had predicted sex discordant with their annotated sex and 1 was ambiguous these were excluded from further analysis.

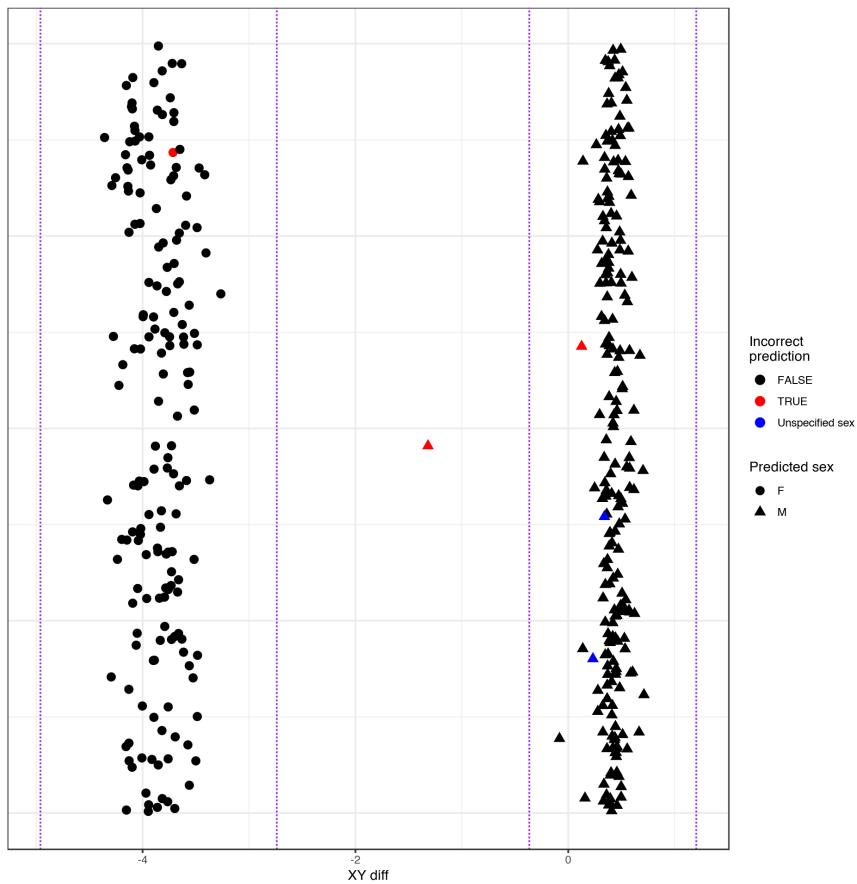


Figure 3.30: Predicted sex of each sample based on the sex chromosome copy numbers inferred from probe intensities for the EPIC array data. Mismatches between the predicted sex and that asserted in the sample annotation metadata are shown in red. Two predicted sex values differ from their annotations. Plot generated by `meffil` QC report.

Arrays: 201516310023 (mavid: 1490), 201516320022 (mavid: 1672), 201530430013 (mavvid: 4090), 201530430015 (mavid: 1903), & 202410280028 (mavid: 2078) were excluded as their median methylated signal was more than  $3\sigma$  from the expected value, (Figure 3.31). No samples were excluded for having a higher than expected proportion of undetected probes (proportion of probes with detection p-value  $> 0.01$  is  $> 0.1$ ) (Figure 3.32). No samples were excluded for having a high proportion of probes with low bead counts (proportion of probes with bead number  $< 3$  is  $> 0.1$ ), (Figure 3.7).

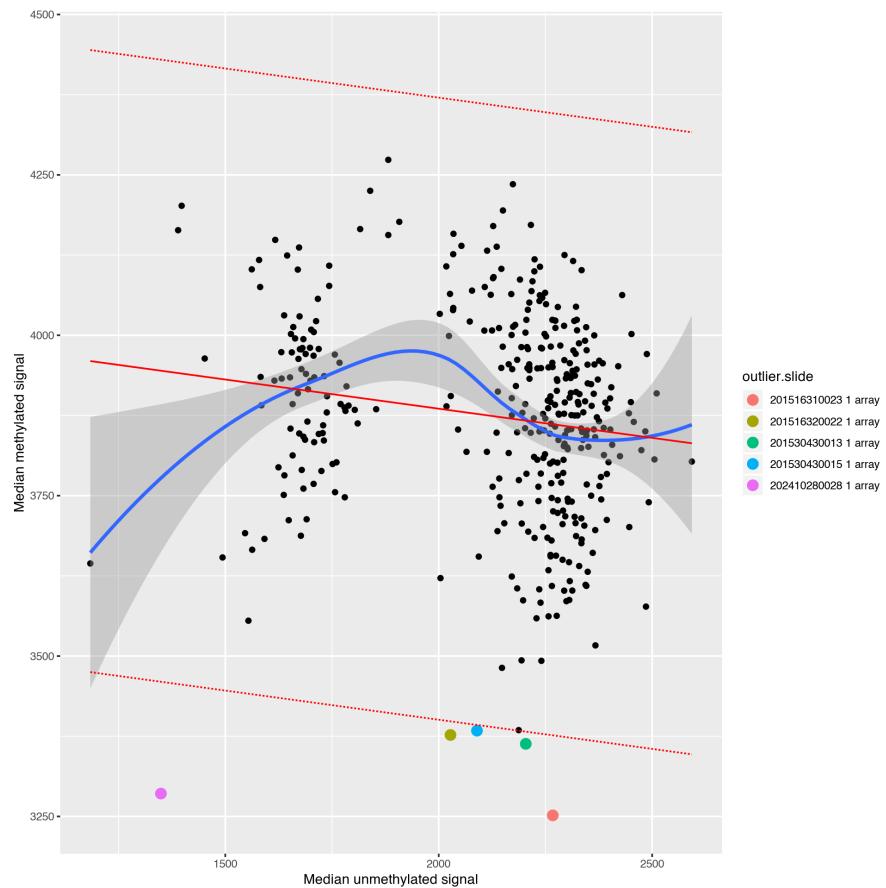


Figure 3.31: Median methylated signal vs unmethylated signal per sample for the EPIC array data, solid red line indicates linear regression of median methylated signal vs median unmethylated signal with dotted red lines representing  $3\sigma$  from the expected mean. Samples outside the expected range are indicated in the legend. Plot generated by `meffil` QC report.

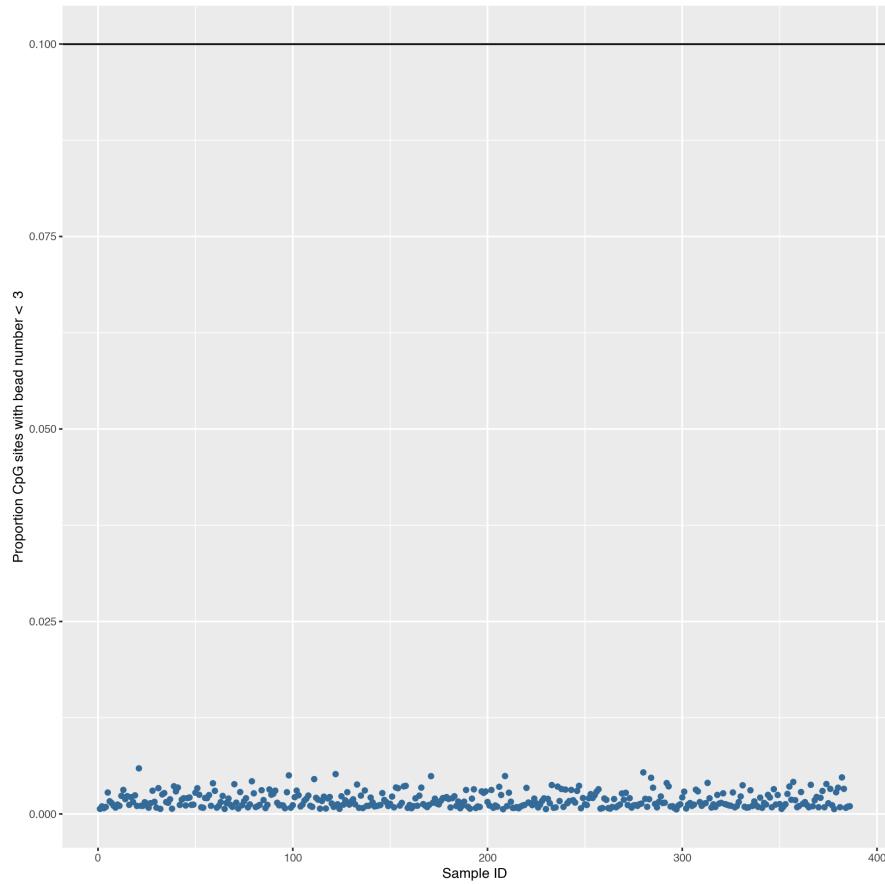


Figure 3.32: Proportion of probes with detection p-values  $>0.01$  by sample for the EPIC array data. Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil` QC report.

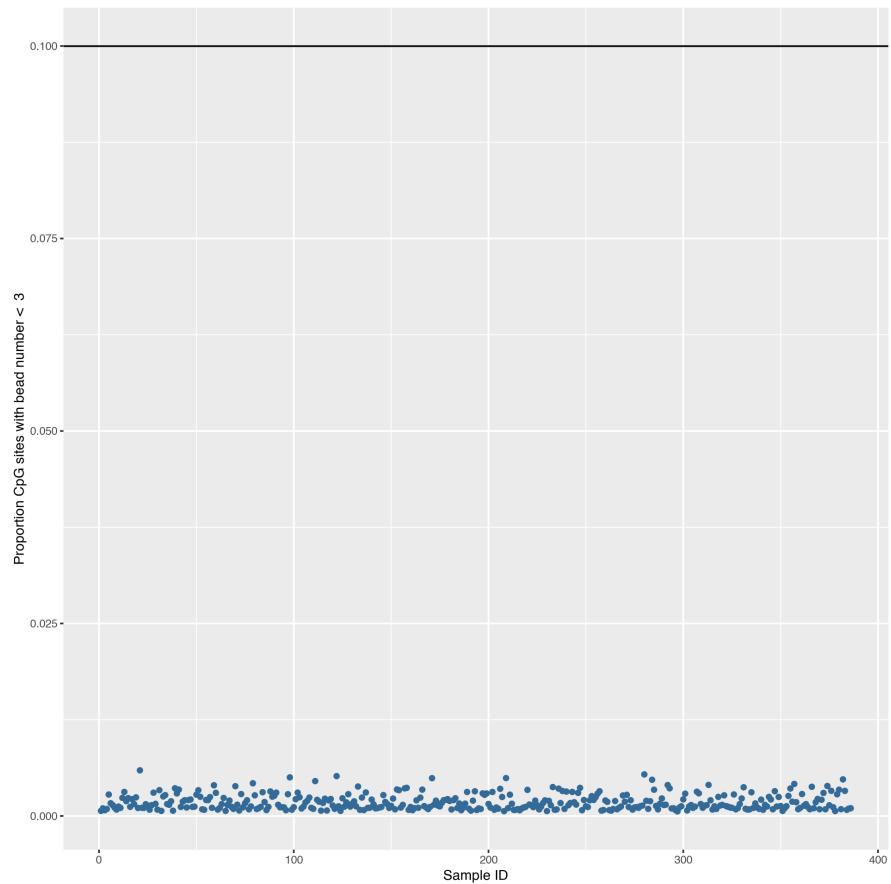


Figure 3.33: Proportion of probes with a bead count of  $< 3$  by sample for the EPIC array data. Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil` QC report.

#### 3.4.2.2 Probe QC

There was one outlier within the control probes, in a non-critical specificity control probe for detecting non-specific methylation detection over an unmethylated background (Figure 3.12). 1,317 probes were excluded for having high background signal in a high proportion of samples (proportion of samples with detection p-value  $> 0.01$  is  $> 0.1$ ), (Figure 3.13). 220 probes were excluded for having low bead count in a high proportion of samples (proportion of samples with bead number  $< 3$  is  $> 0.1$ ), (Figure 3.14). Probes with poor technical quality were excluded from the analysis prior to functional normalisation.

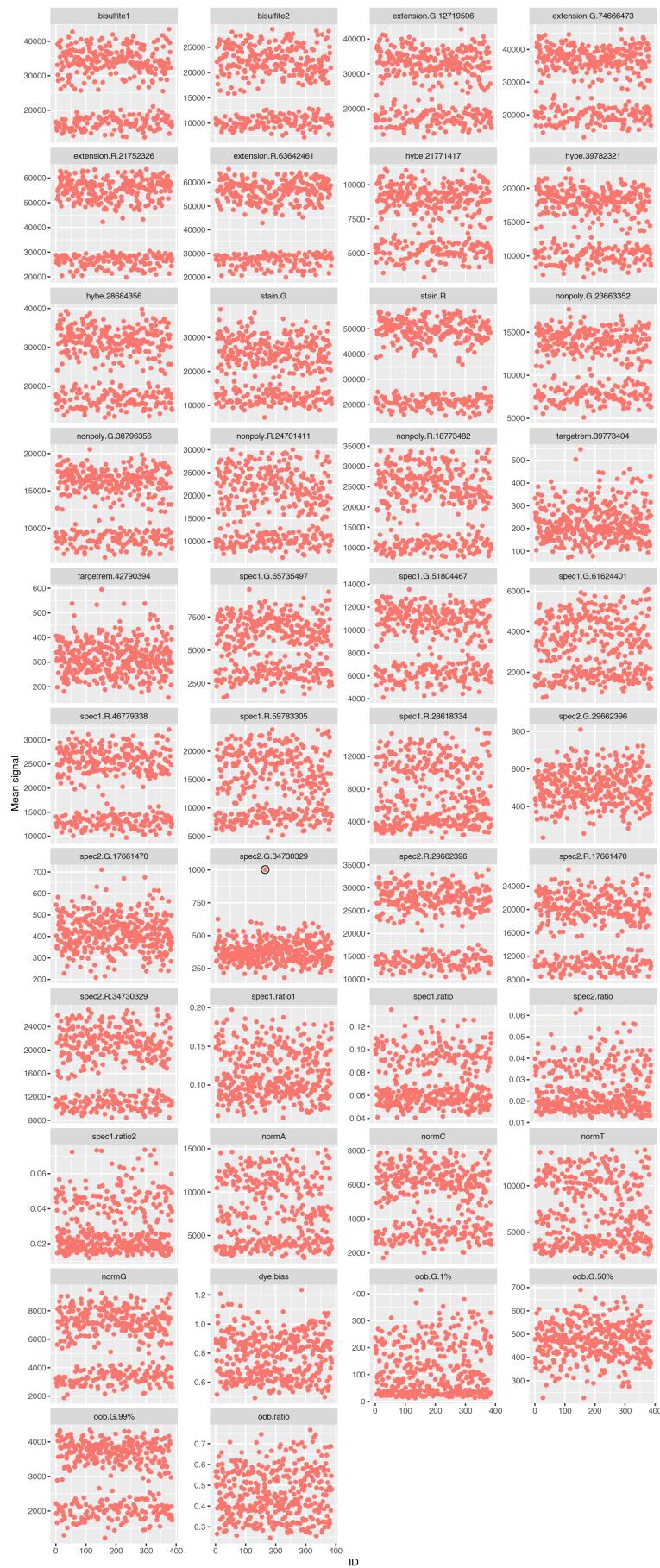


Figure 3.34: Control probe signal by sample for each summary group for the EPIC data. Outliers would be circled in black. Plot generated by `meffil` QC report.

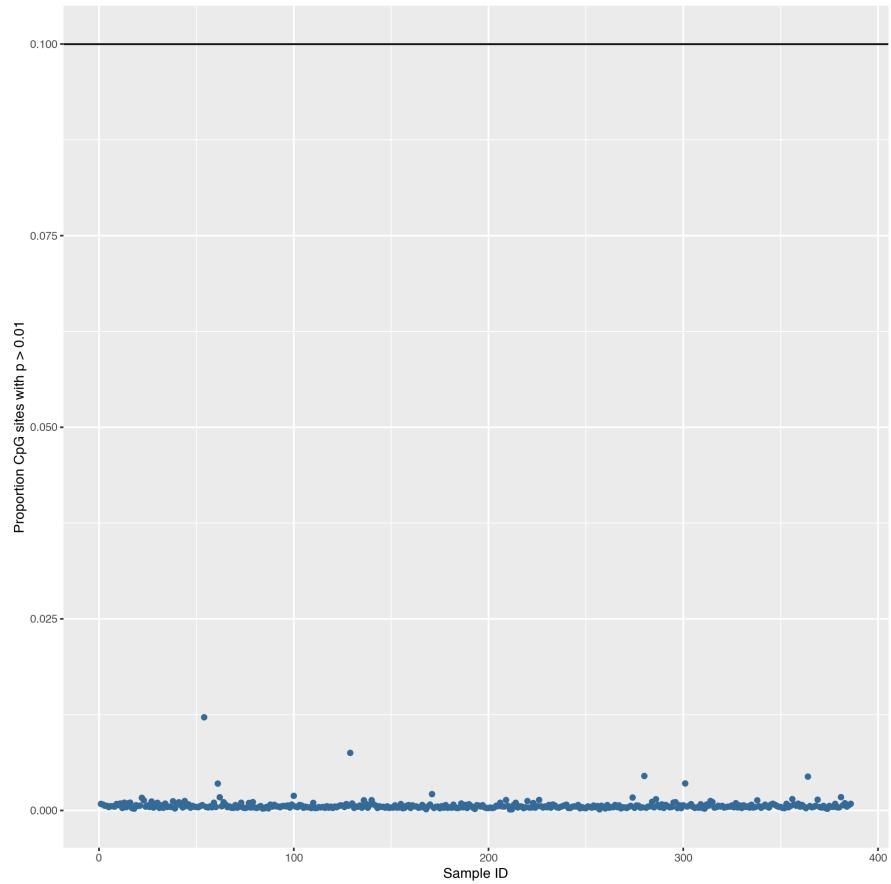


Figure 3.35: Undetectable probes across samples for EPIC data. Manhattan plot showing proportion of samples (y) in which a given probe (x) is not distinguishable from background noise, i.e. a detection p-value of  $> 0.01$ . Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil` QC report.

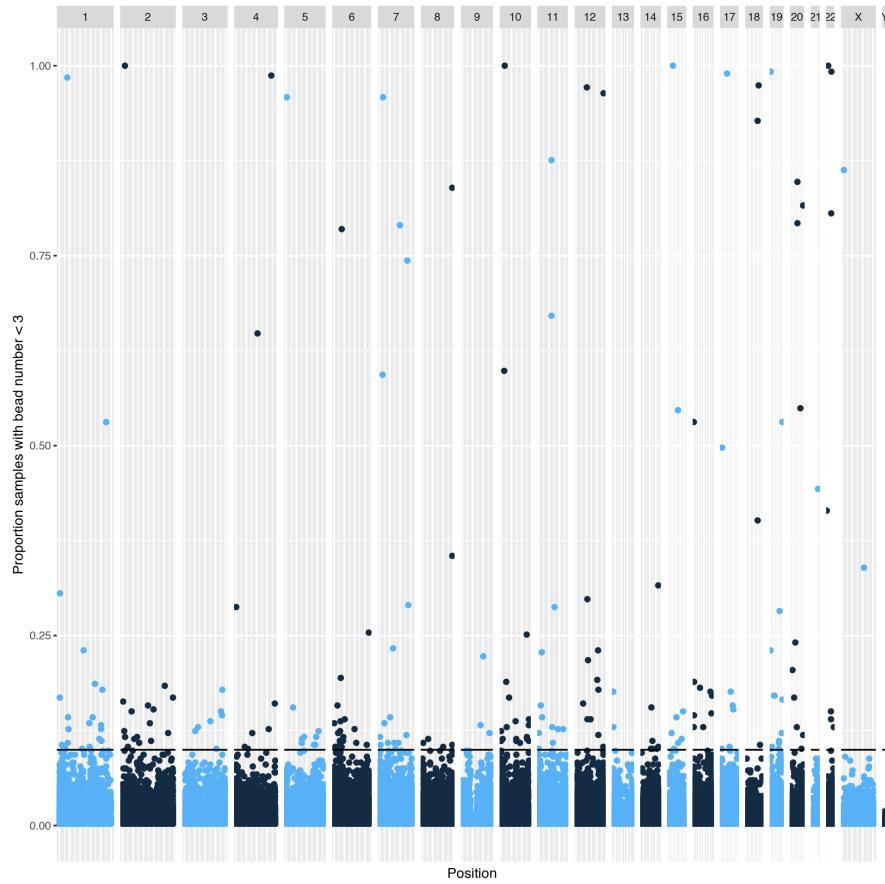


Figure 3.36: Low bead count probes across samples for EPIC data. Manhattan plot showing the proportion of samples (y) in which a given probe (x) has a bead count of  $< 3$ . Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil` QC report.

### 3.4.2.3 Functional Normalisation

An  $m$  of 6 was chosen as this value produced the last steep drop in residual variation, see Figure 3.37.

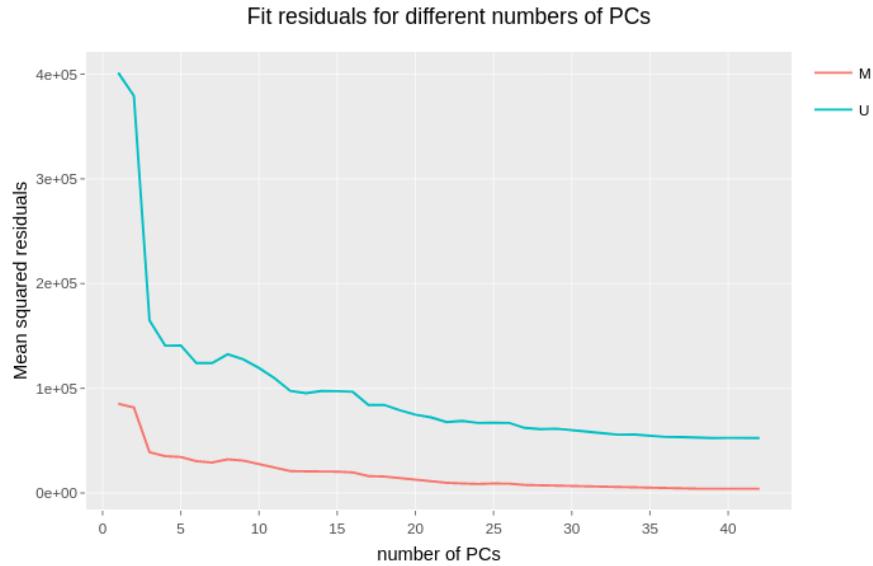


Figure 3.37: Residual variation remaining after functional normalisation of the top 20,000 most variable probes with  $m$  PCs from the control probe summary matrices for the EPIC array samples ( $n=237$ ), for M = methylated and U = unmethylated probes.

#### 3.4.2.4 EWAS

**3.4.2.4.1 Bone Mineral Content at 4 years** The model for this EWAS attempted to predict whole body (minus head) bone mineral content (g) (Figure 3.38) correcting for: Cell type composition (B-cells, CD4+T, CD8+T, Granulocytes, Monocytes, Natural Killer cells, Endothelial cells, Epithelial cells, Stromal cells), Mother's age at birth, Sex, Mother's BMI at 11 weeks gestation, whether the mother smoked during pregnancy, & Gestational Age. The EWAS was performed on 237 individuals. surrogate variable analysis identified 19 significant surrogate variables which were included in the model. No probes fell below the Bonferroni corrected significance threshold for an association between DNAm at that locus and bone mineral content, in any of the models including SVA (Figure 3.39).

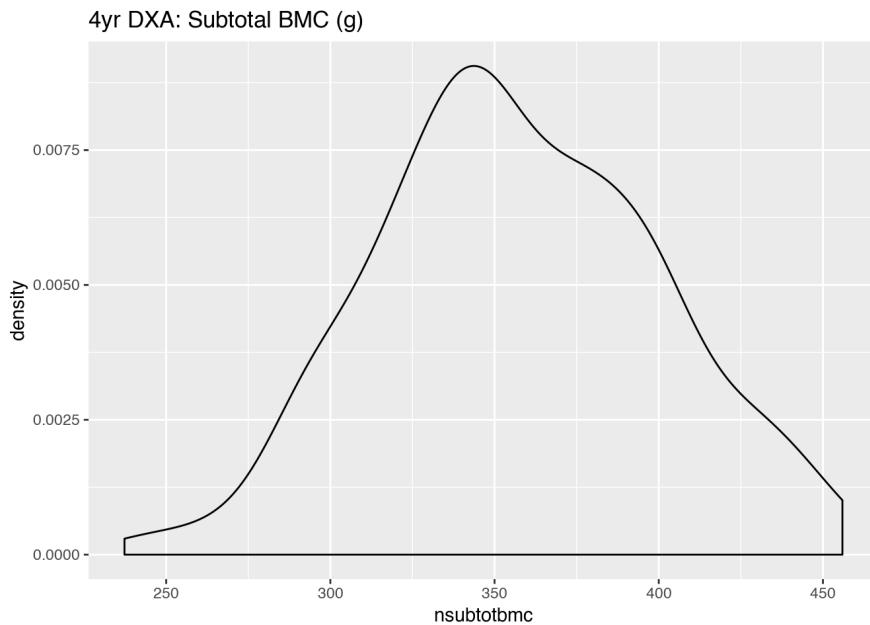


Figure 3.38: Distribution of whole body (minus head) bone mineral content in grams at 4 years of age.

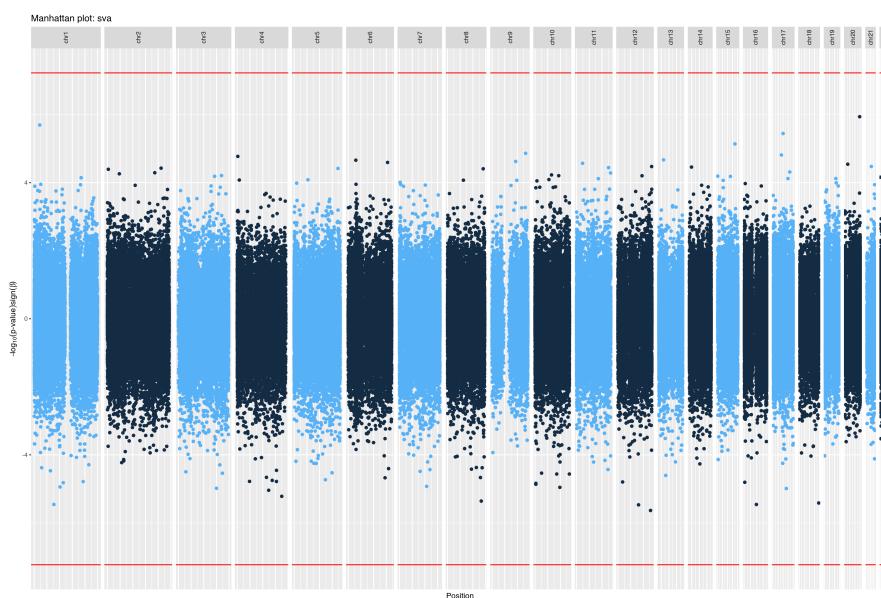


Figure 3.39: Results of EWAS for whole body minus head bone mineral content (g) with SVA model ( $n = 237$ ). Bidirectional Manhattan plot on which  $-\log_{10}(p - \text{value})$  is plotted on the y axis and the sign of this value represents the direction of change. Size and transparency of points increases with  $-\log_{10}(p - \text{value})$  such that the most significant probes are represented by the largest and least translucent points. x axis represents chromosomes and position thereupon. Red line indicates the significance threshold of  $6.18 \times 10^{-8}$  ( $0.05 \div 808,585$ ).

### 3.4.3 Southampton Women's Survey (SWS)

DNA methylation at two probes were significantly associated with total bone mineral content at 6 years and periosteal circumference at 6 years respectively. Probe cg26559250 located at Chr6:157,653,445-157,653,447 at the ZDHHC14 (zinc finger DHHC-type palmitoyltransferase 14) gene showed an increase of DNA methylation with increasing bone mineral content with a significance of  $2.52 \times 10^{-8}$ . Probe cg22570676 located at Chr19:2,527,492-2,527,494 at the GNG7 (G protein subunit gamma 7) gene showed an increase of DNA methylation with increasing periosteal circumference with a significance of  $4.24 \times 10^{-8}$ . Neither probe is flagged for known technical issues or genetic confounding [186].

#### 3.4.3.1 Whole Array QC

The predicted sex of the samples generated using sex chromosome probe intensities was checked against that in the sample annotation and no mismatches were found (Figure 3.40).

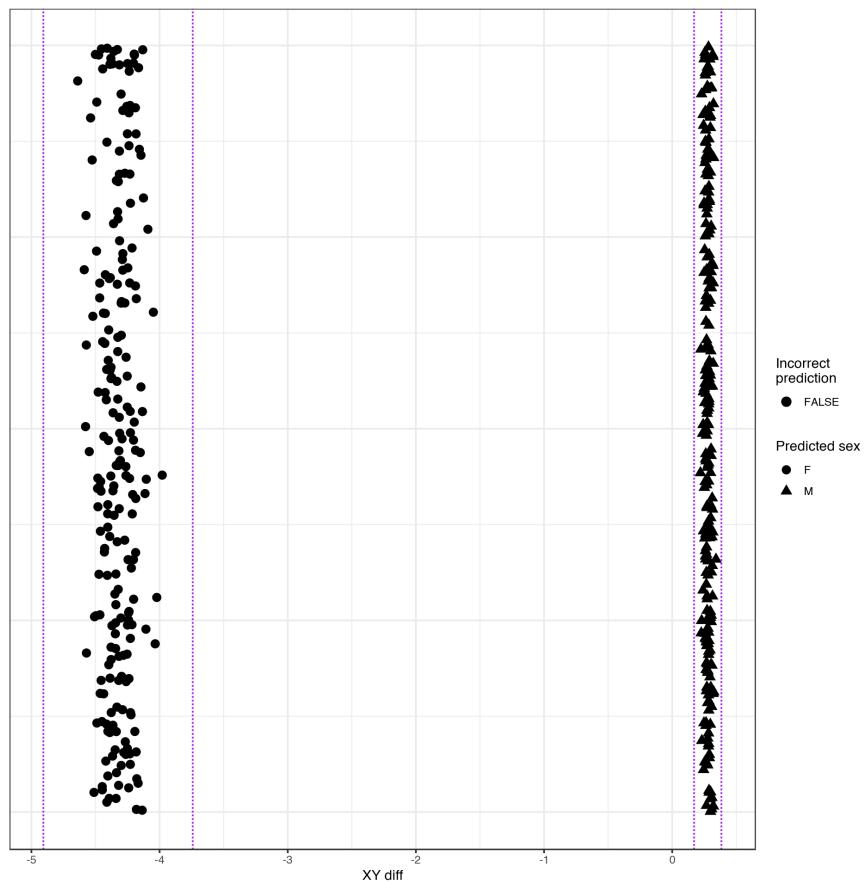


Figure 3.40: Predicted sex of each sample based on the sex chromosome copy numbers inferred from probe intensities for the EPIC array data. Mismatches between the predicted sex and that asserted in the sample annotation metadata are shown in red. Two predicted sex values differ from their annotations. Plot generated by `meffil` QC report.

No samples were excluded for having a median methylated signal that was more than  $3\sigma$  from the expected value, (Figure 3.41). No samples were excluded for having a higher than expected proportion of undetected probes (proportion of probes with detection p-value  $> 0.01$  is  $> 0.1$ ) (Figure 3.42). No samples were excluded for having a high proportion of probes with low bead counts (proportion of probes with bead number  $< 3$  is  $> 0.1$ ), (Figure 3.43).

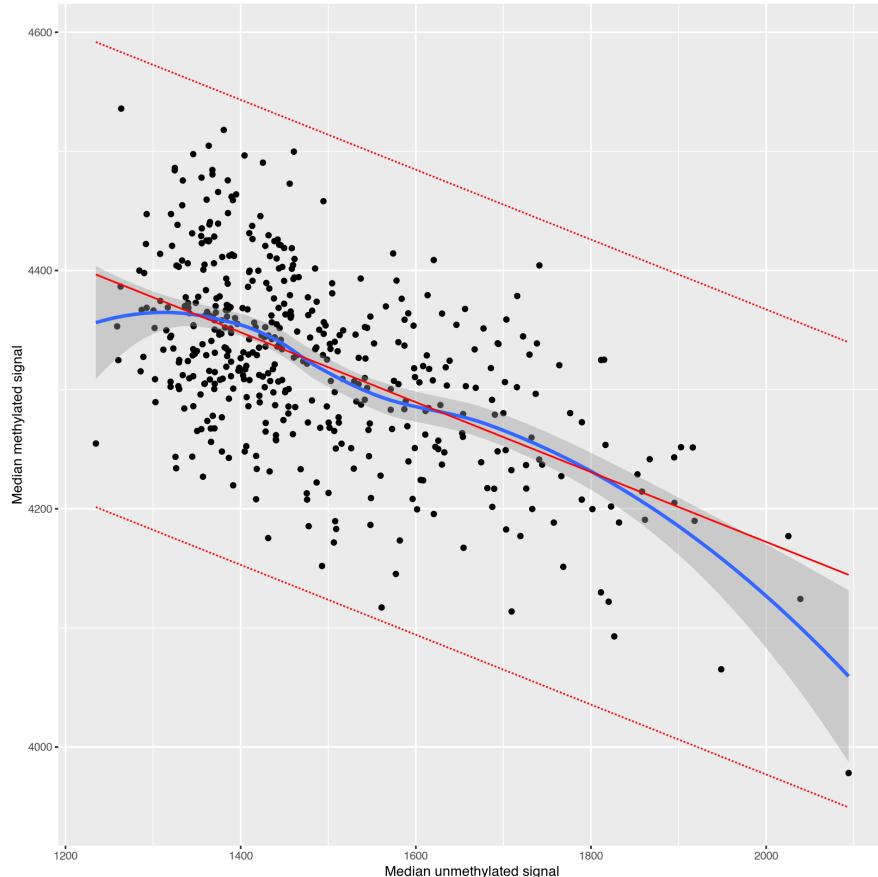


Figure 3.41: Median methylated signal vs unmethylated signal per sample for the EPIC array data, solid red line indicates linear regression of median methylated signal vs median unmethylated signal with dotted red lines representing  $3\sigma$  from the expected mean. Samples outside the expected range are indicated in the legend. Plot generated by `meffil` QC report.

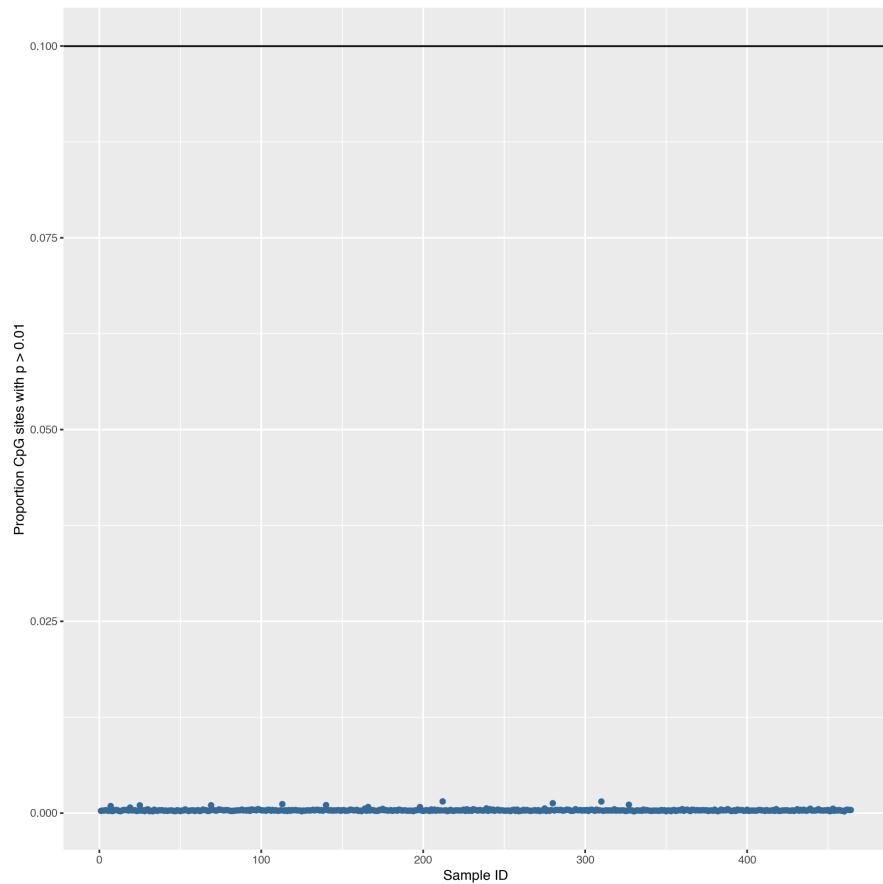


Figure 3.42: Proportion of probes with detection p-values  $>0.01$  by sample for the EPIC array data. Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil` QC report.

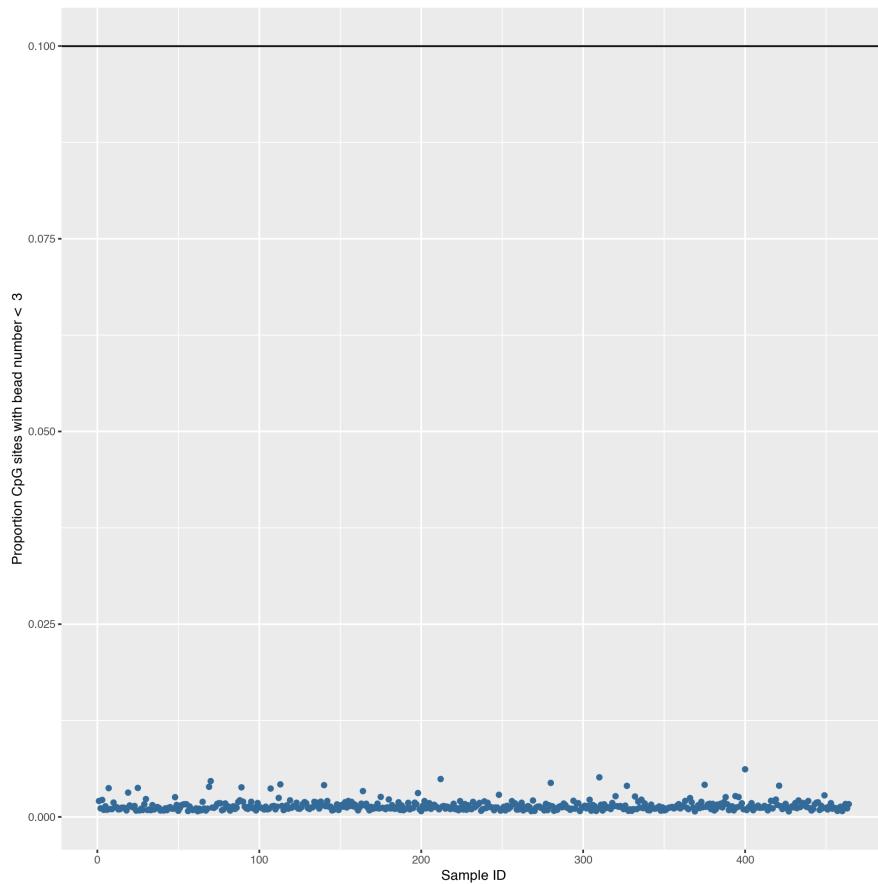


Figure 3.43: Proportion of probes with a bead count of  $< 3$  by sample for the EPIC array data. Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil` QC report.

#### 3.4.3.2 Probe QC

There were no outliers within the control probes (Figure 3.44). 833 probes were excluded for having high background signal in a high proportion of samples (proportion of samples with detection p-value  $> 0.01$  is  $> 0.1$ ), (Figure 3.45). 127 probes were excluded for having low bead count in a high proportion of samples (proportion of samples with bead number  $< 3$  is  $> 0.1$ ), (Figure 3.46). Probes with poor technical quality were excluded from the analysis prior to functional normalisation.

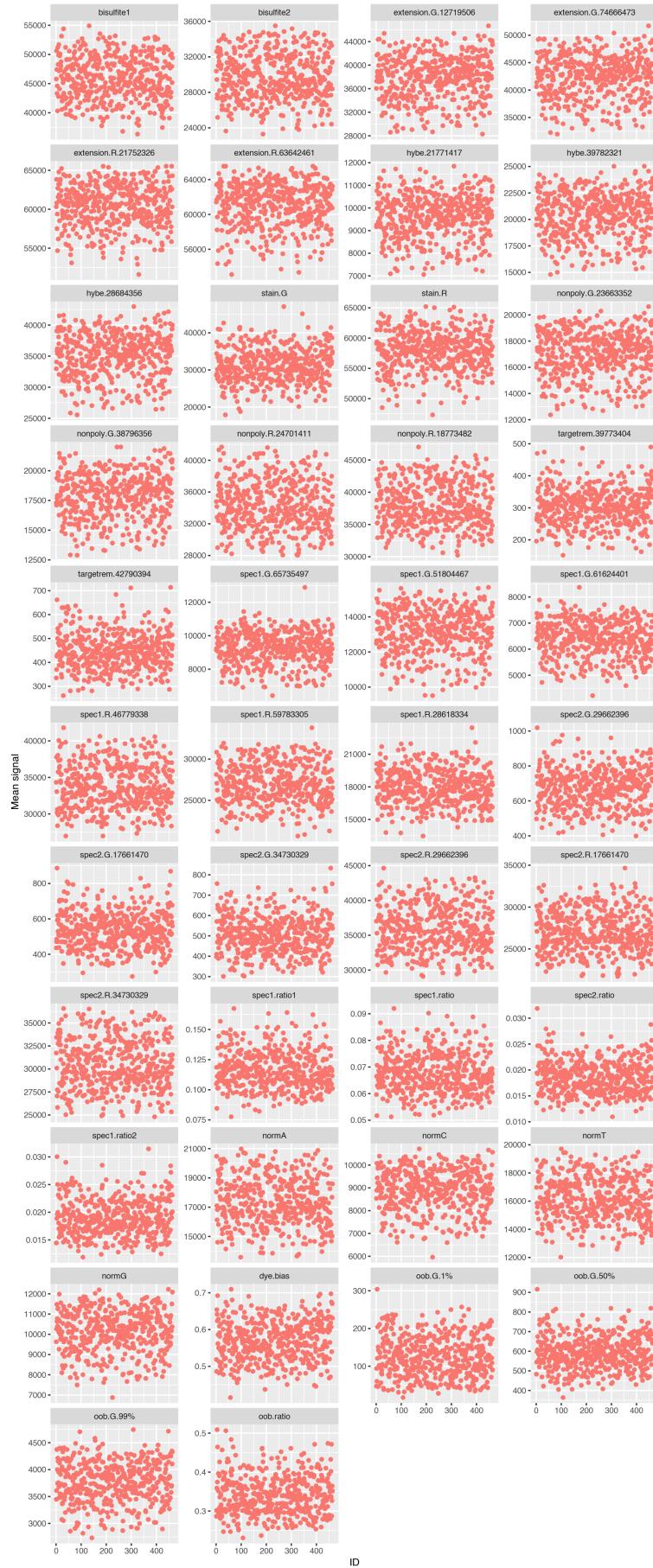


Figure 3.44: Control probe signal by sample for each summary group for the EPIC data. Outliers would be circled in black. Plot generated by `meffil` QC report.

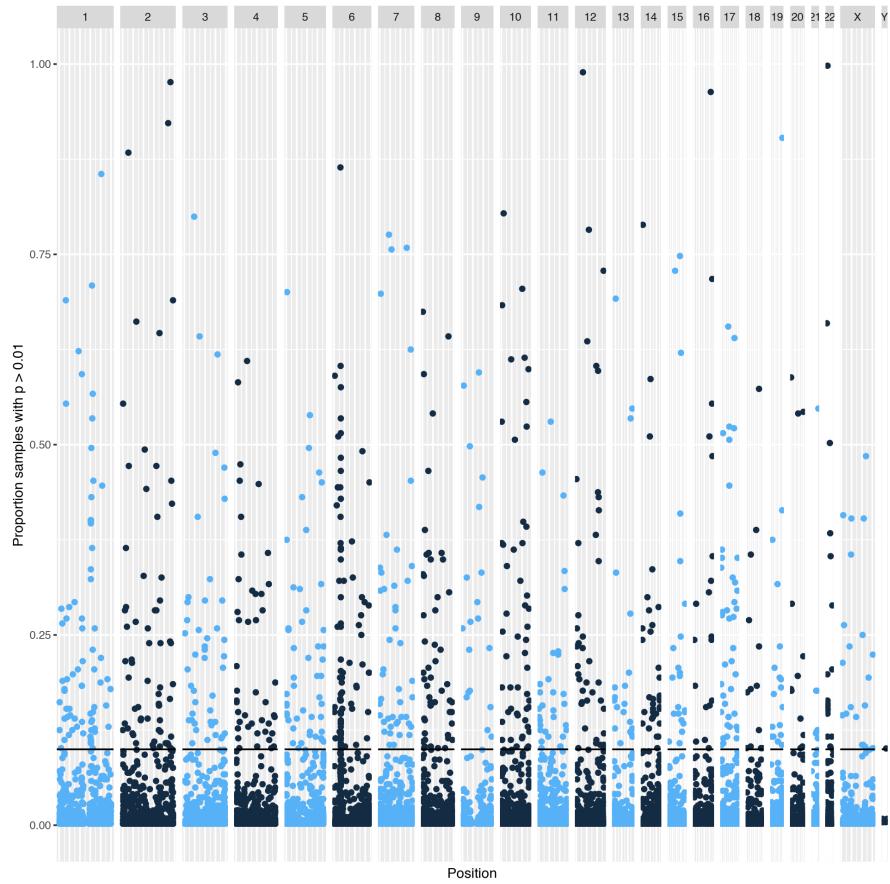


Figure 3.45: Undetectable probes across samples for EPIC data. Manhattan plot showing proportion of samples (y) in which a given probe (x) is not distinguishable from background noise, i.e. a detection p-value of  $> 0.01$ . Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil` QC report.

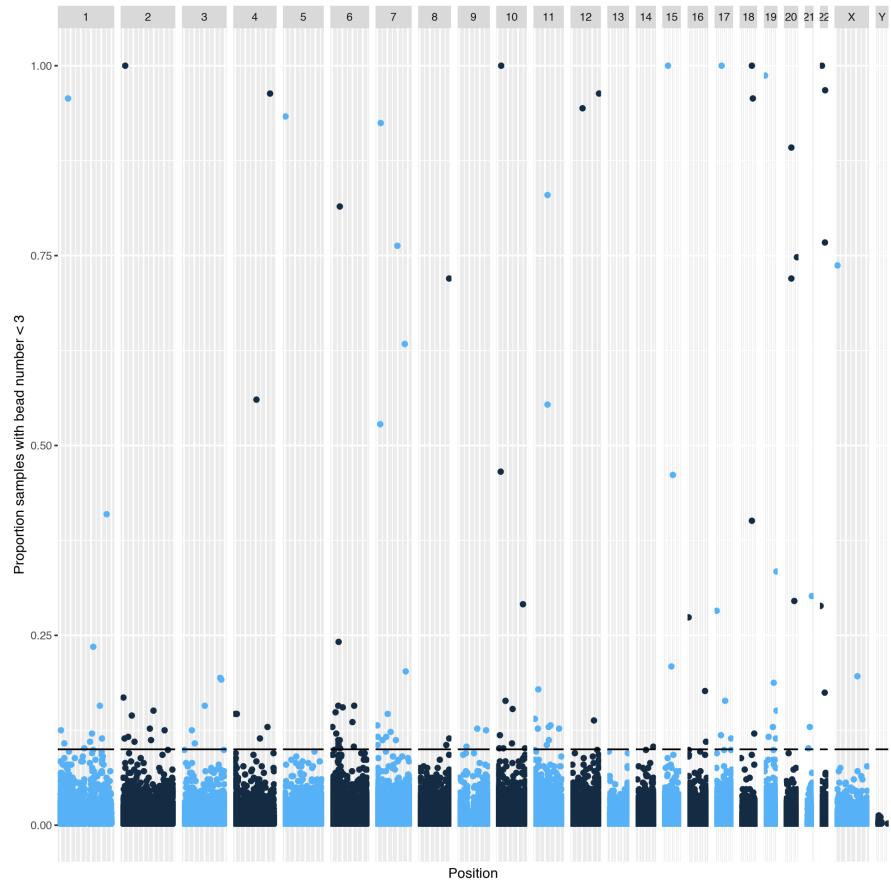


Figure 3.46: Low bead count probes across samples for EPIC data. Manhattan plot showing the proportion of samples (y) in which a given probe (x) has a bead count of  $<3$ . Black line indicates the exclusion threshold of 0.1. Plot generated by `meffil QC report`.

#### 3.4.3.3 Functional Normalisation

An  $m$  of 12 was chosen as this value produced the last steep drop in residual variation, see Figure 3.47.

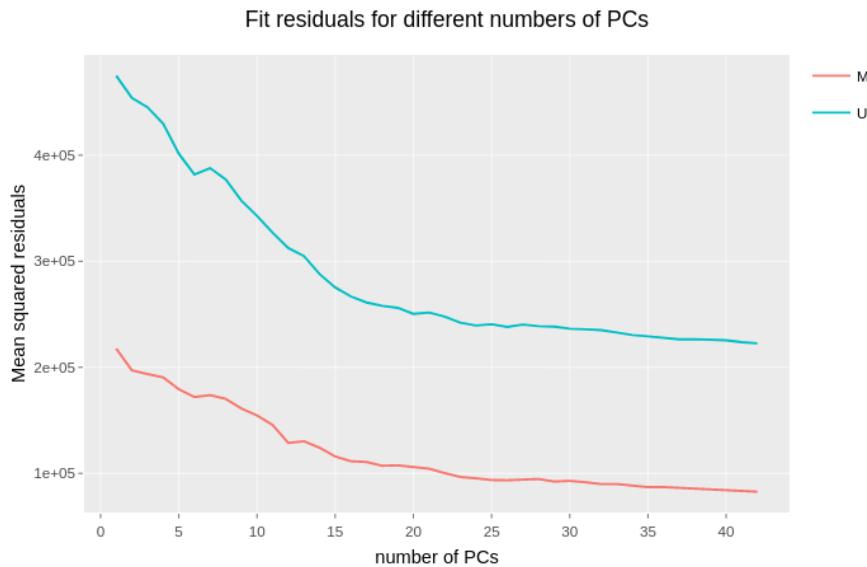


Figure 3.47: Residual variation remaining after functional normalisation of the top 20,000 most variable probes with  $m$  PCs from the control probe summary matrices for the EPIC array samples ( $n=464$ ), for M = methylated and U = unmethylated probes.

#### 3.4.3.4 EWAS

For all the EWAS, blood cell-type counts were estimated using the Houseman method [165] and the cord blood cell-type reference panel from Bakulski et al. [235]. The Cell-types estimated were: B cells, CD4+ T cells, CD8+ T Cells, Granulocytes, Monocytes, Natural Killer cells, & Erythrocytes. In addition to the estimated blood cell counts all models included as covariates: Maternal Age at time of birth (years), Sex, maternal BMI at 11 weeks gestation, parity, whether or not the mother smoked during pregnancy, and gestational age.

**3.4.3.4.1 Total Bone Mineral Content at 8 years** Figure 3.48, illustrates the distribution of bone mineral content at 8 years in the 408 individuals in this EWAS. Surrogate variable analysis identified 95 significant surrogate variables, this is likely an overestimate stemming from small amounts of variation remaining unaccounted for by the manual model thus the Manhattan plots based on the manual model were included. No probes fell below the Bonferroni corrected significance threshold ( $5.92 \times 10^{-8}$ ) for an association between DNA methylation at that locus and total bone mineral content minus head at 8 years adjusted for age and sex, Figure 3.49.

8 yr DXA: Total BMC (kg), without heads,adjusted for sex and age

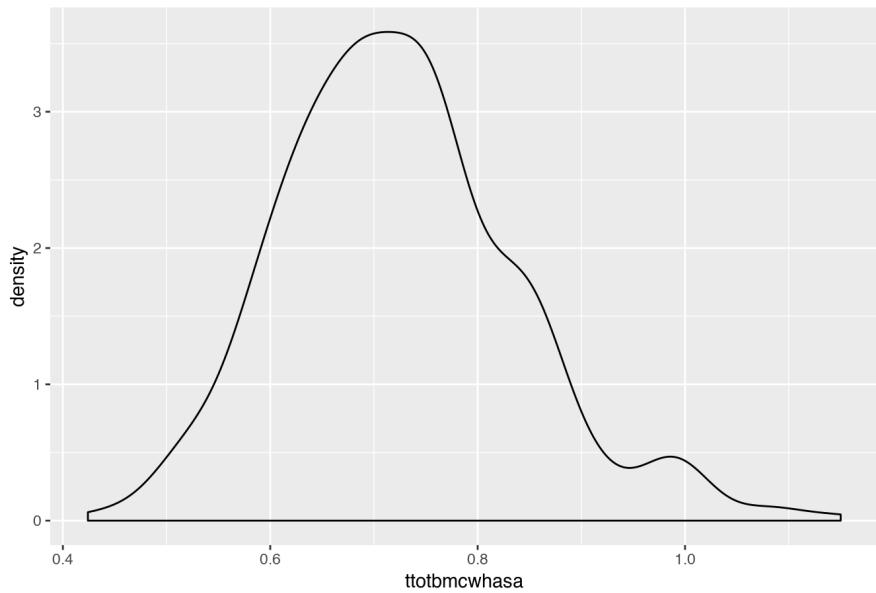


Figure 3.48: Distribution of whole body (minus head) bone mineral content in kg at 8 years of age ( $n = 408$ ), adjusted for sex and age, as measured by DXA.

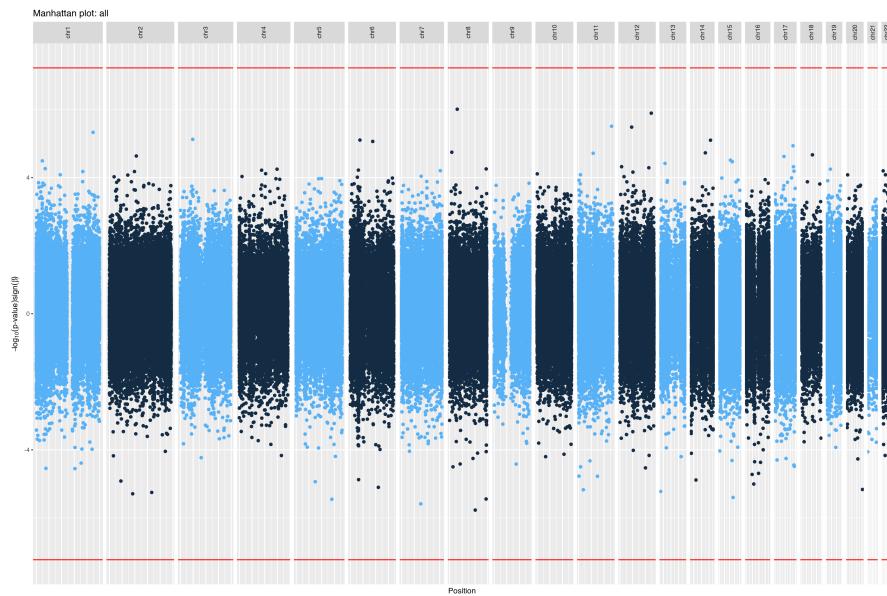


Figure 3.49: Results of EWAS for whole body (minus head) bone mineral content in kg at 8 years of age ( $n = 408$ ), adjusted for sex and age. The 'all' model results are shown here. Bidirectional Manhattan plot on which  $-\log_{10}(p\text{-value})$  is plotted on the y axis and the sign of this value represents the direction of change. The x axis represents chromosomes and position thereupon. Red line indicates the significance threshold of  $5.92 \times 10^{-8}$ .

**3.4.3.4.2 Total Bone Mineral Content at 6 years** Figure 3.50, illustrates the distribution of bone mineral content at 6 years in the 402 individuals in this EWAS. Surrogate variable

analysis identified 97 significant surrogate variables, this is likely an overestimate stemming from small amounts of variation remaining unaccounted for by the manual model thus Manhattan plots based on the manual model have been included.

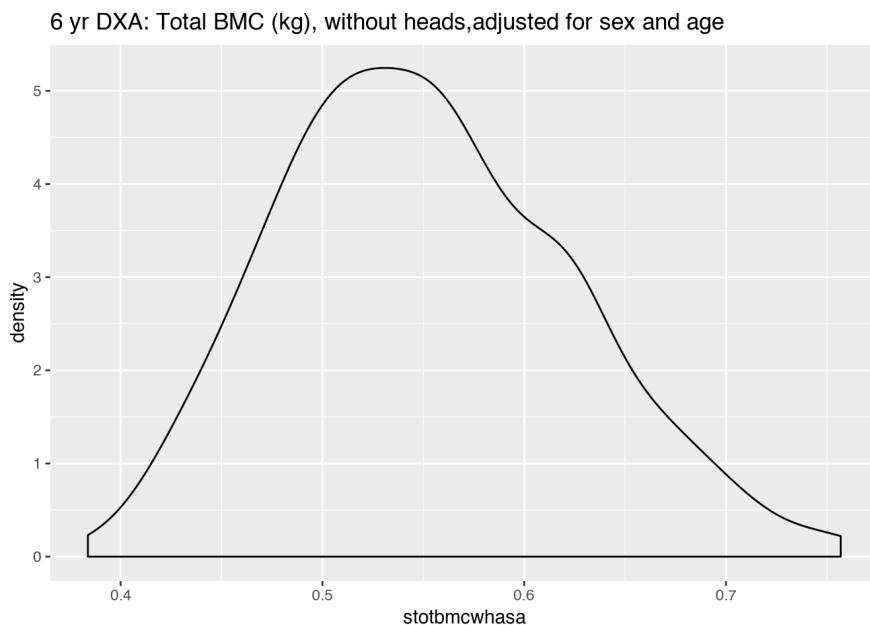


Figure 3.50: Distribution of whole body (minus head) bone mineral content in kg at 8 years of age ( $n = 402$ ), adjusted for sex and age, as measured by DXA.

One probe fell below the Bonferroni corrected significance threshold ( $5.92 \times 10^{-8}$ ) for an association between DNA methylation at that locus and total bone mineral content minus head at 6 years adjusted for age and sex, Figure 3.51. This probe was cg26559250 which is located at Chr6:157,653,445-157,653,447 adjacent to the ZDHHC14 (zinc finger DHHC-type palmitoyltransferase 14) gene. cg26559250 was significant ( $p = 2.52 \times 10^{-8}$ , increase of 1.46% per kg) in the 'all' model and was also Bonferroni significant in the uncorrected model. However, cg26559250 was not significant in the SVA or iSVA models suggesting that it may be attributable to batch or cell-type effects.

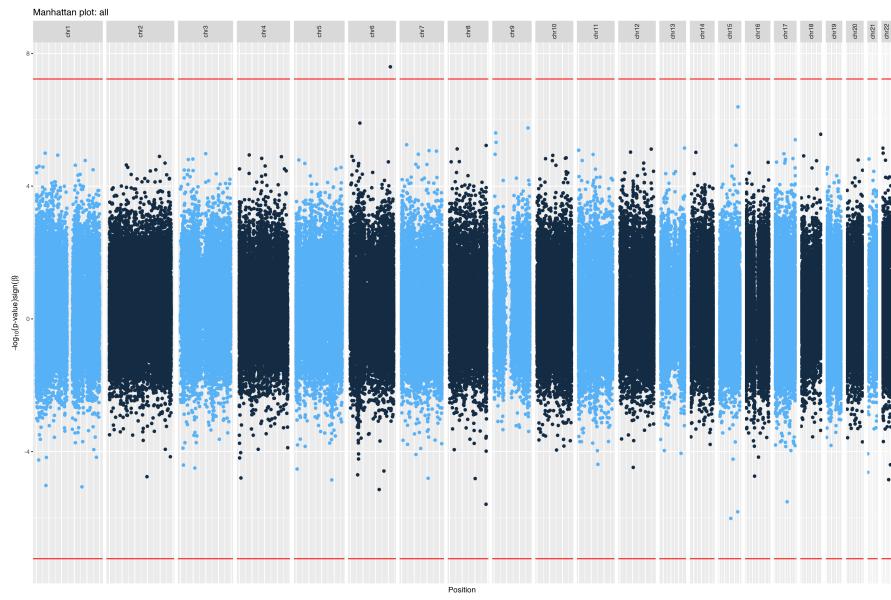


Figure 3.51: Results of EWAS for whole body (minus head) bone mineral content in kg at 6 years of age ( $n = 402$ ), adjusted for sex and age. The ‘all’ model results are shown here. Bidirectional Manhattan plot on which  $-\log_{10}(p\text{-value})$  is plotted on the y axis and the sign of this value represents the direction of change. The x axis represents chromosomes and position thereupon. Red line indicates the significance threshold of  $5.92 \times 10^{-8}$ .

**3.4.3.4.3 Periosteal Circumference at 6 years** Figure 3.52, illustrates the distribution of periosteal circumference at 38% from the distal end of the tibia at 6 years (mm) in the 141 individuals in this EWAS. Surrogate variable analysis identified 37 significant surrogate variables, this is likely an overestimate stemming from small amounts of variation remaining unaccounted for by the manual model thus Manhattan plots based on the manual model have been included.

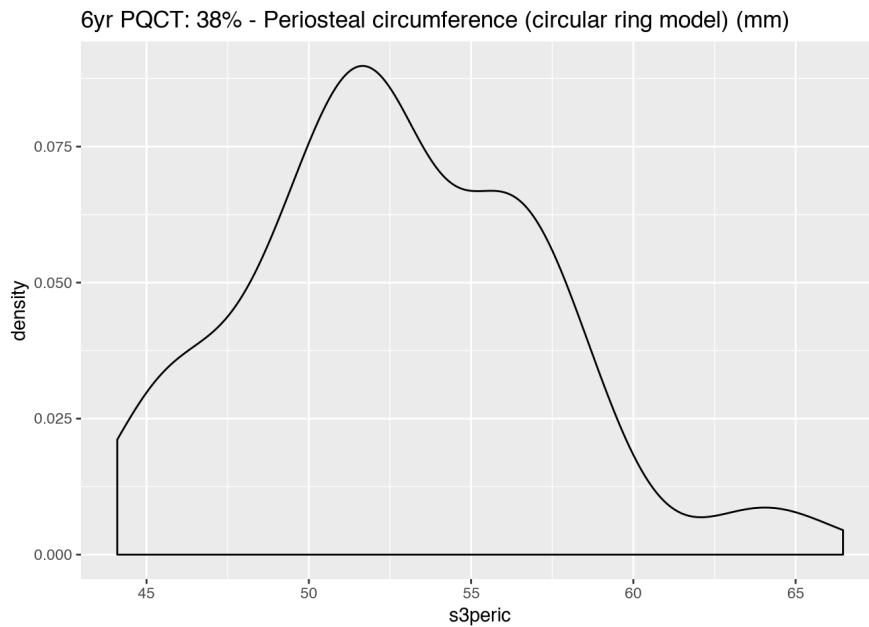


Figure 3.52: Distribution of periosteal circumference at 38% from the distal end of the tibia (mm) at 6 years of age ( $n = 141$ ), adjusted for sex and age, as measured by PQCT.

One probe fell below the Bonferroni corrected significance threshold for an association between DNA methylation at that locus and periosteal circumference at 38% from the distal end of the tibia at 6 years (mm) adjusted for age and sex, Figure 3.53. This probe was cg22570676 which is located at Chr19:2,527,492-2,527,494 at the GNG7 (G protein subunit gamma 7) gene. cg22570676 was significant ( $p = 4.24 \times 10^{-8}$ , increase of 0.370% per mm) in the ‘all’ model and was also Bonferroni significant in the uncorrected model. However, cg22570676 was not significant in the SVA or iSVA models suggesting that it may be attributable to batch or cell-type effects.

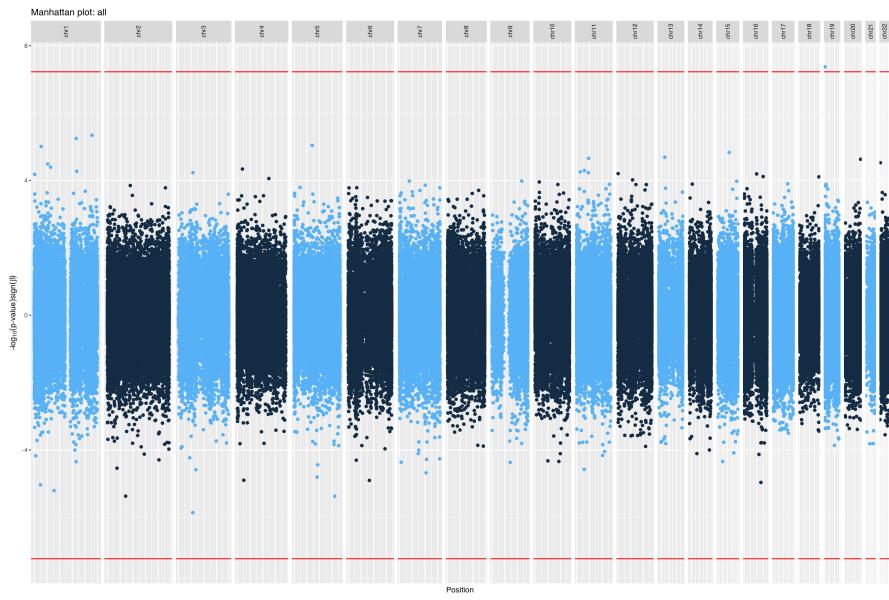


Figure 3.53: Results of EWAS for periosteal circumference at 38% from the distal end of the tibia (mm) at 6 years of age ( $n = 141$ ), adjusted for sex and age. The ‘all’ model results are shown here. Bidirectional Manhattan plot on which  $-\log_{10}(p\text{-value})$  is plotted on the y axis and the sign of this value represents the direction of change. The x axis represents chromosomes and position thereupon. Red line indicates the significance threshold of  $5.92 \times 10^{-8}$

**3.4.3.4.4 Cortical Denisty at 6 years** Figure 3.54, illustrates the distribution of cortical density at 38% from the distal end of the tibia at 6 years ( $mg\ cm^{-3}$ ) in the 141 individuals in this EWAS. No probes fell below the Bonferroni corrected significance threshold ( $5.92 \times 10^{-8}$ ) for an association between DNA methylation at that locus and cortical density at 6 years, Figure 3.55. Surrogate variable analysis identified 37 significant surrogate variables, this is likely an overestimate stemming from small amounts of variation remaining unaccounted for by the manual model thus Manhattan plots based on the manual model.

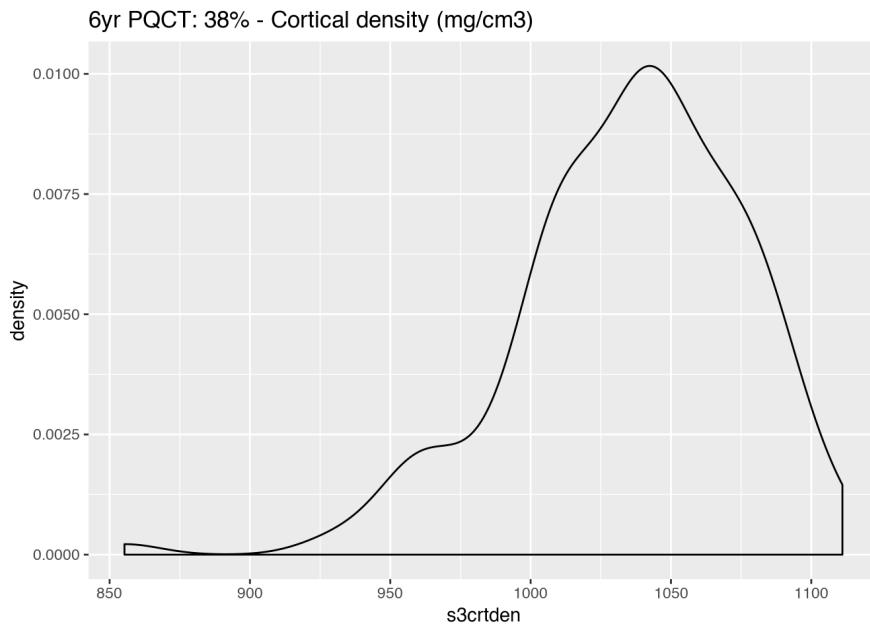


Figure 3.54: Distribution of cortical density at 38% from the distal end of the tibia ( $mg\ cm^{-3}$ ) at 6 years of age (n = 141), as measured by PQCT.

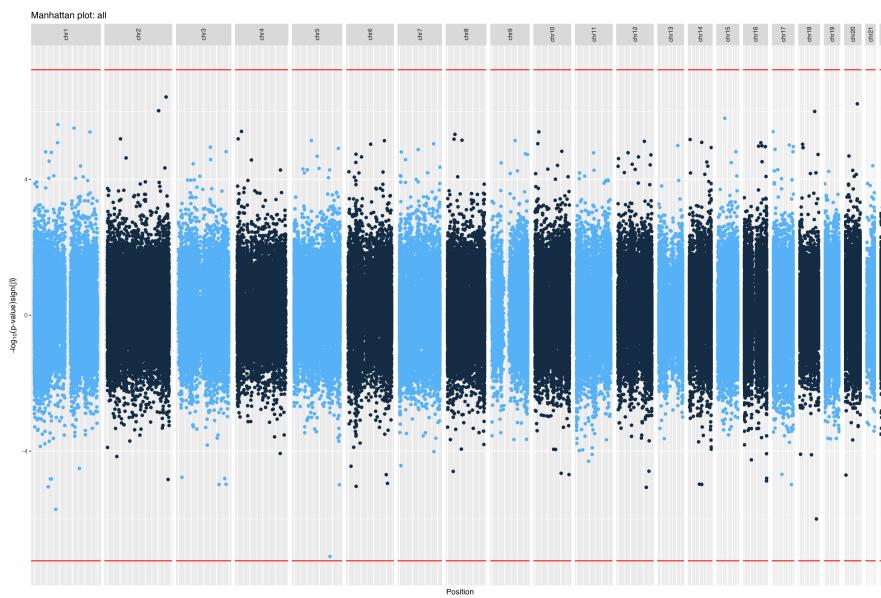


Figure 3.55: Results of EWAS for cortical density at 38% from the distal end of the tibia ( $mg\ cm^{-3}$ ) at 6 years of age (n = 141). The 'all' model results are shown here. Bidirectional Manhattan plot on which  $-\log_{10}(p - value)$  is plotted on the y axis and the sign of this value represents the direction of change. The x axis represents chromosomes and position thereupon. Red line indicates the significance threshold of  $5.92 \times 10^{-8}$

### 3.5 Discussion

EWAS for 9 outcomes were carried out across 3 sets of samples from MAVIDOS and SWS using the EPIC and 450K array platforms. No significant results were found in either the first or second phase of the MAVIDOS analysis. Two possible results for bone outcomes at 6 years were identified in the SWS data but did not remain in models including surrogate variables for possible confounding effects. The effect observed at probe cg26559250, adjacent to the ZDHHC14 gene, in EWAS for bone mineral content at 6 years was not seen at the 8 year time point. Furthermore both this result and the finding for probe cg22570676 and periosteal circumference were of small effect sizes 1.46% per kg, 0.370% per mm respectively. Consequently these results should be treated with considerable caution, they are in need of replication before they can be considered reliable. No Bonferroni significant DNA methylation changes at the *CDKN2A* and *RXRA* loci which have been previously associated with maternal vitamin D and bone phenotypes were identified, despite the presence of 95 and 75 probes annotated as being in the vicinity of these genes respectively.

Whilst the calculation of power for EWAS is a complex and rather understudied problem [242–244] it is possible to achieve some approximations using Cohen’s methods [245]. To achieve the modest goal of 80% power for a small effect size ( $r^2 = 0.02$ ) in a linear regression analysis (F-test) with 7 covariates at a significance level suitable for the EPIC array of  $p = 5.92 \times 10^{-8}$  an n of 2607 is needed. Seven was the smallest number of covariates used in an SVA model in these analyses, models in the SWS analyses had 13 manually specified variables. When considering what is for EWAS a very large effect size ( $r^2 = 0.15$ ) with the 13 covariates used in the SWS models it is possible to achieve 80% power for an n of 374. This is an effect size in line with the effects of smoking on DNA methylation at some loci [246,247]. This would make the two largest EWAS performed here for BMC at 8 and 6 years (n = 408, n = 402 respectively) powered only to find large effect sizes with just over 80% probability. The most generous set of parameters (80% power,  $r^2 = 0.15$ ,  $p = 5.92 \times 10^{-8}$ , & 2 covariates) yield an n of 259, more realistic numbers (90% power,  $r^2 = 0.02$ ,  $p = 5.92 \times 10^{-8}$ , & 13 covariates) yield an n of 3370. EWAS have identified biologically relevant changes associated with environmental exposures in DNA methylation with magnitudes of less than a single percentage point, and percentage changes in the low single digits are not uncommon in EWAS [247]. This makes all of the EWAS performed here underpowered to identify small DNA methylation changes which might be expected to occur. A collaboration with colleagues in Bristol is underway to perform a meta-analysis to include these data with similar results from other cohorts in order to increase the power of these analyses. The covariates included in the models for the second phase of MAVIDOS analysis and the SWS analysis are matched to those being used by our collaborators to maximise the comparability of our results.

Given that the effect of maternal vitamin D on neonatal bone mass appears to be seasonal, with only babies born in the winter months showing statistically significant benefits of supplementa-

tion [172], it would be interesting to perform seasonally stratified EWAS were sufficient numbers available to do so with reasonable power.

Attempting to identify small changes in the overall DNA methylation state of complex populations of cells like blood and umbilical cord tissue which are associated with phenotypes such as circulating maternal vitamin D is a technically challenging undertaking. This work provides two candidates for further analysis for associations between DNA methylation and bone health outcomes at 6 years of age. Furthermore, these results have contributed to a larger meta analysis with greater power to detect associations with metrics of bone health.

## Chapter 4

# The Genomic Loci of Specific Human tRNA Genes Exhibit Ageing-Related DNA Hypermethylation

### 4.1 Abstract

Understanding how the epigenome deteriorates with age and subsequently impacts on biological function may bring unique insights to ageing-related disease mechanisms. As a central cellular apparatus, tRNAs are fundamental to the information flow from DNA to proteins. Whilst only being transcribed from ~46kb (<0.002%) of the human genome, their transcripts are the second most abundant in the cell. Furthermore, it is now increasingly recognised that tRNAs and their fragments also have complex regulatory functions. In both their core translational and additional regulatory roles, tRNAs are intimately involved in the control of metabolic processes known to affect ageing. Experimentally DNA methylation can alter tRNA expression, but little is known about the genomic DNA methylation state of tRNAs.

Here, we find that the human genomic tRNA loci are enriched for ageing-related DNA hypermethylation. We initially identified DNA hypermethylation of 44 and 21 specific tRNA genes, at study-wide ( $p < 8.36 \times 10^{-5}$ ) and genome-wide ( $p < 4.34 \times 10^{-9}$ ) significance, respectively, in 4,350 MeDIP-seq peripheral blood DNA methylomes (16 - 82 years). This starkly contrasted with 0 hypomethylated at both these significance levels. Further analysing the 21 genome-wide results, we found 3 of these tRNAs to be independent of major changes in cell-type composition

(tRNA-iMet-CAT-1-4, tRNA-Ser-AGA-2-6, tRNA-Ile-AAT-4-1). We also excluded the ageing-related changes being due to the inherent CpG density of the tRNAome by permutation analysis (1,000x, Empirical p-value  $< 1 \times 10^{-3}$ ). We additionally explored 79 tRNA loci in an independent cohort using Fluidigm deep targeted bisulfite-sequencing of pooled DNA (n=190) across a range of 4 timepoints (aged ~4, ~28, ~63, ~78 years). This revealed these ageing changes to be specific to particular isodecoder copies of these tRNAs (tRNAs coding for the same amino acid but with sequence body differences) and included replication of 2 of the 3 genome-wide tRNAs. Additionally, this isodecoder-specificity may indicate the potential for regulatory fragment changes with age.

In this study we provide the first comprehensive evaluation at the genomic DNA methylation state of the human tRNAome, revealing a discreet and strongly directional hypermethylation with advancing age.

## 4.2 Introduction

Ageing is implicated as a risk factor in multiple chronic diseases [248]. Understanding how the ageing process leads to deteriorating biological function is now a major research focus. This field has hopes of increasing human longevity and ‘healthspan’ whilst ameliorating the extensive physical, social and economic costs of these ageing-related disorders [249]. Epigenetic processes, which influence or inform cell-type specific gene expression, are altered with age and are a fundamental hallmark of this progression, indeed they are arguably a hub mediating other hallmarks including stem cell exhaustion, cellular senescence, and mitochondrial dysfunction [11,13].

DNA methylation (DNAm) is the most common epigenetic modification of DNA and age-associated changes in this mark were recognised in mammalian tissues as early as 1983 [109]. Changes in DNAm with age are extensive with thousands of loci affected. Many of these changes represent ‘drift’ [119] arising from the imperfect maintenance of methylation state. Specific genomic regions show distinct directional changes, with loss of DNA methylation in repetitive or transposable elements [250], and gains in certain promoters, including the targets of polycomb repressor complex [127] as well as bivalent domains [125]. The advent of high-throughput DNAm arrays [127,137,251] has elucidated more detailed patterns of ageing related changes in DNAm. The identification of precise individual CpG sites that exhibit consistent changes with age enabled the construction of predictors of chronological age known as epigenetic or DNAm ‘clocks’ [137,138,157,158,252]. Furthermore, it was observed that ‘acceleration’ of this DNAm-derived measure is a biomarker of ‘biological’ ageing due to associations with morbidity and mortality [168,253]. In a previous investigation of ageing-related DNAm changes within common disease-associated GWAS regions, Bell et al. identified hypermethylation of the specific transfer RNA gene, tRNA-iMet-CAT-1-4 [170]. The initiator methionine tRNA possesses certain unique properties [254–256], including its capacity to be rate limiting for translation [257], association with the translation initiation factor eIF2 [258], and ability to impact the expression of other tRNA genes [259].

tRNAs are fundamental in the translation process for all domains of life and are thus evolutionarily ancient [260]. This translation machinery and the regulation of protein synthesis are controlled by conserved signalling pathways shown to be modifiable in longevity and ageing interventions [261]. The human tRNAome, comprising all of the genomic locations of tRNA genes, represents an extremely small portion of the genome [262]. There are 416 high confidence tRNA genes when additionally considering tRNA pseudo genes, nuclear encoded mitochondrial tRNA genes and possibly some closely related repetitive sequences the number of sequences closely resembling tRNAs is 610 (gtRNADB [263]) this extended set covers <46 kb (including introns) and represents <0.002% of the human genome. Despite their small genomic footprint, and the observation that approximately half of all tRNA genes are transcribed at negligible levels if at all [264], these genes produce the second most abundant RNA species next to ribosomal RNA

[265] and are required for the production of all proteins. Mature tRNAs have an L-shaped three dimensional structure arising from a ‘clover-leaf’ shaped two dimensional structure comprised of three hairpin stem-loop structures (Figure 4.1).

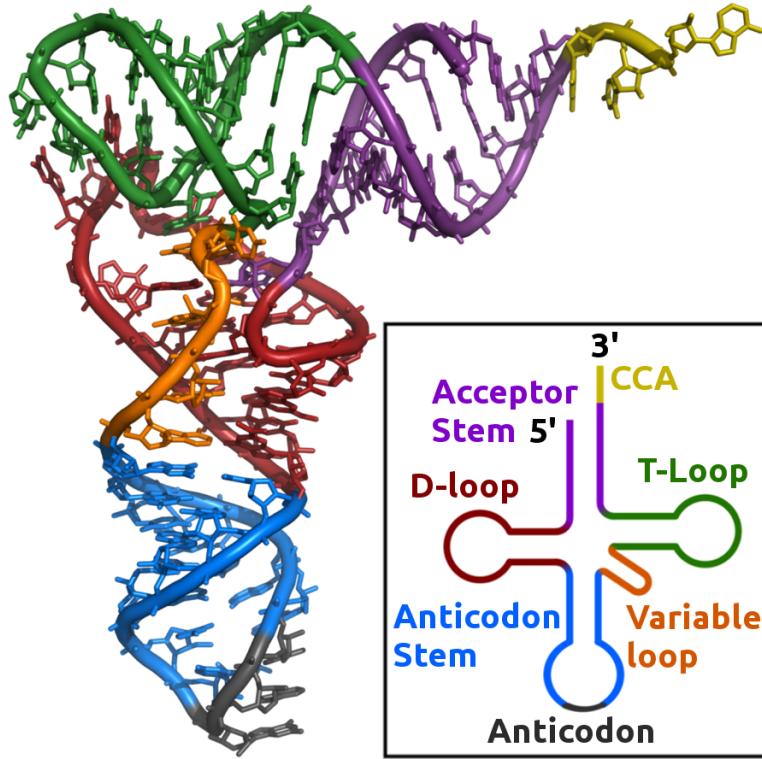


Figure 4.1: **Structure of a mature tRNA** Two and three dimensional representations of tRNA structure with matching colour coding. Adapted from the wikimedia foundation structure based on PDBID: 1ehz. tRNAs are ‘charged’ when an amino acid is attached at the CCA site at the 3’ end.

tRNA genes are transcribed by RNA polymerase III (polIII) [266] and have type II polIII promoters which contain A and B-box internal promoter elements bound by the complex TFIIIC, followed by TFIIIB, and polIII [267] (Figure 4.2). Transcription is terminated by a simple run of Ts and proceeds in rounds of fast re-initiation where the same polIII molecule is preferentially re-used [268]. tRNA gene expression is modulated by the polIII specific transcription factor Maf1 a highly conserved factor which represses tRNA transcription [269–271]. The activity of Maf1 is modulated by the Target of Rapamycin Kinase Complex 1 (TORC1) [272], a highly conserved hub for signals that modulate ageing [273]. Several general transcription factors also influence tRNA gene expression; the tumour suppressors p53 [274] and Rb [275] both negatively regulate tRNA expression and c-Myc upregulates tRNA gene expression, all act via TFIIIB [276]. tRNAs are dysregulated in cancer, and have potential utility as prognostic markers [277]. tRNA regulation may play an important role in cancer [278]. DNAm is able to repress the expression of tRNA genes experimentally [279], in plasmid expression systems, but may also represent co-ordination

with the local repressive chromatin state [280]. In addition, recent results from Gerber et al. [281] show a mechanism by which polII can function in the regulation of certain polIII transcribed loci including several tRNA genes.

## tRNA Transcription

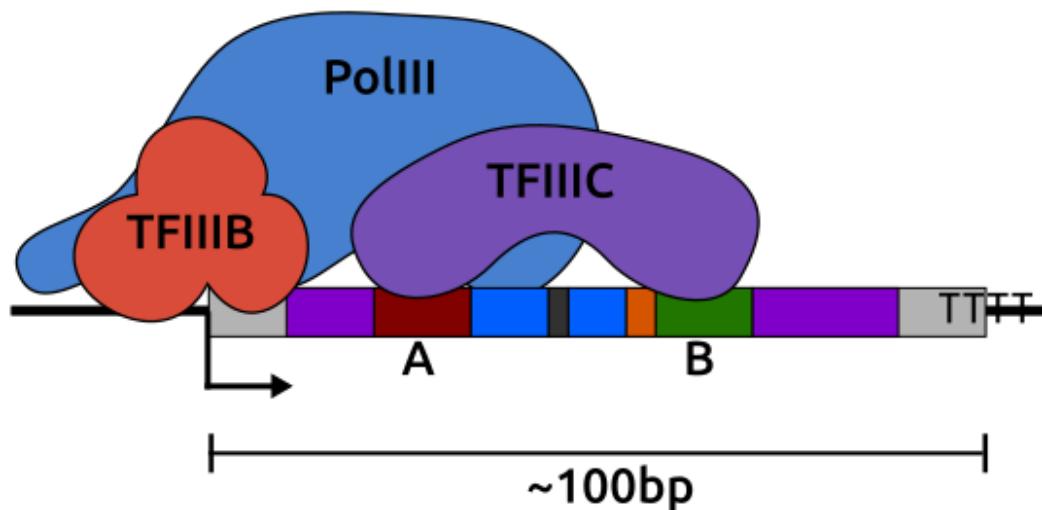


Figure 4.2: **tRNA Transcription** Cartoon representation of the RNA polIII transcription initiation complex and structure of the type II RNA polymerase III promoter. Colour coding here corresponds to that in figure 4.1 illustrating that promoter is internal as the A box corresponds approximately to the D-loop and the B box to the T-loop in the tRNA structure.

Assuming similar rates of transcription one would expect that the more frequently an amino acid is used in the exome, the more copies of that tRNA gene there would be in the genome [282]. Indeed, tRNA gene dosage is quite closely matched to amino acid usage frequency in the human exome, though the correlation is less strong for codon usage (Figure 4.3). The imperfect nature of this correlation suggests that there may be regulation of tRNA expression beyond simply having copy numbers proportionate to usage frequency.

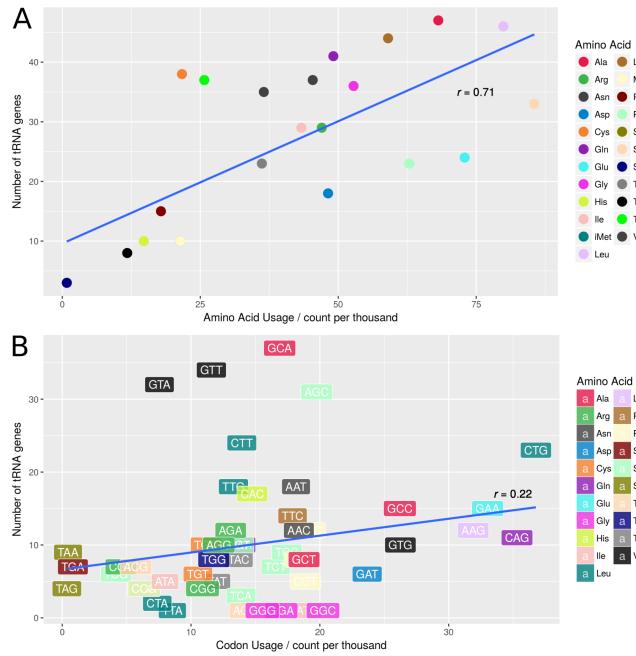


Figure 4.3: tRNA gene copy number is imperfectly correlated with amino acid and codon usage frequency **A**) Amino acid usage frequency in the human exome vs tRNA gene copy number **B**) Codon usage frequency in the human exome vs tRNA gene copy number Usage Frequency Data [283], tRNA gene count data from GtRNADB [263].

A variety of observations have indicated that many tRNA genes are expressed in a tissue specific fashion in diverse organisms [284],[285]. These include the tissue specific manifestations of diseases caused by mutations involving tRNAs and their usage [286],[287]. Though detailed characterisation of tissue specific tRNA expression patterns are still lacking. Considering the broad similarity of transcriptome codon usage across time and tissues and the purported functional equivalence in translation of tRNAs with the same anticodon, there is an open question as to the reason for tissue specificity of tRNA expression. The tissue specific expression of tRNAs whilst substantially maintaining codon ratios would require highly coordinated tRNA gene regulation to obtain similar expression levels of putatively functionally equivalent genes. This situation is strongly suggestive of as yet undiscovered function.

Additionally, beyond their core role in the information flow from DNA to protein sequence, tRNAs can fragment into numerous tRNA-derived small RNAs (tsRNAs). There are a number of types of tRNA derived small RNAs (tsRNAs) [288] which serve as signalling molecules [289] (Figure 4.4). Four main types of tsRNAs have been identified, these are 5' tRNA halves, 3' tRNA halves, and the tRNA derived fragments (tRFs) of ~18 - ~22bp sizes [290–292]. tRNAs and pre-tRNAs also give rise to piRNAs [293–295] and internal tRNA fragments or itRFs [296]. The terminology and ontology of tsRNAs is still stabilising as our understanding of these RNA species evolves. Some of these tsRNAs feedback on protein synthesis by regulating ribosome

biogenesis [297], others have diverse regulatory functions such as targeting transposable element transcripts [298]. The extent of the functional significance of tsRNAs is also an open question [299].

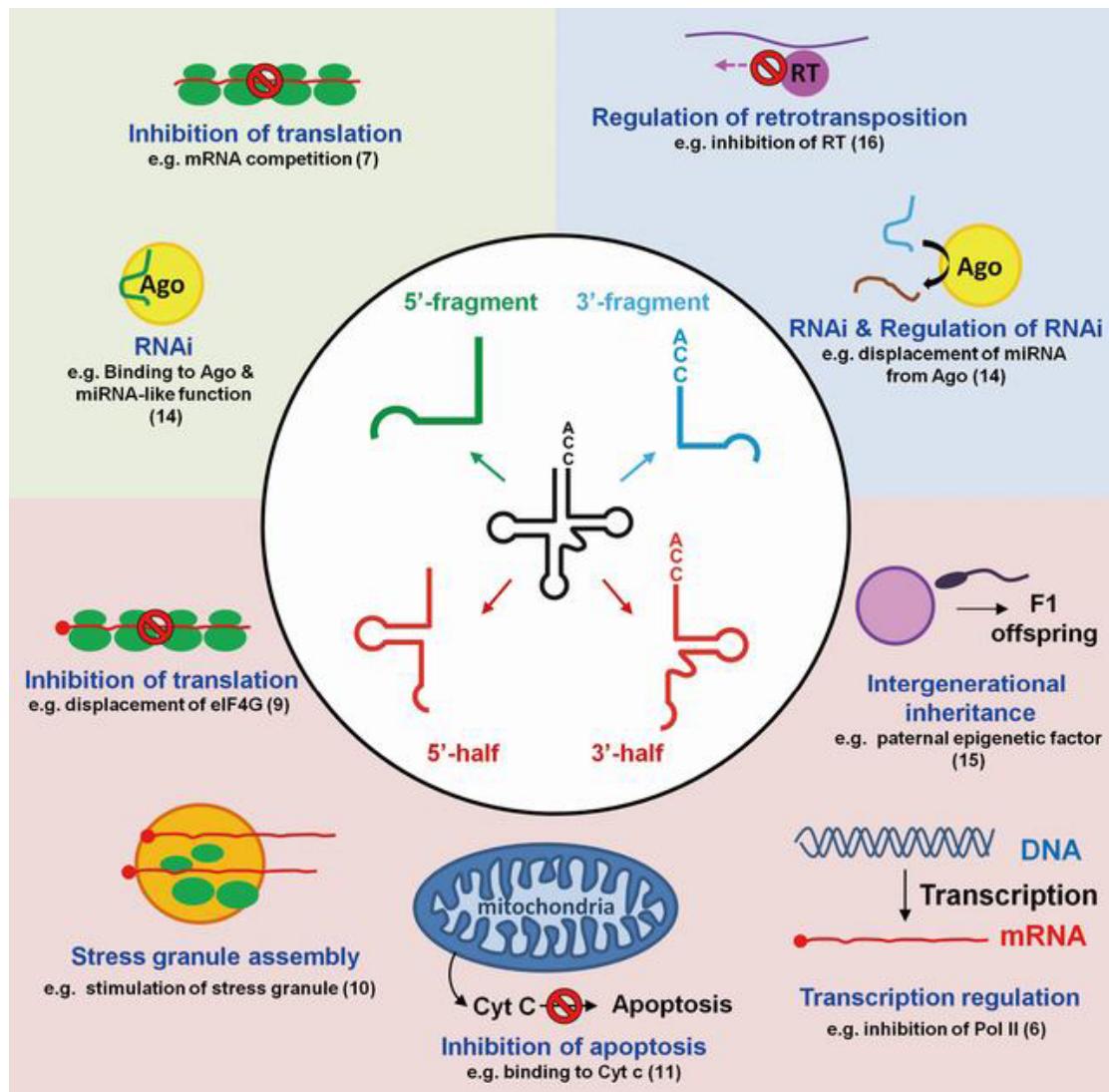


Figure 4.4: **The Types and Functions of tRNA derived small RNAs** Reproduced from Cristodero et al. [300].

Ageing is linked to core aspects of metabolic regulation, with nutrient sensing and stress response acting as major modulators of ageing. tRNAs as well as tsRNAs are integral to the regulation of protein synthesis and stress response. Protein synthesis represents a substantial proportion of total cellular energy expenditure and this fraction can vary considerably with nutrient availability [301]. Metabolic processes are also recognised to modulate the age estimates of DNAm clocks [302]. Partial inhibition of translation increases lifespan in multiple model organisms [303] and polIII inhibition increases longevity acting downstream of TORC1 [304]. Furthermore, 5' tRNA halves circulating in serum are modulated by ageing and caloric restriction [305].

DNA is able to repress the expression of tRNA genes on methylated expression vectors [279]. Whilst the broader chromatin milieu affects tRNA transcription, tRNA genes are unusual in that they are sufficiently short to fit within a single nucleosome. tRNA genes are generally ‘nucleosome free’ and precise placement of the nucleosome immediately upstream of the transcriptional start site is of importance for their expression [306]. A recent review of the epigenetic regulation of the polIII transcriptome [307] noted the very limited CpG methylation data available at polIII loci. Transcription by RNA polymerase III at SINE loci is suppressed by histone methylation but not by DNA methylation [280], indicating that DNA methylation may not directly influence the expression of tRNA genes but may do so by influencing the surrounding chromatin state. Standard RNA-seq is of limited utility in examining tRNA expression due to the issues of mappability, in addition much RNA-seq data is generated with size and polyA selection methods which would exclude tRNA derived transcripts. Also, tRNAs are a major target for ‘epitranscriptomic’ modification, with an average of 11-13 modifications per tRNA [299,308], some of which stop polymerases from elongating or alter base pairing, creating further mapping challenges. Thus, variants on standard RNA-seq procedures have been developed [309–311]. polIII ChIP-seq has been used as a proxy for tRNA gene expression and, unlike RNA-seq based methods, generates reads from uniquely identifiable flanking regions which map to a known locus of origin [267]. This same advantage exists for the MeDIP-seq data used in this study.

tRNA gene loci may also play a role in large scale genome organisation. tRNA gene clusters act as insulators [312], and have extensive long-range chromatin interactions with other tRNA gene loci [313]. The coordinated transcription of tRNAs at subnuclear foci may represent an organising principle for 3D-chromatin by providing spatial constraints. In both budding and fission yeast tRNA genes localise to the nucleolus. In fission yeast, a subset of B-box sequence elements are bound by TFIIIC and not polIII serving as chromatin anchors to the nuclear periphery and acting as boundaries between euchromatic and heterochromatic regions [314]. It is unclear to what extent tRNA genes may play similar roles in large scale chromatin organisation in other organisms.

In this study ageing-related changes in the epigenetic DNA methylation state of the entire tRNAome were directly investigated, facilitated by the availability of a large-scale MeDIP-seq dataset. Arrays poorly cover this portion of genome. The 450k and EPIC arrays have 110 robust probes covering 84 tRNAs and 129 robust probes covering 89 tRNAs respectively, thus even the latest EPIC arrays cover <15% of the tRNA genes, with robust probes, and in total only ~4.7% of all the tRNA gene CpGs [186].

tRNA genes sit at the heart not only of the core biological process of translation but at a nexus of signalling networks operating in several different paradigms, from small RNA signalling to large scale chromatin organisation [313]. In summation tRNA biology, protein synthesis, nutrient sensing, stress response and ageing are intimately interlinked. The newly described

ageing-related observation of tRNAome DNA hypermethylation has here been identified as and independently replicated.

## 4.3 Methods

Code for analyses performed in this chapter can be found here: [https://github.com/RichardJAtton/tRNA\\_paper\\_code](https://github.com/RichardJAtton/tRNA_paper_code)

### 4.3.1 Participants

#### 4.3.1.1 MeDIP-seq DNA methylomes

Participants in the ‘EpiTwins’ study are adult volunteers from the TwinsUK Register. The participants were aged between 16 and 82 years, with a median of ~55 years (cohort profile [195]). Ethics for the collection of these data were approved by Guy’s & St Thomas’ NHS Foundation Trust Ethics Committee (EC04/015—15-Mar-04) and written informed consent was obtained from all participants.

#### 4.3.1.2 Targeted Bisulfite sequencing

Participants for our targeted bisulfite sequencing of select tRNA loci were drawn from two studies. Samples from participants aged 4 and 28 years are from the MAVIDOS [171] study and participants aged 63 and 78 years are from the Hertfordshire cohort study [315]. Due to a limited number of available samples, the two 4 year old pools contained DNA from 20 individuals each, with all other pools having 25 contributing individuals. Pool 1, the first 4 year old pool used DNA from all male samples, with all other pools using all female samples. Thus, the total number of participants was 190 (see Table 4.7). Samples from the 28 year old time point are all from pregnant women at ~11 weeks gestation.

### 4.3.2 tRNA annotation information

Genomic coordinates of the tRNA genes were downloaded from GtRNAdb [263]. The 2 tRNAs located in chr1\_g1000192\_random are tRNA-Gly-CCC-8-1 & tRNA-Asn-ATT-1-2 ([Supplementary File S1](#)). The 213 probes overlapping tRNA genes were derived from intersecting the tRNA gene annotation data from gtrNAdb with the Illumina 450k array manifest annotation for the hg19 genome build using bedtools v2.17.0 [194]. 107 tRNA genes were excluded from blacklisted regions of hg19 [316].

#### 4.3.2.1 tRNA Gene Clustering

To explore the genomic spatial distribution of the tRNA genes, the tRNA loci were clustered by grouping together all tRNAs within 5Mb of one another using the `bedtools merge` tool v2.17.0 [194]. A command of the form: `bedtools merge -c 4 -o collapse -d <N> -i hg19-tRNAs.bed` was used, where `<N>` is the binsize. The binsize was varied and 5Mb was selected as it is at approximately this size that number of clusters with more than one tRNAs

exceeds the number of singleton tRNAs (Figure 4.5). The further requirements that these groupings contain at least 5 tRNA genes with a density of at least 5 tRNA genes per Mb were added for these groups to be considered clusters.

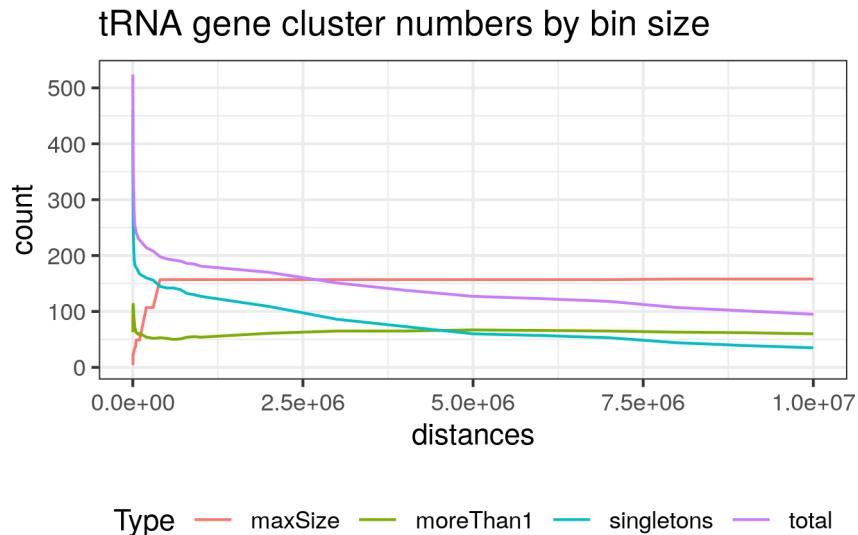


Figure 4.5: **tRNA gene cluster numbers at different bin sizes** total: total number of tRNA clusters. singletons: number of tRNAs in clusters alone. moreThan1: number of tRNAs in clusters with more than one tRNA. maxSize: the number of tRNAs in the largest cluster.

#### 4.3.2.2 tRNA gene mappability assessment

To assess the mappability of tRNA gene regions a value for mappability score density was computed to facilitate comparisons of regions of the genome. Mappability score density is computed as the area under the encode mappability tracks [317] over the length of the region.

### 4.3.3 DNA methylome data

#### 4.3.3.1 TwinsUK MeDIP-seq methylomes

The Methylated DNA Immunoprecipitation sequencing (MeDIP-seq) data was processed as previously described [84,170] and detailed in [Methods 2.2](#). These processed data are available from the European Genome-phenome Archive (EGA) (<https://www.ebi.ac.uk/ega>) under study number EGAS00001001910 and dataset EGAD00010000983 and were generated by BGI Shenzhen for TwinsUK. The dataset used in this work consists of 4350 whole blood methylomes with age data. 4054 are female and 270 male. 3001 have full blood counts. There are 3652 individuals in this data set. These individuals originate from 1933 unique families. There are 1234 monozygotic (MZ) twin pairs (2468 individuals), and 458 dizygotic (DZ) twin pairs (916 individuals). MeDIP-seq data from regions of interest was extracted using Bedtools v2.17.0 [194].

#### 4.3.3.2 Analysis of DNA methylome data for Significant Ageing-related changes

All analysis was performed in R/3.5.2. Linear models were fitted to age using the MeDIP-seq DNA methylome data, as quantile normalised RPM scores at each 500bp window. Quantile normalisation was performed with the `qqnorm` R function with the theoretical quantiles of the RPM values at each window used in subsequent analysis.

Models were fitted with:

1. No covariates
2. Batch information as a fixed effect
3. Blood cell-type counts for neutrophils, monocytes, eosinophils, and lymphocytes as fixed effects
4. Batch and Blood Cell counts as fixed effects
  - Models 1 & 2 were fitted on the full set of 4350 as batch information was available for all samples but blood cell count data was only available for a subset of 3001 methylomes.
  - Models 1 & 2 fitted in the n=3001 subset were similar to those fitted in the complete set of 4350.
  - Models 3 & 4 were fitted in the n=3001 subset with full covariate information and sets of significant tRNAs identified at study-wide and genome wide levels in model 4 were used in subsequent analyses.

Models were also fitted for two unrelated subsets created by selecting one twin from each pair (Monozygotic or Dizygotic), yielding sets with n = 1198 & 1206 DNA methylomes. One additional model was fitted for longitudinal analysis, samples were selected by identifying individuals with a DNA methylome at more than one time point and filtering for only those with a minimum of 5 years between samples. This yielded 658 methylomes from 329 individuals with age differences of 5-16.1 years, median 7.6 years. Models for this set included participant identifier as a fixed effect in addition to blood cell counts and batch information.

#### 4.3.3.3 Permutation Analysis for Enrichment with Age-related Changes

Permutation analysis was performed to determine whether the CpG distribution of sets of the tRNAome was the principle driver of the ageing-related changes observed. Windows overlapping tRNAs have a higher proportion of windows with a greater CpG density than their surrounding sequences (see Figure 4.6). CpGs residing within moderate CpG density loci are the most dynamic in the genome [63] and CpG dense CpG island regions include specific ageing-related changes [125,127,170]. For comparison the permutation was also performed in the CGI regions from the Polycomb group protein target promoters in Teschendorff *et al.* [127] and bivalent loci from ENCODE ChromHMM ‘Poised Promoter’ classification in the GM12878 cell-line [318]. A random set of 500bp windows representing an equivalent CpG density distribution of the feature

set in question were selected from the genome-wide data. Above a certain CpG density there are insufficient windows to sample without replacement within a permutation. Furthermore, above  $\sim \geq 18\%$  CpG density CpG Islands become increasingly likely to hypomethylated [319]. Therefore, all windows with a CpG density of  $\geq 18\%$  (45 CpGs per 500bp) were grouped and sampled from the same pool. i.e. a window overlapping a tRNA gene which had a 20% density could be represented in permutation by one with any density  $\geq 18\%$ . This permutation was performed 1,000 times to determine an Empirical p value by calculating the number of times the permutation result exceeded the observed number of significant windows in the feature set.  $Empirical\ p-value = \frac{r+1}{N+1}$ , where r is the sum of significantly hypermethylating windows in all permutations and N is number of permutations [320].

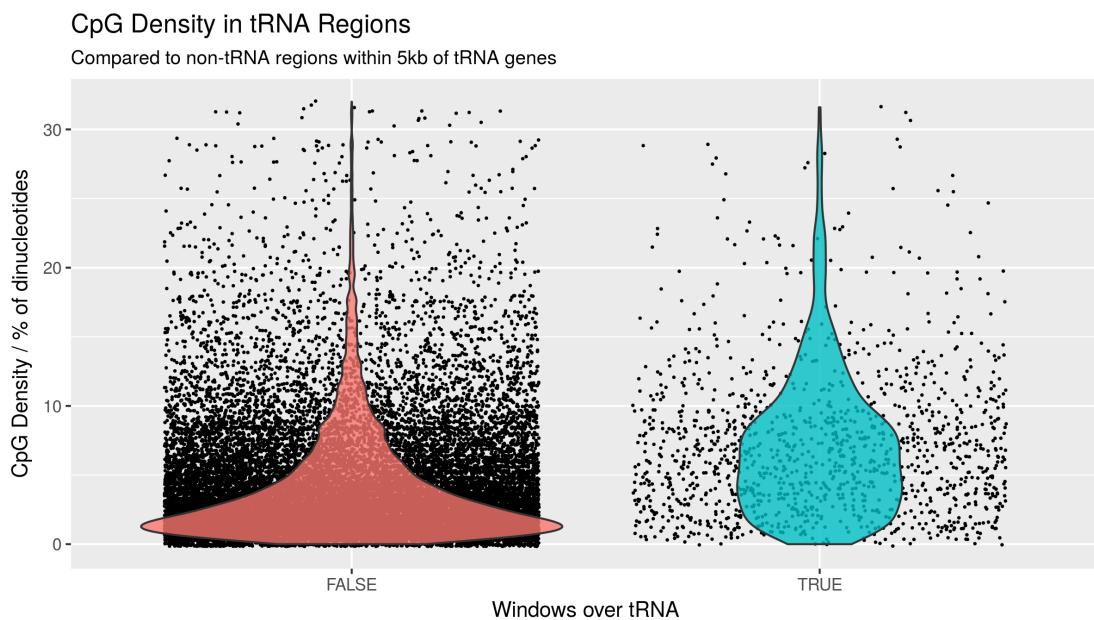


Figure 4.6: CpG Density in windows overlapping tRNA genes compared to that of non-tRNA overlapping windows in arbitrary flanking sequences (+/-5kb)

#### 4.3.3.4 Neonate and Centenarian Whole Genome Bisulfite Sequencing

DNA methylation calls were downloaded from GEO: [GSE31263](#) and intersected with tRNA genes using bedtools v2.17.0 [194].

#### 4.3.3.5 Sample pooling and EPIC array

An Illumina Infinium DNA methylation EPIC array ((C) Illumina) and targeted bisulfite sequencing of select tRNA gene loci were performed. Here DNA was extracted from whole blood and pooled into 8 samples from unrelated individuals at 4 time-points with 2 pools at each time-point. The individual were age matched at approximately 4, 28, 63, and 78 year timepoints (Table 4.7) Using the EPIC array DNAm age was estimated using the Horvath clock [138] and

blood cell-type-composition with the Houseman method [165].

Pool	Mean Age	Sex	Min Age	Max Age	n
Pool 1	4.07	Male	3.99	4.38	20
<b>Pool 2</b>	4.09	Female	3.99	4.36	20
Pool 3	28.07	Female	25.87	29.80	25
<b>Pool 4</b>	28.23	Female	26.05	30.01	25
Pool 5	63.40	Female	62.80	63.80	25
<b>Pool 6</b>	63.26	Female	62.70	63.70	25
Pool 7	77.96	Female	75.50	80.50	25
<b>Pool 8</b>	77.22	Female	74.40	80.10	25

Figure 4.7: Summary information on participants in each pool.

#### 4.3.3.6 Targeted Bisulfite Sequencing

tRNA loci were selected for targeted sequencing which exhibited changes and DNA methylation with age along with closely related tRNAs in which changes were not observed. Primer design was performed using ‘methPrimer’ [200] ([Supplementary File S2](#)). A total of 84 tRNA loci were targeted in 2 rounds of sequencing, 79 subsequently generated reliable results post-QC. The targeted tRNAs covered a total of 723 CpGs with a median of 8 CpGs per tRNA (range 1-13), data passing QC was generated for 458 CpGs, median 6 (range 1-9) per tRNA.

Initial QC for readcount, base call quality and base composition was carried out with [FastQC](#) [192] and [multiqc](#) [321] for combined visualisation of QC outputs. Quality based trimming was performed with [Trim Galore](#) [322] and target specific primer trimming with [cutadapt](#) [323] and custom [perl](#) 5 scripts. Alignment and methylation calling was performed with [Bismark](#) (v0.20.0) [199] making use of [bowtie2](#) [324].

The alignment was performed against both the whole hg19 genome and just the tRNAome +/- 100bp to assess the possible impact of off-target mapping. Mapping to the whole genome did produce purported methylation calls at a larger number of loci than mapping just to the tRNAome (683,783 vs 45,861 respectively). Introducing a minimum coverage threshold of 25 reads dramatically reduced this and brought the number of sites into line with that in the tRNAome set (36,065 vs 33,664 respectively) suggesting a small number of ambiguously mapping reads. All subsequent analysis was performed using the alignment to just the tRNAome with a minimum coverage of 25 reads.

Pairwise differential methylation analysis of the tRNA genes at the different time points was performed using [RnBeads](#) [325] with [limma](#) [326] and a minimum coverage of 25 reads. Linear regression was used to predict age from DNA methylation at the targeted tRNA sites, permitting us to compare rates of increase with age. For the linear regression, only CpG sites with more than 25 reads mapped to the regions of the genome targeted for amplification were used.

#### 4.3.3.7 TwinsUK Illumina 450k array methylomes

Illumina Infinium DNA methylation 450k arrays ((C) Illumina) were also performed on TwinsUK participants, in 770 Blood-derived DNA samples which had matched MeDIP-seq data (Methods 2.1). These data were available for this analysis in a pre-processed form, Methylation ‘beta’ values subject to beta-mixture quantile normalisation (BMIQ) as previously described [84,170]. Cell-type correction was performed using cell-count data and the following model: `lm(age ~ beta + eosinophils + lymphocytes + monocytes + neutrophils)`.

#### 4.3.4 Chromatin Segmentation Data

Epilogos chromatin segmentation data [318,327] was downloaded for the tRNA gene regions +/- 200bp from [https://explore.altius.org/tabix/epilogos/hg19.15.Blood\\_T-cell.KL.gz](https://explore.altius.org/tabix/epilogos/hg19.15.Blood_T-cell.KL.gz) using the tabix utility. The data used was the ‘Blood & T-cell’ 15 State model based on segmentation of 14 cell-types. This data was manipulated and visualised with R and ggplot2. tRNAs were assigned a predominant chromatin state base on the state with the highest score over that tRNA gene.

#### 4.3.5 Isolated Blood Cell Type Specific Data

Data from 7 cell-type fractions from 6 Male individuals was downloaded from GSE35069 [328] using GEOquery [329]. Five of the 6 top age hypermethylating tRNAs are covered by this array dataset.

#### 4.3.6 Cancer and Tissue Specific Methylation Data

Data was downloaded from the TCGA (The Cancer Genome Atlas) via the GDC (genomic data commons) data portal [330] using the GenomicDataCommons R package. Data from foetal tissue [331,332] was downloaded from GEO (GSE72867, GSE30654). From the TCGA, samples were selected for which DNAm data was available from both the primary site and normal solid tissue, and for which an approximate age could be inferred (within one year). Selecting those probes overlapping tRNA genes yielded 73,403 data points across 19 tissues with an age range of 15-90yrs (median 63.4) (Supplementary File S3)

#### 4.3.7 Assaying tRNA expression in blood with MINTmap

Small RNA-seq data from sorted blood cell fractions [333] was used. (GSE100467) and the MINTmap [296] tRNA fragment alignment tool. This dataset covered 42 individuals aged 21-63. A customised MINTmap reference designed to include only fragments which unambiguously map to a single tRNA gene locus and which overlap the 5' or 3' end of the genomic tRNA sequence by at least one base with no mismatches was produced. This reference is intended to capture pre-tRNAs prior to processing and CCA addition operating under the assumption that the levels

of pre-tRNAs will be informative about the amount of transcription taking place at the tRNA loci. This approach provides at most a many to one mapping of tRNA fragments to a tRNA genes.

Assaying the expression of tRNA genes presents numerous difficulties [290], and usually requires variants on standard RNA-seq protocols. This custom MINTmap reference build yielded 383 fragments mapping to 92 distinct tRNA loci in this data. To control quality, only fragments with more than 20 total instances in the dataset, and present in more than 20 individuals were considered.

The maximum length of a fragment was limited to 50nt, due to the read length of the small RNA-seq data.

#### 4.3.8 Mouse RRBS Analysis

Methylation calls and coverage information resulting from RRBS performed by Petkovich *et al.* [334] were downloaded from GEO using GEOquery [329] GSE80672. These data from 152 mice covered 68 tRNA and 436 CpGs after QC requiring >50 reads per CpG and >10 data points per tRNA. 5 tRNAs were excluded for being located within blacklisted regions of mm10 [316]. After QC there were 58 tRNA genes and 385 CpGs. Simple linear modelling to predict age (in months) from methylation level at each tRNA and each CpG were performed in R.

## 4.4 Results

### 4.4.1 DNA Methylation of Specific tRNA Gene Loci Changes with Age

Due to tRNAs critical role in translation and evidence of their modulation in ageing and longevity-related pathways, these genes were interrogated for evidence of ageing-related DNA methylation changes. The discovery set was a large-scale peripheral blood-derived DNA methylome dataset comprising of 4,350 samples (see Figure 4.8).

Blood			Other Tissues
Discovery	Validation	Replication	Tissue Specificity
Method: MeDip-Seq tRNAs: 598 N = 4,350 Ages: 19 - 82 yrs Source: Twins UK	Method: 450k array tRNAs: 158 N = 587 Ages: 18 - 81 yrs Source: Twins UK	Method: Targeted Bisulfite Sequencing tRNAs: 79 N = 190 in 8 pools Ages: 4 - 80 yrs Source: MAVIDOS / Hertfordshire	Method: 27k/450k array tRNAs: 43-115 N = 733 Ages: 0 - 90 yrs Source: TCGA/GDC/GEO 19 Tissues matched Normal and Tumour, 11 Fetal

Figure 4.8: **Study Structure** tRNAs differentially methylated with age initially identified in MeDIP-seq, validated (where covered in 450k array) and replicated in targeted bisulfite sequencing of pooled samples. Tissue specificity of these effects was explored in TCGA and foetal tissue data.

This sequencing-based dataset had been generated by Methylated DNA Immunoprecipitation (MeDIP-seq) [98], which relies on the enrichment of methylated fragments of 200-500 bp to give a regional DNAm assessment (500 bp semi-overlapping windows, see [Methods 4.3.3.1](#)). In total the extended human tRNAome is comprised of 610 tRNAs and closely related sequences (gtRNAdb)(see Figure 4.9), though only 492 are autosomal and do not reside in blacklisted regions of the genome [316].

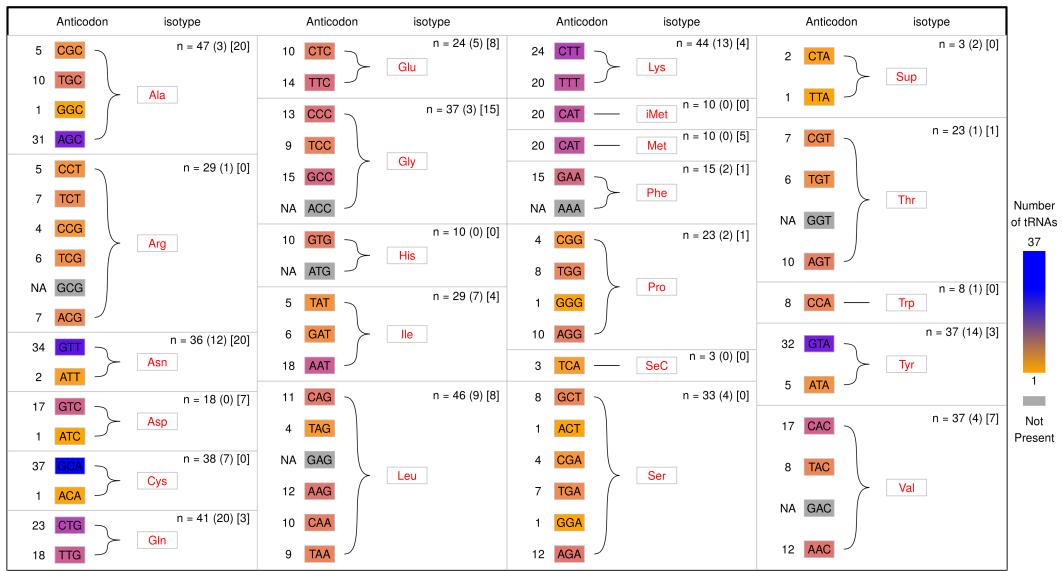


Figure 4.9: The genetic code as instantiated in the human tRNAome. The triplet genetic code leads to the incorporation of specific amino acids into an elongating protein via corresponding tRNAs. n is the number of tRNA genes which encode a given amino acid, the number in parentheses is how many of those may be pseudogenes based on their tRNAscan score [335], and the number in square brackets is the number in blacklisted regions [316]. There are a total of 610 tRNAs and closely related sequences in GtRNAdb [263], 416 of which are high confidence tRNAs, 116 of which are potential pseudogenes, and 107 are in blacklisted regions [316]. Notably 7 of the 61 non-STOP codons are missing from the human tRNAome therefore these codons are handled by wobble base matching (*e.g.* GCG Arg, ACC Gly). Also of note are the suppressor and selenocysteine tRNAs. The 20 methionine tRNAs are split equally between initiator methionine and internally incorporating methionine tRNAs, which are structurally distinct [258]. There are also 23 nuclear encoded mitochondrial tRNAs.

Due to the small size of these tRNAs (60-86bp, median 73bp, excluding introns which are present in ~30 tRNAs with sizes from 10-99bp, median 19bp), this fragment-based method enabled a robust examination of the epigenetic state of these highly similar sequences. This was supported by a mappability assessment. The median mappability score density for the tRNAome was 0.90 for 50mers when considering tRNA genes ±500bp reflecting the regional nature of the MeDIP-seq assay Methods 4.3.2.2. In contrast the 50mer mappability density is 0.68 for the tRNA gene sequences alone without flanking sequences. Excluding the flanking region is representative of the mappability of reads generated using a technique such as whole-genome bisulfite sequencing. This is because there is no IP fragment so reads mapping to adjacent more mappable sequences do not convey information about the methylation state of sites in the same fragment but to which it is harder to map (see Figures 4.10 & 4.11).

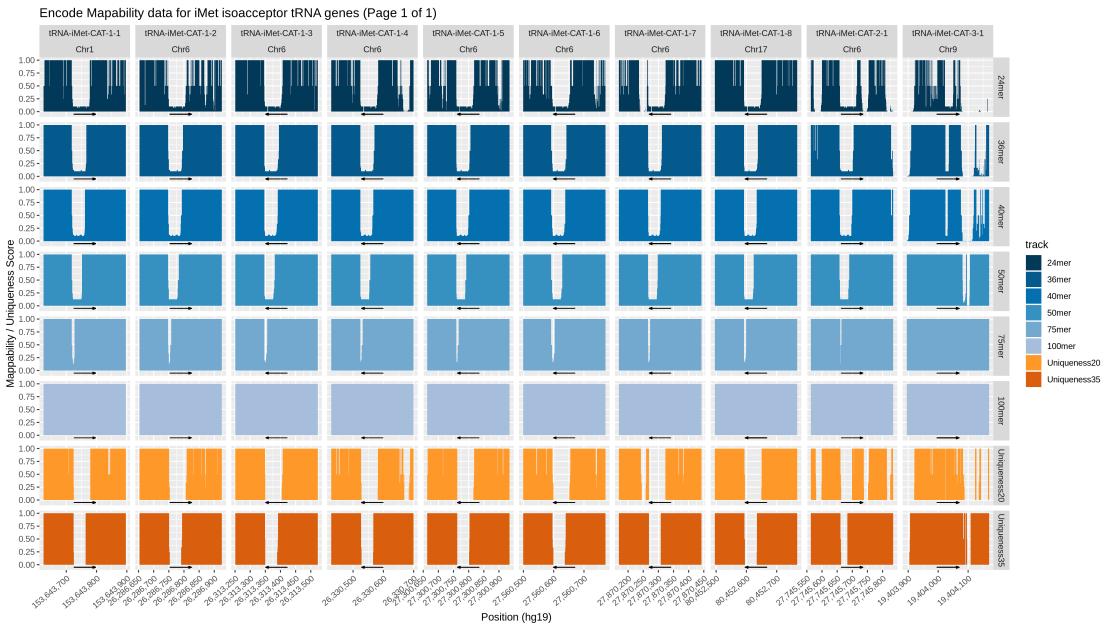


Figure 4.10: Example of mappability data from the encode mappability tracks [317] for the initiator methionine tRNA genes.

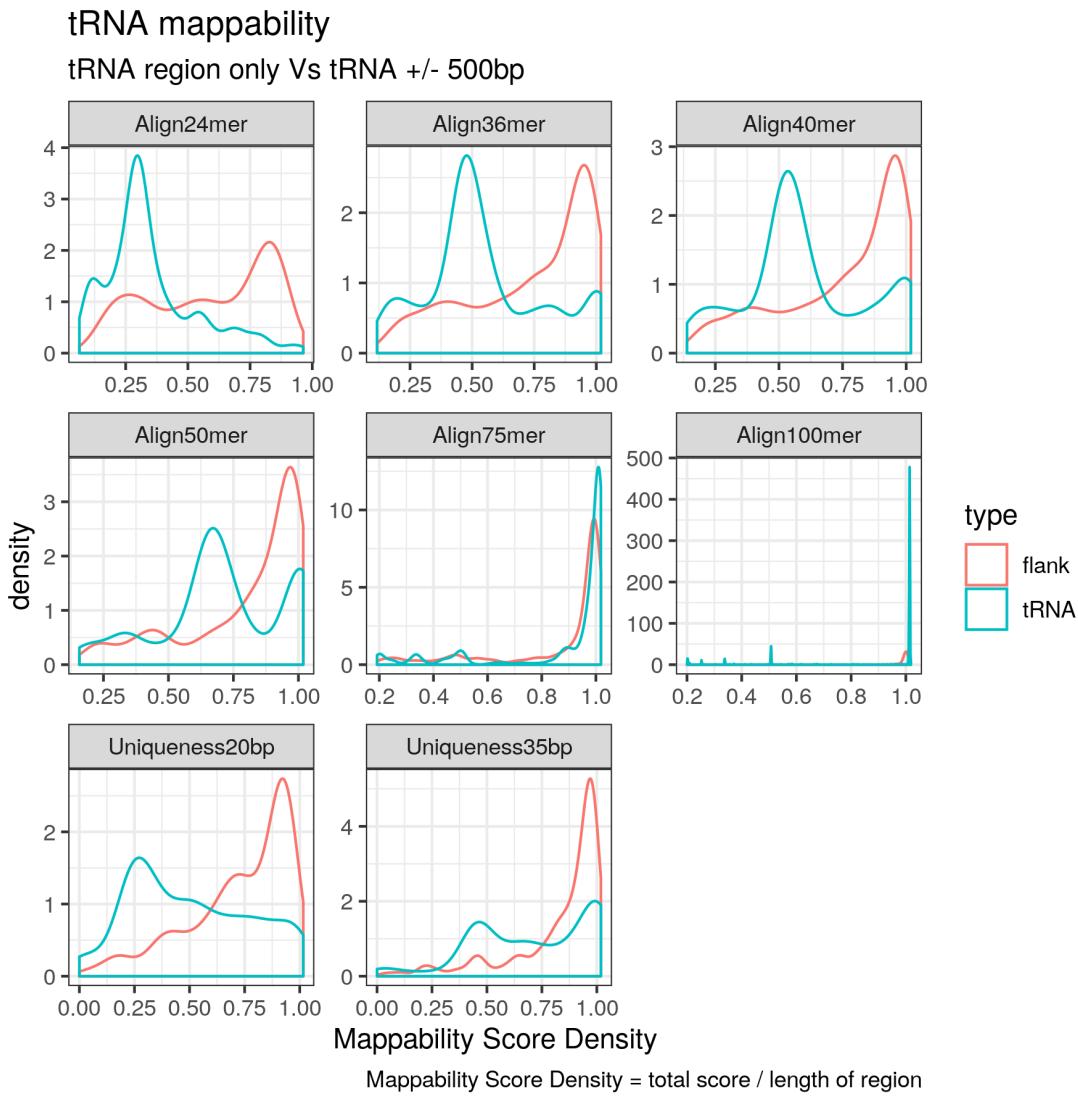


Figure 4.11: Mappability score density of the tRNAome increases with read length and is greater when flanking regions ( $\pm 500\text{bp}$ ) are included. Mappability score density is computed as the area under the encode mappability tracks [317] over the length of the region.

21 genome-wide significant and 44 study-wide significant results were identified via linear regression, ( $p < 4.34 \times 10^{-9}$  and  $8.36 \times 10^{-5}$ , respectively; see [Methods 4.3.3.2](#), batch corrected  $n=4350$ ). Study-wide significance was calculated conservatively using the Bonferroni method for all 598 autosomal tRNAs. There was a strong directional trend with all results at both significance levels being due to increases in DNA methylation. Age-related changes in cell type proportion are strong in heterogeneous peripheral blood, and include a myeloid skew, loss of naive T cells and increases in senescent cells [163]. A subset of 3 genome-wide and 16 study-wide significant hypermethylation results remained significant even after correcting for potential cell-type changes by including lymphocytes, monocytes, neutrophils and eosinophil cell count data ( $n=3001$ , Listed in [Table 4.12](#), Red in [Figure 4.13](#)). The 3 genome wide significant tRNAs

are: tRNA-iMet-CAT-1-4, tRNA-Ile-AAT-4-1, and tRNA-Ser-AGA-2-6. tRNA-iMet-CAT-1-4 is located on chromosome 6. tRNA-Ile-AAT-4-1 and tRNA-Ser-AGA-2-6 are neighbours and are located on chromosome 17 within the 3' UTR of *CTC1* (CST Telomere Replication Complex Component 1). Going forward these most robustly corrected sets of 3 and 16 tRNA genes are referred to as the genome-wide and study-wide significant tRNA genes respectively.

tRNA	Window	MeDIP		450k Array		Targeted BiS-seq	
		Slope	p-value	Slope	p-value	Slope	p-value
tRNA-Gln-CTG-7-1	Chr1:147,800,750-147,801,250	0.84	2.60e-05				
	Chr2:131,094,500-131,095,000	1.11	4.64e-09				
tRNA-Glu-TTC-1-1	Chr2:131,094,250-131,094,750	1.00	1.12e-07				
	Chr2:131,094,750-131,095,250	1.00	3.28e-07				
tRNA-His-GTG-1-2	Chr1:146,544,500-146,545,000	0.92	1.38e-06				
tRNA-His-GTG-2-1	Chr1:149,155,750-149,156,250	1.05	2.98e-08				
	Chr1:149,155,500-149,156,000	0.83	1.37e-05				
tRNA-Ile-AAT-10-1	Chr6: 27,251,500- 27,252,000	1.07	1.45e-08			1.30	1.22e-03
	Chr6: 27,251,750- 27,252,250	0.90	1.86e-06			1.30	1.22e-03
tRNA-Ile-AAT-4-1	Chr17: 8,130,000- 8,130,500	1.19	2.98e-10	19.63	8.92e-06	-0.74	6.88e-04
	Chr17: 8,130,250- 8,130,750	0.77	3.99e-05	19.63	8.92e-06	-0.74	6.88e-04
tRNA-Ile-TAT-2-2	Chr6: 26,987,750- 26,988,250	0.97	7.25e-07	4.16	1.17e-02	-0.60	3.84e-01
tRNA-iMet-CAT-1-4	Chr6: 26,330,500- 26,331,000	1.28	2.83e-11	13.01	6.07e-06	4.54	9.35e-04
	Chr6: 26,330,250- 26,330,750	1.13	2.89e-09	13.01	6.07e-06	4.54	9.35e-04
tRNA-Leu-TAG-2-1	Chr14: 21,093,250- 21,093,750	1.04	9.38e-08			2.49	8.77e-03
	Chr14: 21,093,500- 21,094,000	0.94	8.50e-07			2.49	8.77e-03
tRNA-Pro-AGG-2-2	Chr6: 26,555,500- 26,556,000	1.04	3.97e-08				
	Chr6: 26,555,250- 26,555,750	1.01	9.58e-08				
tRNA-Ser-ACT-1-1	Chr6: 27,261,250- 27,261,750	0.97	3.53e-07			0.66	1.45e-01
tRNA-Ser-AGA-2-6	Chr17: 8,129,750- 8,130,250	1.21	1.16e-10	20.87	6.72e-05	0.62	4.28e-02
	Chr17: 8,130,000- 8,130,500	1.19	3.03e-10	20.87	6.72e-05	0.62	4.28e-02
tRNA-Ser-TGA-2-1	Chr6: 27,513,000- 27,513,500	0.90	3.58e-06	87.21	1.38e-04	-0.25	5.74e-01
tRNA-Val-AAC-1-2	Chr5:180,590,750-180,591,250	0.91	3.28e-06				
tRNA-Val-AAC-4-1	Chr6: 27,648,500- 27,649,000	1.07	1.25e-08	40.06	9.90e-03		
	Chr6: 27,648,750- 27,649,250	0.95	4.31e-07	40.06	9.90e-03		
tRNA-Val-CAC-2-1	Chr6: 27,247,750- 27,248,250	0.85	2.33e-05	59.16	5.05e-06		

Figure 4.12: Study-wide significantly hypermethylating tRNAs in blood cell-type and batch corrected model MeDIP-seq. With Corresponding results in Twins UK 450k array (blood cell-type corrected) and targeted bisulfite sequencing results. Age models for array and targeted BiS-seq were calculated using all probes / CpGs overlapping the indicated tRNA gene. ‘Slope’ corresponds to the beta value for methylation in the linear model, orange colouring indicates hypermethylation and blue hypomethylation. p-values coloured such that low values are dark blue and high values are yellow. blank grey cells indicate missing data.

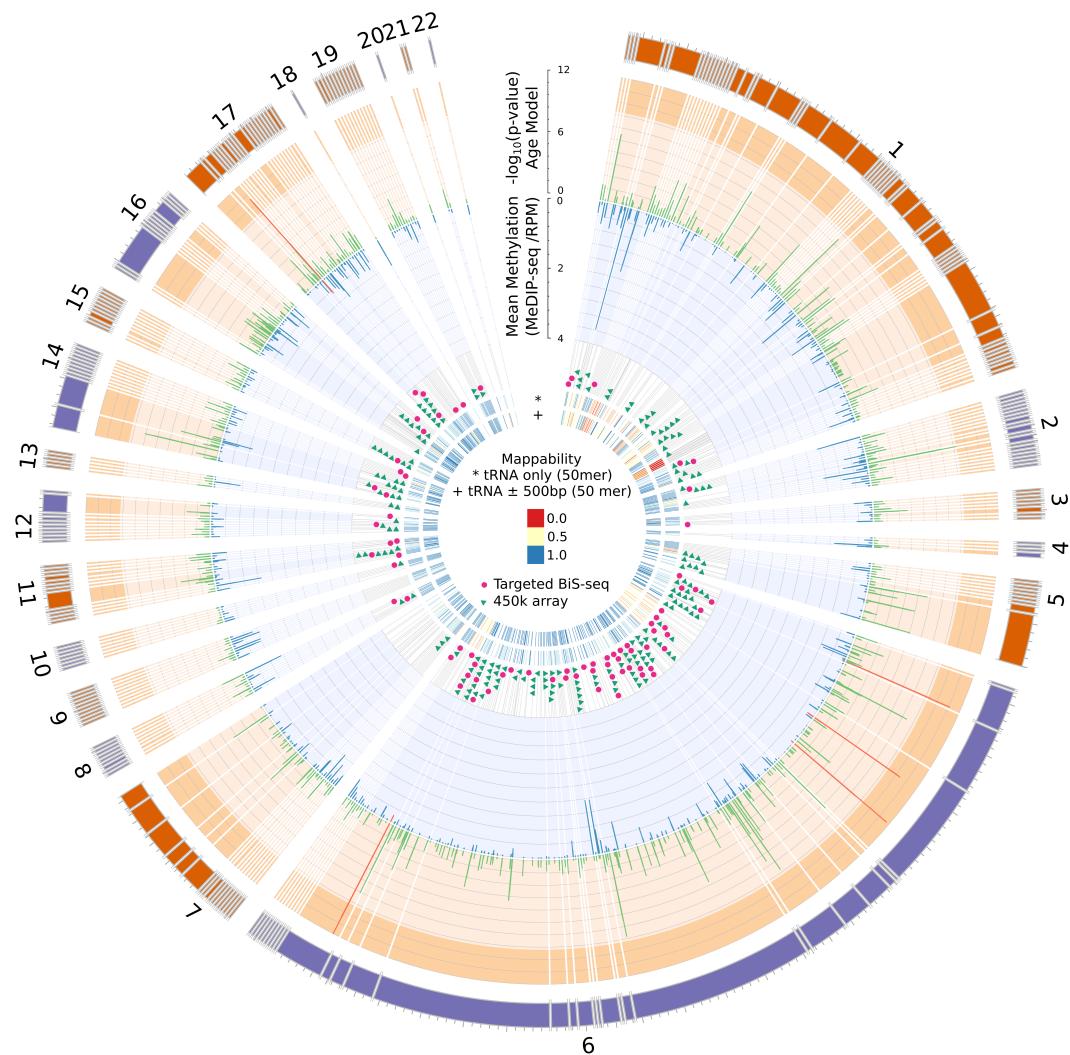


Figure 4.13: Human tRNAome overview. From the outside in: Chromosome ideograms scaled by the number of tRNA genes (total = 598), as excludes chromosome X (10), Y (0) and contig chr1\_gl000192\_random (2; see [Methods 4.3.2](#)). tRNA genes within 20kbp of one another are grouped with breaks inserted between these clusters. Radial grey lines represent the location of tRNA genes in the genome.  $-\log_{10}(p - \text{value})$  for the blood cell-type and batch corrected age model are shown for each window overlapping a tRNA gene in green. Mean methylation across all samples ( $n=3001$ ) in RPM (reads per million base pairs) is shown in blue. Genome-wide significant cell-type & batch corrected ( $p < 4.34 \times 10^{-9}$ ) tRNAs show in red. The 158 Loci covered by 213 probes on the 450k array which directly overlap a tRNA gene are shown with green triangles. The 84 loci targeted for bisulfite sequencing in this study are indicated in magenta. Mappability score density is computed as the area under the encode mappability tracks [317] over the length of the region.

Due to the related nature of these twin samples, these data were also analysed in two subsets

of n=1198 & 1206 by selecting one twin from each pair into the separate sets. This analysis also included correction for Batch and Blood Cell counts. Whilst in these smaller datasets no tRNAs were genome-wide significant, 5 and 7 tRNA genes, respectively, reached study-wide significance, with consistent hypermethylation. In these sets 5/5 and 6/7 of these were present in 16 study-wide significant tRNA genes.

Furthermore, a subset of samples with longitudinal data were examined (n=658 methylomes from 329 individuals, median age difference 7.6 yrs). At the nominal significance threshold ( $p < 0.05$ ) this yielded a split of 41 hypermethylating tRNA genes and 22 hypomethylating tRNA genes. Of these hypermethylated tRNAs, 2 are in the previously identified genome-wide significant set of 3 (with tRNA-iMet-CAT-1-4 ranked 3rd by p-value) and 9 are in the study-wide significant set of 16.

#### 4.4.2 tRNA Genes are Enriched for Age Related DNA Hypermethylation

Whilst ageing changes are pervasive throughout the DNA methylome, a strong enrichment for such changes occurs within the discrete tRNAome (Fisher's Exact Test  $p = 1.05 \times 10^{-27}$ ) (see Figure 4.14). This is still significant if the 6 of the study-wide significant 16 tRNAs that overlap polycomb or bivalent regions are excluded ( $p = 4.66 \times 10^{-15}$ )

CpG density is known to have a clear impact on the potential for variability of the DNA methylome as well as ageing-related changes [63,124]. To assess whether this hypermethylation finding was being driven merely by the inherent CpG density of the tRNAome, a CpG density matched permutation analysis was performed (1,000X, see Methods 4.3.3.3). This supported the specific nature of these age-related DNAm changes within the functional tRNAome (Empirical p-value  $< 1.0 \times 10^{-3}$ , Figure 4.14 b). As a point of comparison for this genomic functional unit, the same permutation analysis was performed for the known age-related changes in the promoters of genes that are polycomb group targets [127] and those with a bivalent chromatin state [125]. The enrichment of the polycomb group targets and bivalent regions (Empirical p-value  $< 1.0 \times 10^{-3}$ ) was reproduced in this dataset.

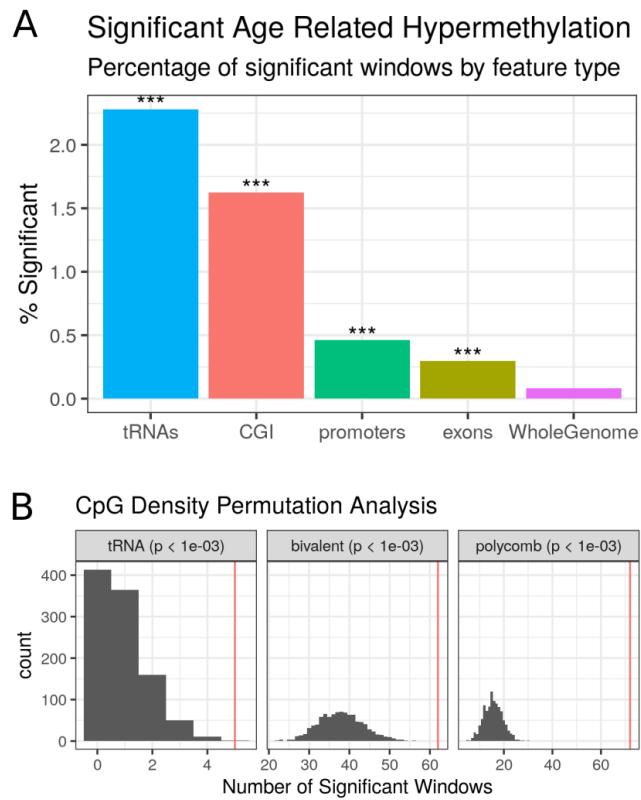


Figure 4.14: **A)** tRNA genes are enriched for age-related hypermethylation compared to the genomic background, (Fisher's Exact Test  $p < 1.05 \times 10^{-27}$ ,  $n = 3001$ ). **B)** tRNA genes show more significant hypermethylations than CpG Density matched permutations. Each permutation represented a random set of windows matching the CpG density of the functional unit (bivalent domains, polycomb group target promoters & the tRNA genes). These are subsequently assessed for significant age-related DNAm changes (see [Methods 4.3.3.3](#)). The red line is the observed number of significant loci.

tRNA-iMet-CAT-1-4 is located in the largest tRNA gene cluster in the human genome at chr6p22.2-1 ([Methods 4.3.2.1](#)). This cluster contains 157 tRNA genes spanning the 2.67Mb from tRNA-iMet-CAT-1-2 to tRNA-Leu-AAG-3-1, and also hosts a histone gene microcluster. tRNA-Ile-AAT-4-1 and tRNA-Ser-AGA-2-6 are neighbours and are located on chromosome 17 (Figure 4.15). Notably tRNA-Ile-AAT-4-1 and tRNA-Ser-AGA-2-6 have a third close neighbour tRNA-Thr-AGT-1-2 which does not show significant age-related hypermethylation. A similar pattern of sharp peaks of significance closely localised around the other loci was observed in the study-wide significant set. GENCODE 19 places tRNA-Ile-AAT-4-1 in the 3'UTR of a Nonsense-mediated decay transcript of *CTC1* (CST Telomere Replication Complex Component 1, CTC1-002, ENST00000449476.2) and not of its primary transcript.

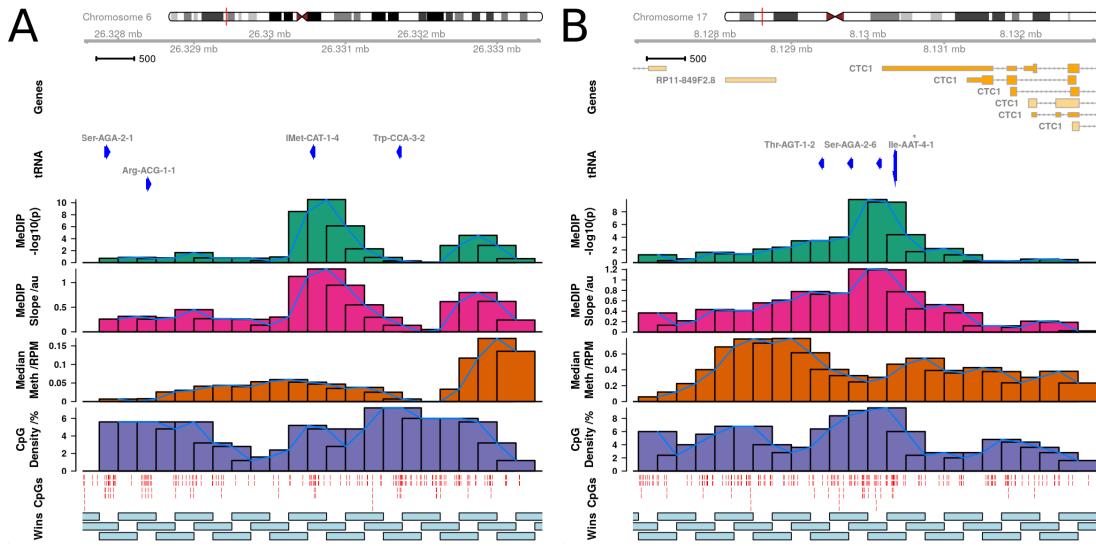


Figure 4.15: **A)** tRNA-iMet-CAT-1-4 as well as **B)** tRNA-Ser-AGA-2-6 and tRNA-Ile-AAT-4-1 exhibit age related DNA hypermethylation in MeDIP-seq data. The sharp peaks suggest that this effect is localised to individual tRNA genes. One of the windows overlapping tRNA-Ile-AAT-4-1 also partially overlaps tRNA-Ser-AGA-2-6. Results shown are from the blood cell-type and batch corrected model. Median Methylation in reads per million is calculated across all samples in the model. The window structure track (Wins) illustrates the tiled nature of the MeDIP-seq windowing. Some annotations of the 3' UTR of CTC1 extend further than illustrated here covering tRNA-Thr-AGT-1-2. CpG density and CpG position are included to illustrate that similarly CpG dense tRNA genes are exhibiting differing age related DNAm change patterns.

#### 4.4.2.1 tRNA gene clustering

To place these hypermethylating tRNA genes in their genomic context the clustering for the extended set of 44 non-blood cell-type corrected study-wide significant tRNAs was examined. The tRNA genes were clustered by grouping together all tRNAs within 5Mb of one another and then required that a cluster contain at least 5 tRNA genes with a density of at least 5 tRNA genes per Mb (see [Methods 4.3.2.1](#)). This yielded 12 major tRNA gene clusters containing a total of 353 tRNA genes, 42 of the 44 study-wide significant tRNA genes in the non cell-type corrected age model reside within these clusters (figure 4.16 a). The hypermethylating tRNA genes are spread evenly among these clusters proportionately to their size, (no significant difference in a one-way ANOVA of percentage of significant tRNAs by cluster).

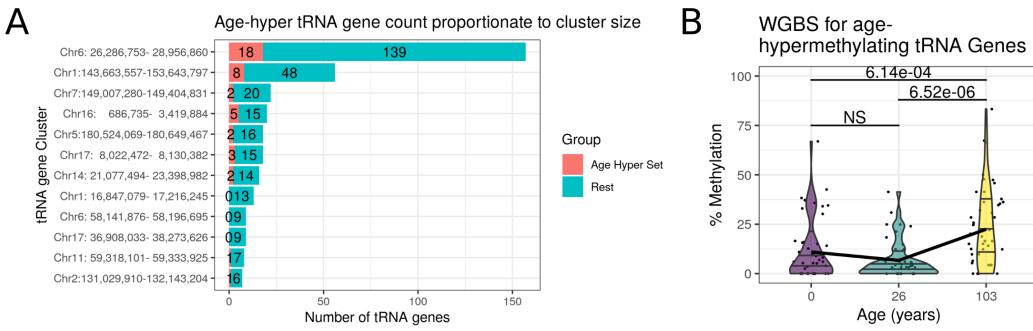


Figure 4.16: **A)** Almost all tRNA genes (42/44, counts shown in red) study-wide significant in the MeDIP batch corrected age model, and all of the genome-wide significant and blood-cell type corrected sets of tRNA genes reside in one of 12 major tRNA gene clusters. (Defined by joining all tRNAs within 5Mb of one another and requiring at least 5 tRNA genes per cluster with a density of at least 5 tRNA genes per Mb). **B)** Available study-wide significant (SWS) tRNA genes ( $n = 14$ ) are more methylated in a centenarian than in a neonate or a 26 year old. Whole Genome Bisulfite Sequencing Data in a newborn, as adult and a centenarian. Each point represents the methylation level at an individual CpG within a tRNA gene.

#### 4.4.2.2 Age-related tRNAome DNA Hypermethylation is even observed in one Newborn versus one Centenarian

Whole genome bisulfite data was sought out for two reasons: 1) The greater positional resolution of WGBS permits one to check if the age-related hypermethylation could be seen at CpGs located within the tRNA genes themselves. 2) As DNA methylation levels are estimated based on the ratio of methylated to unmethylated reads using WGBS, seeing an increased DNA methylation with age is not likely to be due to a copy number expansion or differential read distribution with age which could confound the MeDIP-seq signal. An available Whole Genome Bisulphite sequencing (WGBS) dataset from Heyn *et al.* [132] was examined (see Methods 4.3). These data consisted of blood-derived DNA WGBS in one newborn child and one 26 year old, and centenarian (103 years). In their analysis, the centenarian was found to have more hypomethylated CpGs than the neonate across all genomic compartments, including promoters, exonic, intronic, and intergenic regions. However, even in this examination of 3 individuals of 3 different ages in the 55% of tRNA genes that possessed coverage, DNA hypermethylation with age was observed among the study-wide significant hypermethylating tRNA genes. The centenarian was significantly more methylated in this set of tRNAs than the neonate (Wilcoxon rank sum test, 6.14% increase (95% CI -Inf - 4.31),  $W = 717$ ,  $p = 6.14 \times 10^{-4}$ , see Figure 4.16 b).

#### 4.4.2.3 Age-related Changes Independently Replicated with Targeted Bisulfite Sequencing

In order to further robustly support these-ageing related changes, an attempt was made to replicate these findings in an independent ageing dataset. Furthermore, a different technology was employed, Targeted bisulfite (BiS) sequencing, to further validate the MeDIP-seq-derived results. These data provide individual CpG resolution to identify what may be driving the regional DNAm changes observed.

This targeted BiS-seq was performed in blood-derived DNA from 8 pools of age-matched individuals at 4 time-points (~4, ~28, ~63, ~78 years) from a total of 190 individuals, as detailed in Table 4.7. A total of 79 tRNA loci generated reliable results post-QC (see Methods 4.3.3.6). These tRNAs covered a total of 458 CpGs with a median of 6 CpGs per tRNA (range 1-9). Median Coverage per site across pools, technical replicates and batches was 679 reads (mean 5902).

Firstly, the 8 Pooled samples were run on the Illumina EPIC (850k) array to confirm that this pooling approach was applicable for DNAm ageing-related evaluation. This showed an  $R^2 = 0.98$  between pool mean chronological age and Horvath clock DNAm predicted age [138](see Figure 4.17). This approach permitted the assay tRNA gene DNA methylation across a large number of individuals whilst requiring a minimal amount of DNA from each individual (80-100ng), and costing ~1/24th as much as performing sequencing individually. Therefore, this confirmed the utility of this novel pooling approach. These DNA methylation array derived data were also used to estimate the major blood cell proportions for each of these pools with the Houseman algorithm [165].

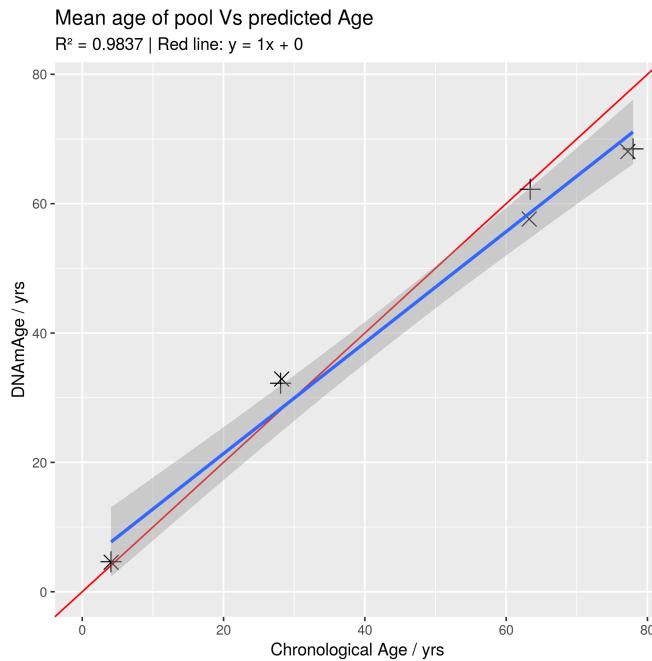


Figure 4.17: Mean chronological age is tightly correlated with DNAm Horvath clock [138] predicted age for the 8 pooled samples. (See Table 4.7 for pool details).

It was observed that individual tRNA loci exhibiting age-related changes in DNAm had duplicate or isodecoder (same anticodon but body sequence variation) sequences in the genome, which despite exact or near sequence identity did not show similar changes. tRNA-iMet-CAT-1-4 for instance is 1 of 8 identical copies in the genome and was the only locus that showed significant changes. The results of pairwise differential methylation tests between age groups for the 6 top tRNAs from the MeDIP-seq models are listed in Table 4.19.

Of the 3 top hits in MeDIP-seq, tRNA-iMet-CAT-1-4 (Figure 4.20c) and tRNA-Ser-AGA-2-6 (Figure 4.20i) exceeded nominal significance ( $p$ -values =  $9.35 \times 10^{-4}$  &  $4.28 \times 10^{-2}$ , respectively). tRNA-Leu-TAG-2-1 from the study-wide significant set also showed nominally significant hypermethylation with age (Figure 4.20u). Also, four of the individual CpGs in tRNA-iMet-CAT-1-4 exhibited nominally significant increases in DNAm with Age (Figure 4.18). However, tRNA-Ile-AAT-4-1 (Figure 4.20n) showed a nominal decrease in DNAm with age.

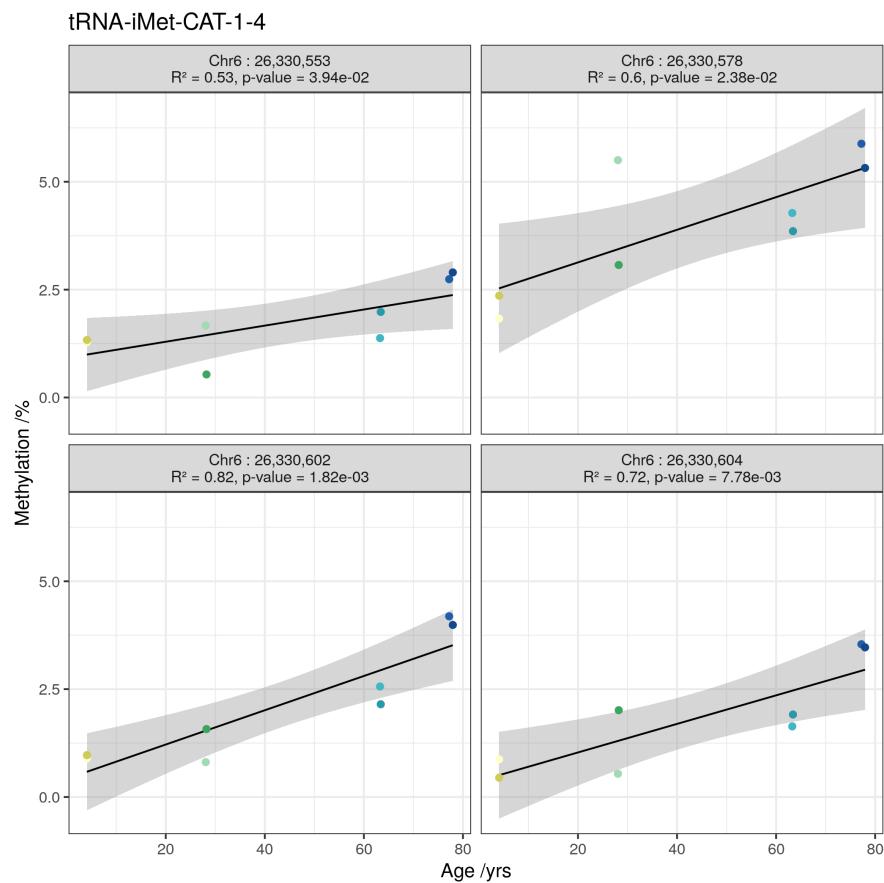


Figure 4.18: Individual CpG methylation increases (nominally significant  $p < 0.05$ ) in tRNA-iMet-CAT-1-4.

tRNA	num. CpGs	comparison	p-value	delta
tRNA-Ile-AAT-4-1	8	4 vs. 28	1.518e-01	-0.2
		4 vs. 63	1.774e-01	-0.234
		4 vs. 78	3.060e-01	0.0113
		28 vs. 63	7.152e-01	-0.0334
		28 vs. 78	1.553e-01	0.212
		63 vs. 78	2.057e-01	0.245
tRNA-iMet-CAT-1-4	5	4 vs. 28	8.403e-02	0.0116
		4 vs. 63	1.716e-01	0.0125
		4 vs. 78	1.997e-04*	0.0368
		28 vs. 63	3.943e-01	0.000869
		28 vs. 78	1.724e-02*	0.0252
tRNA-Ser-AGA-2-6	9	4 vs. 28	4.222e-01	0.0573
		4 vs. 63	3.968e-01	0.0274
		4 vs. 78	4.651e-01	0.0423
		28 vs. 63	1.095e-01	-0.0299
		28 vs. 78	2.126e-01	-0.015
		63 vs. 78	2.201e-01	0.0149

Figure 4.19: Pairwise Differences in Methylation between Age groups by tRNA. p-values are for pairwise methylation differences (see [Methods 4.3.3.6](#))[\[325\]](#).

**4.4.2.3.1 Select Duplicates & Isodecoders of Hypermethylating tRNA loci remain unchanged** A selection of these duplicate and isodecoder loci were targeted for bisulfite sequencing in order to confirm that the identified DNAm changes are specific to a given locus and not general to related tRNAs. Examining the tRNA-iMet-CAT-1 family, only the previously identified 1-4 version confirmed significant hypermethylation (not 1-2, 1-3 or 1-5)(Figure 4.20a-e). Likewise the tRNA-Ser-AGA-2-6 version was supported compared to 2-1,2-4 and 2-5(Figure 4.20f-j)).

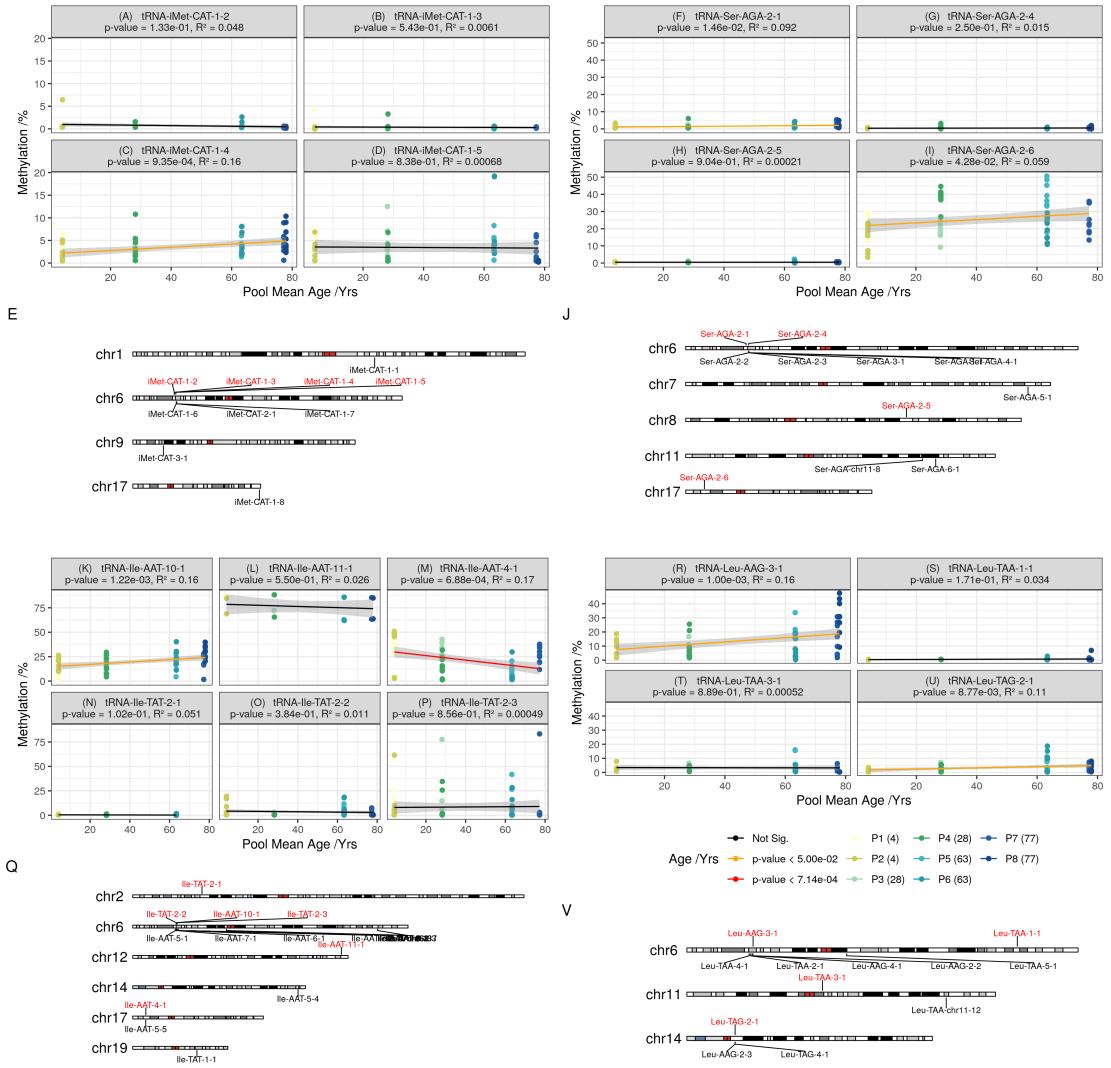


Figure 4.20: Combined CpGs within tRNA loci results (experiment-wide Bonferroni  $p = 7.14 \times 10^{-4}$ ); **(A-D)** Comparison of select tRNA-iMet-CAT loci: Hypermethylation is specific to iMet-CAT-1-4 (c) not other isodecoders (A, B, & D); **(F-I)** Comparison of select tRNA-Ser-AGA loci: Hypermethylation is specific to Ser-AGA-2-6 A (viii) and to a lesser extent Ser-AGA-2-1 (F), whilst not other isodecoders (G, H); **(K-P)** Comparison of select tRNA-Ile loci: Hypermethylation is specific to Ile-AAT-10-1 (K), Ile-AAT-4-1 (M) displays hypomethylation contrary to previous MeDIP findings, Ile-TAT-2-2 & 2-3 lack hypermethylation (previously non-significant in blood-corrected MeDIP, although significant in uncorrected), whilst no change in Ile-AAT-11-1 (L) and Ile-TAT-2-1 (N); **(R-U)** Comparison of select tRNA-Leu loci: Hypermethylation in Leu-AAG-3-1 (R) consistent with 450k and Leu-TAG-2-1 (U) consistent with MeDIP, whilst no change in Leu-TAA-1-1 (S) & Leu-TAA-3-1 (T).

#### 4.4.2.4 DNA methylation 450k Array Data Validates the MeDIP-seq Results

Although DNA methylation arrays poorly cover the tRNAome, an attempt was made to ascertain if this bisulfite conversion-based but differing and well-established technology was supportive at all of the previous DNA hypermethylation findings. TwinsUK had available 450k array on 587 individuals, and this platform includes 143 probes, covering 103 tRNAs. All the 3 top tRNAs in the MeDIP-seq results were covered by this data set, and 7 of the 16 study-wide significant set. 9 tRNAs show significant ( $p < 4.58 \times 10^{-4}$ ) increases in DNA methylation with age in models corrected for blood cell counts including all 3 of the 3 tRNAs identified in the MeDIP-seq as genome-wide significant and 5 of the 7 study-wide significant set present on the array (Figure 4.21). Although it should be noted that 56 of these 143 probes are within the non-robust set of Zhou et al. [186], including 1 of the genome-wide, and 1 of the study-wide results (covering tRNA-Ile-AAT-4-1 & tRNA-Val-AAC-4-1), respectively.



Figure 4.21: Volcano-like plot. tRNAs are labelled if they are significant here or were in the MeDIP-seq data (Red). Model slope: the model coefficient for the methylation values. Unfilled circles indicate those probes in the general mask generated by Zhou et al. [186]. Significance threshold:  $0.05/103 \approx 4.58 \times 10^{-4}$  (the number of tRNA genes examined).

#### 4.4.2.5 Ageing-Related tRNA Loci show increased Enhancer-Related Chromatin Signatures

The activity of the tRNAome was further explored using public chromatin segmentation data in blood (Epilogs Blood & T-cells set) [327]. This shows proportionally more Enhancer-related (Enh, EnhBiv & EnhG) chromatin states at tRNA genes hypermethylating with age than the stronger Promoter-related (TSS) in other tRNAs. (Figure 4.22). Whereas these characteristics are less frequently predominant in the rest of the tRNAs (Figure 4.22). Age-hypermethylating tRNA are enriched for enhancer chromatin states compared to the rest of the tRNAome (Fisher's Exact test  $p = 0.01$ ).

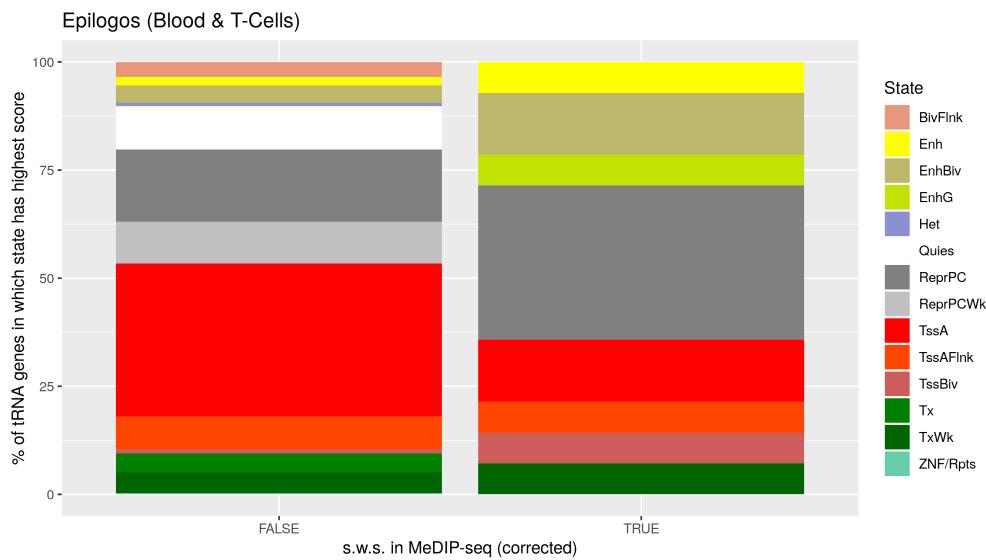


Figure 4.22: Chromatin segmentation data from the Epilogs [327] ‘Blood & T-cell’ 15 State model (tRNA genes +/- 200bp). Frequency with which a model state was the predominant state at a given tRNA. Proportions of predominant tRNA state for the 14 study-wide significant age-hypermethylating tRNAs covered compared to other 371 available tRNAs.

#### 4.4.2.6 Age Hypermethylating tRNAs are more methylated in Lymphoid than Myeloid cells

Three tRNA genes remained genome-wide significant and 16 study-wide significant following correction for major cell-type fraction. This is suggestive of either cell-type independent change or, presumably less likely, a very large effect in a minor cell-type fraction. tRNAs have exhibited tissue-specific expression [284,285,336] and blood cell-type populations change with age. Specifically, there is shift to favour the production of cells in the myeloid lineage [163]. These points lead us to examine tRNA gene DNAm in sorted cell populations. A publicly available 450k array dataset [328]) that has been used in the construction of cell-type specific DNAm references for cell-type fraction prediction using the Houseman algorithm [165] was used (see [Methods 4.3.5](#)).

This consists of data from 6 individuals (aged  $38 \pm 13.6$ /yrs) from seven isolated cell populations (CD4+ T cells, CD8+ T cells, CD56+ NK cells, CD19+ B cells, CD14+ monocytes, neutrophils, and eosinophils). It was found that tRNA gene DNAm could separate myeloid from lymphoid lineages (Figures 4.23 & 4.24).

Of the eight study-wide significant tRNAs with array coverage, it was identified that collectively these eight are significantly more methylated in the lymphoid than the myeloid lineage (1.1% difference, Wilcoxon rank sum test  $p = 1.50 \times 10^{-6}$  95% CI  $0.7\%-\infty$ ). Thus, any age related increases in myeloid cell proportion would be expected to dampen rather than exaggerate the observed age-related hypermethylation signal. In addition tRNA-Ile-AAT-4-1 and tRNA-Ser-AGA-2-6 have the highest variance in their DNAm of all 129 tRNAs covered in this dataset. This could represent ageing-related changes as these samples range across almost 3 decades. Another possibility may be that these loci as well as hypermethylating are also increasing their variability with age in a similar fashion to those identified by Slieker *et al.* [120]. In that study they identified that those loci accruing methylomic variability were associated with fundamental ageing mechanisms.

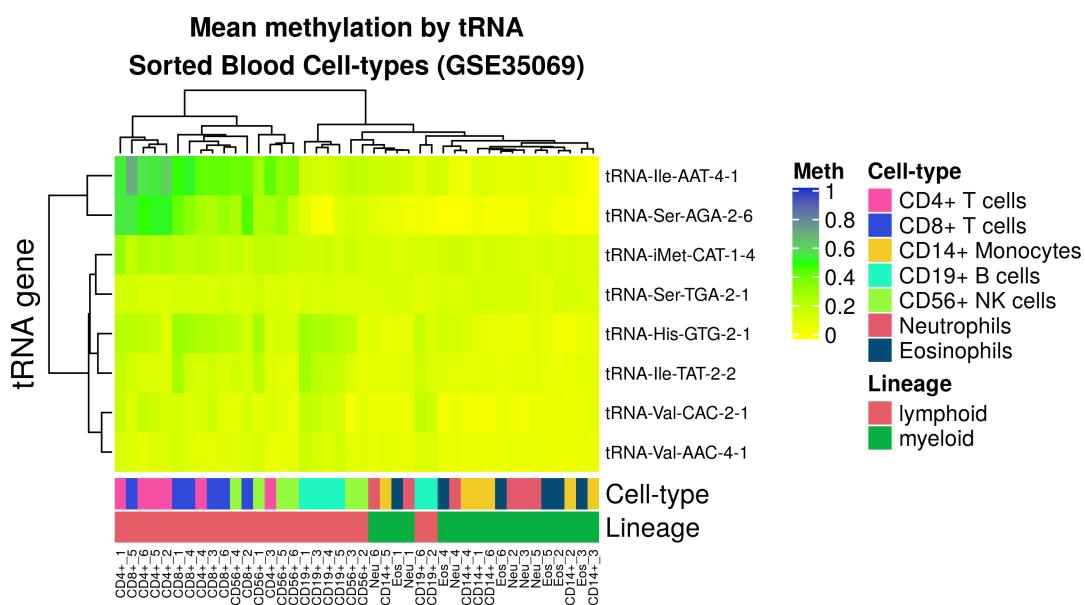


Figure 4.23: Heatmap [337] of mean methylation of probes covering each tRNA in 7 cell-type fractions from 6 Male individuals. Data from GSE35069 [328]. Of the 16 study-wide significant hypermethylating tRNAs, 8 are covered by this dataset.

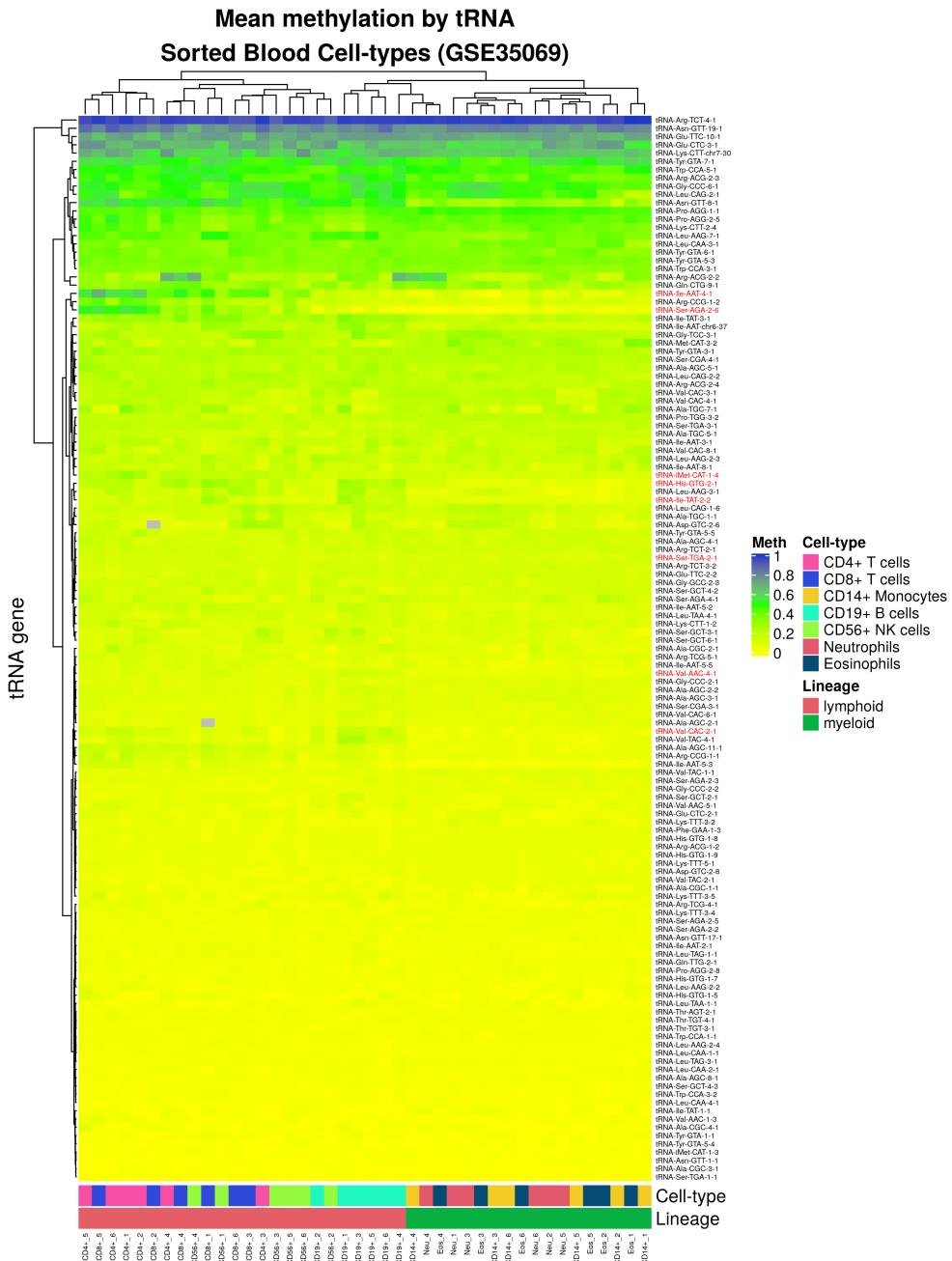


Figure 4.24: Heatmap Mean Methylation of probes covering each tRNA in 7 cell-type fractions from 6 Male individuals. Showing all 150 tRNAs covered by 213 probes on the Illumina 450k array. Data from [GSE35069](#) [328] downloaded using [GEOquery](#) [329]. Generated with the [ComplexHeatmap](#) R package [337].

#### 4.4.3 tRNA Gene DNA Methylation in Other Tissues

Some tRNA gene expression has been shown to be highly tissue specific [284,285,336]. It follows that our observations of changes in DNAm with age in blood might be specific to that tissue. A mix of 450k and 27k array data from ‘solid tissue normal’ samples made available by TCGA

(The Cancer Genome Atlas) and data from foetal tissue [331,332] were downloaded from GEO for use in this analysis (see [Methods 4.3.6](#)). The samples from TCGA range in age from 15-90 (n = 733). Only 43 tRNA genes had adequate data to compare across tissues in this dataset and 115 in the foetal tissue data.

#### 4.4.3.1 tRNA Genes also Hypermethylate with Age in Solid Tissue

Only 2 of the 3 tRNA genes identified as genome-wide significant and a further 1 of the study-wide significant tRNA genes are present in the set of 45 tRNA genes in the TCGA data, thus limiting our ability to draw conclusions about the tissue specificity of these results. Solid tissue samples have a strong preponderance for low levels of methylation consistent with the active transcription of many tRNA genes and show slight increasing methylation with age but age accounts for very little of the variance (linear regression slope estimate = 1.52;  $R^2 = 0.0002$ ; p-value  $1.34 \times 10^{-3}$  ([Figure 4.27d](#)). In a pan-tissue analysis it was found that 10 tRNA genes showed changes in DNAm with age, 9 of which were hypermethylation (p-value  $< 1.1 \times 10^{-3}$ ). One of these tRNA genes, tRNA-Ser-TGA-2-1 was also present in study-wide significant set of tRNA genes. Figures [4.25](#) & [4.26](#) illustrate minimal tissue specific differences. Interestingly, however, tRNA-iMet-CAT-1-4 and tRNA-Ser-AGA-2-6 appeared more variable in methylation state than many other tRNAs in the TCGA normal tissue samples ([Figure 4.25](#)) and indeed have the highest variance in DNA methylation across tissues ([Figure 4.27c](#)).

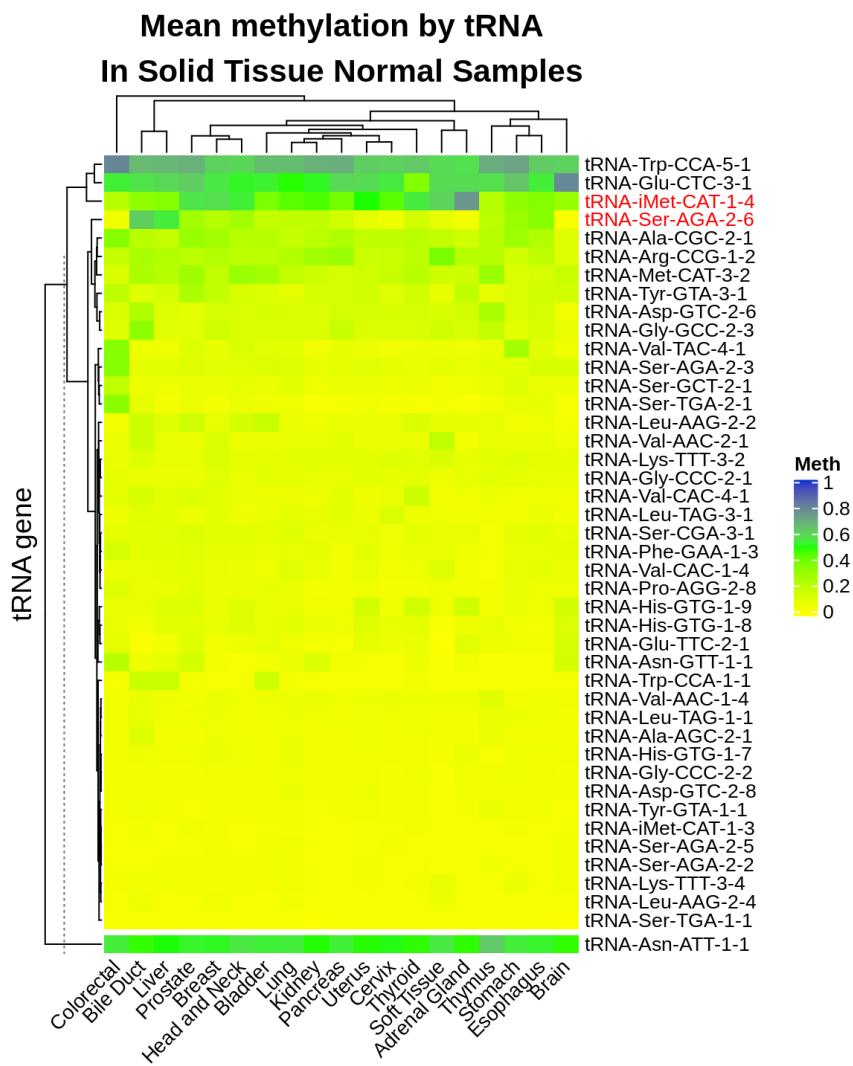


Figure 4.25: Mean Methylation of 43 tRNAs in 19 tissues. Possible pseudogene (tRNA-Asn-ATT-1-1) is shown in a separate cluster beneath the main heatmap [337].

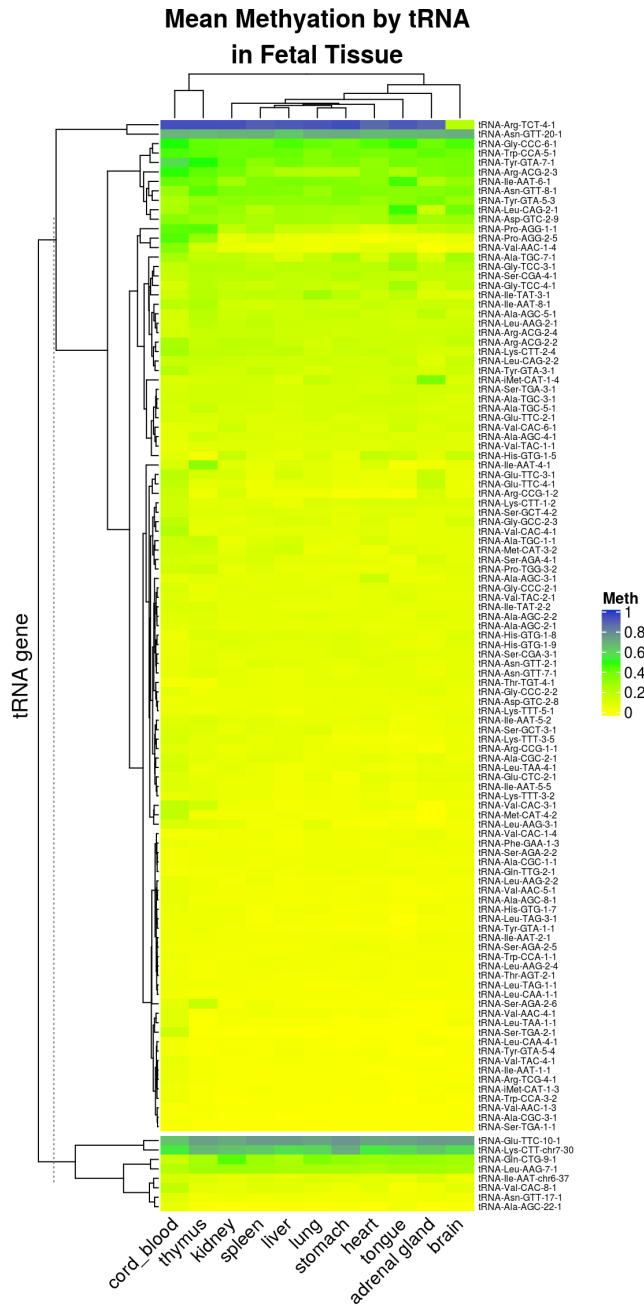


Figure 4.26: Mean Methylation of 115 tRNAs in 11 tissues. Possible pseudogenes are shown in a separate cluster beneath the main heatmap [337].

## 4.5 tRNA Gene Methylation in Cancer

A coarse grained examination of the tRNA DNA methylome in tumour samples with matched normal tissue data was performed, considering the 45 tRNA genes and 19 tissues for which data were available. Across multiple tissues it was found that both tumour and ‘normal’ samples have a strong preponderance for low levels of methylation consistent with the active transcription of

many tRNA genes.

Tumour samples have a higher mean methylation,  $0.27\%$  ( $95\%ci 2.28 \times 10^{-3} - \infty$ ) greater than that of normal samples (Wilcoxon signed-rank test p-value =  $9.1 \times 10^{-20}$ ), and a greater variance (Levene test p-value <  $1 \times 10^{-3}$ ) (Figure 4.27b). Interestingly visual inspection of methylation values in normal tissues (Figure 4.27b & d) suggests a slight bimodality peaking around 50% methylation that is not present in the tumour samples. These differences are small in absolute terms and have very wide error margins but are consistent with a picture of dysregulation in cancer cells [278].

Both Primary Tumour and Solid Tissue Normal samples show slight increasing methylation with age but age accounts for very little of the variance (linear regression slope estimates = 1.65, 1.52;  $R^2 = 0.0006, 0.0002$ ; p-values  $1.20 \times 10^{-6}, 1.34 \times 10^{-3}$  respectively)(Figure 4.27d).

Consistent with Figure 4.25 the tRNAs tRNA-iMet-CAT-1-4 and tRNA-Ser-AGA-2-6 which have the highest ranked variance in normal tissue and also rank highly in tumour samples (Figure 4.27c).

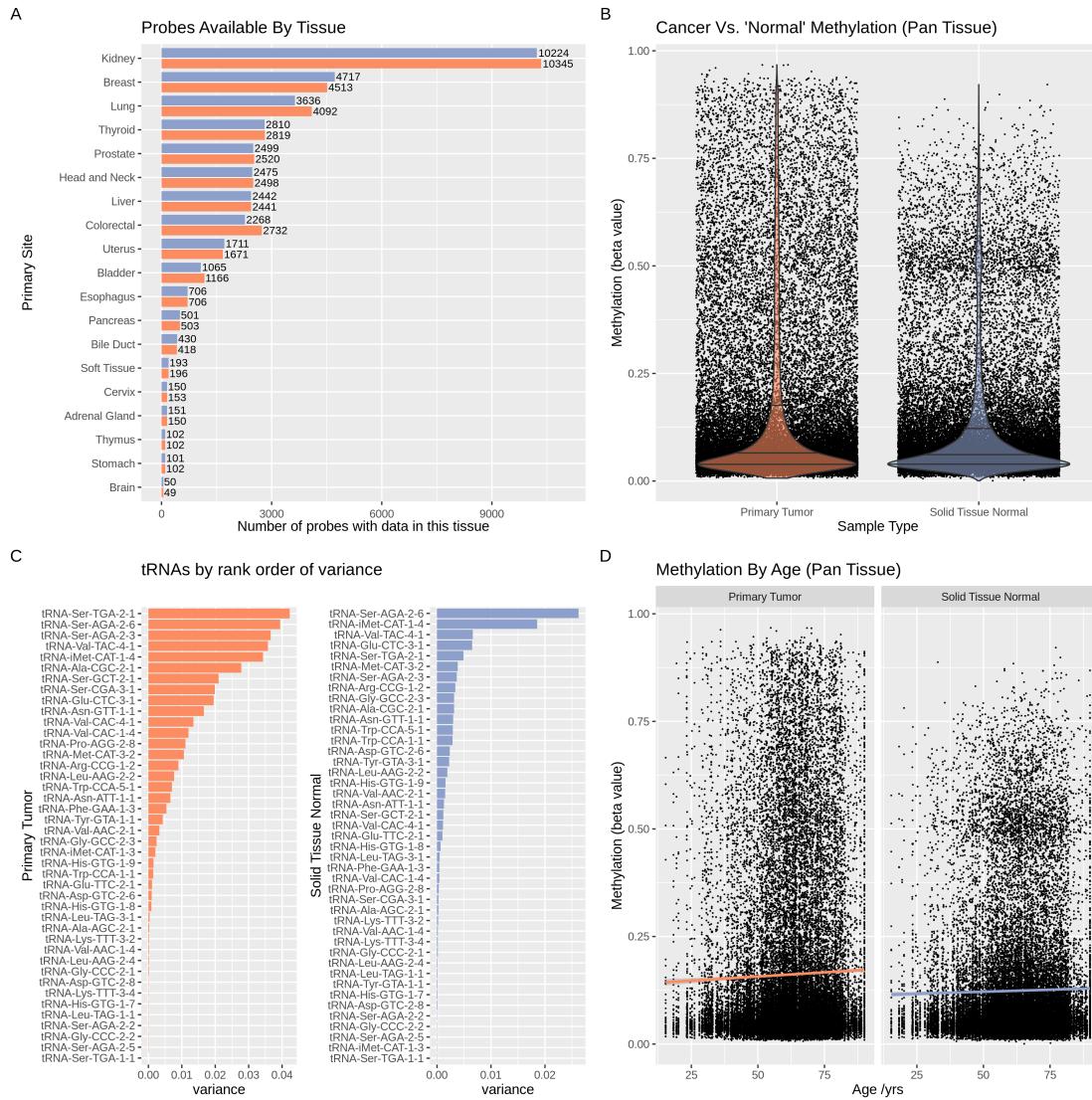


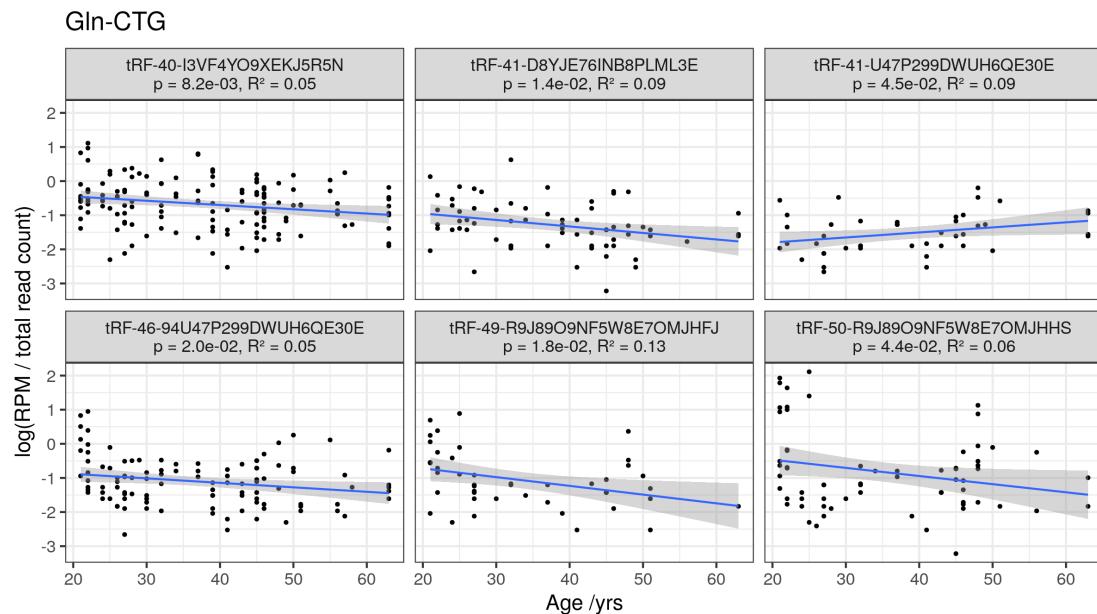
Figure 4.27: Global properties of tRNA methylation data for 45 tRNA genes across 19 tissues with matched normal and tumour samples from 733 cases in TCGA [331,332].

#### 4.5.1 Expression of tRNAs in Blood with Age

Having observed specific tRNA gene isodecoders hypermethylating with age it followed that the expression of tsRNA in blood cell-types would be worth exploring to ascertain if the expression levels of such transcripts declined as DNA methylation levels of the loci which encode them increase. A bioinformatic approach was devised to attempt to assay tRNA transcription in order to use standard publicly available small RNA-seq datasets. A customised MINTmap [296] reference was designed to include only fragments which unambiguously map to a single tRNA gene locus and which overlap the 5' or 3' end of the genomic tRNA sequence by at least one base with no mismatches. This reference is intended to capture pre-tRNAs prior to processing and CCA addition operating under the assumption that the levels of pre-tRNAs will be informative

about the amount of transcription taking place at the tRNA loci (see [Methods 4.3.7](#)). This custom MINTmap reference build yielded 383 fragments mapping to 92 distinct tRNA loci in this data. The lack of coverage of age hypermethylating tRNAs by uniquely attributable RNA-seq reads prevented us from drawing any strong conclusions about the relationship between DNAm changes and changes in tRNA transcription.

Using the original MINTmap reference optimised to detect tRNA fragments derived from mature tRNAs there were 5384 unique fragments derived from as many as 417 tRNA loci. However, the mapping between fragments and loci in this reference is many to many, with each tRNA gene able to give rise to many fragments and each fragment attributable to at least 1 and usually many tRNA genes. The examination of these fragments was limited to those with a length of greater than or equal to 40nt to capture reads more likely to be derived from mature tRNAs rather than tRFs or tRNA halves ([Figure 4.29](#)). 48 tsRNAs with nominally significant expression changes ( $p < 0.05$ ) were identified, 8 increased and 40 decreased in abundance with age. For example 5 of 6 fragments showing significant age-related expression changes derived from the Gln-CTG family of tRNAs are decreasing with age ([Figure 4.28](#)). This is suggestive that expression of some tRNA genes may decline with age but this possibility is in need of additional tRNA expression data before it can be asserted with confidence.



[Figure 4.28](#): tRNA fragments derived from the Gln-CTG family of tRNAs, selected as tRNA-Gln-CTG-7-1 is one of the 16 study-wide significant age-hypermethylating tRNA genes. Panel titles contain the MINTbase Plates, unique identifiers of the tRNA fragments [[288](#)].

#### 4.5.1.1 MINTmap reference Fragment distribution

In the original MINTmap reference (Figure 4.29b) there are peaks at around 18nt, 22nt and 32nt. This is consistent with the expected tRNA fragment size distributions with ‘tRNA halves’ at 30-33nt and other tRFs at 18nt and 22nt. In this custom reference (Figure 4.29a) whilst there is still a peak at ~18nt, with suggestions of peaks near 22nt and 32nt the tRNA fragment length distribution is somewhat different from that of the standard MINTmap reference. There are larger peaks at ~28 and ~40nt consistent with the longer fragments expected given that this reference aimed to target fragments derived from pre-tRNAs not tRFs derived from mature tRNAs.

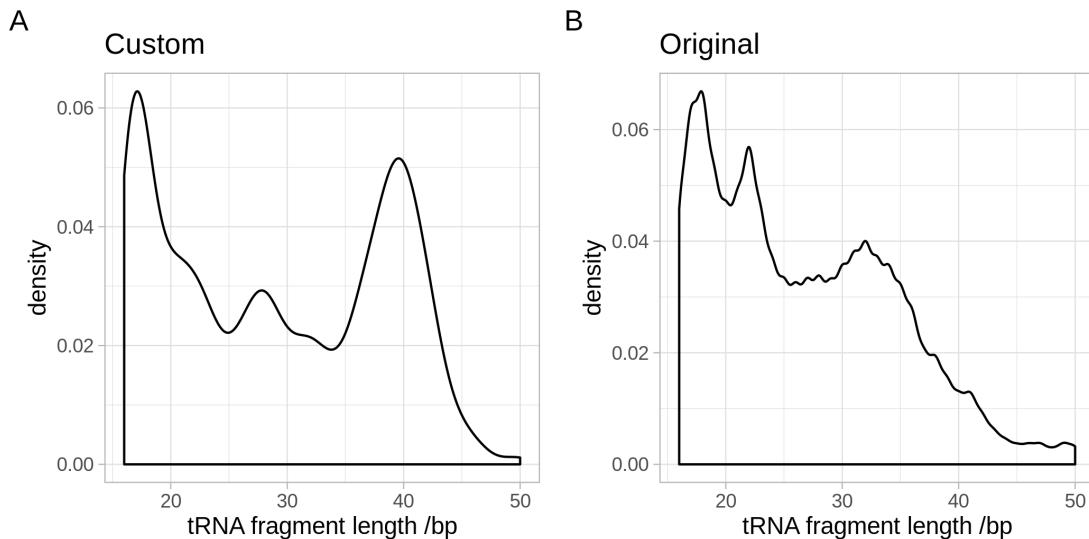


Figure 4.29: Comparison of the fragment size distributions between our custom reference **A**) and the original the MINTmap reference **B**).

#### 4.5.2 Mice also show age-related tRNA gene DNA hypermethylation

To see if this age-related tRNA gene DNA hypermethylation was present in another mammal, the DNA methylation of the mouse tRNA genes was examined using data from a reduced representation bisulfite sequencing (RRBS) experiment performed by Petkovich et al. [334]. These data from 152 mice covered 51 tRNA genes and 385 CpGs after QC (see [Methods 4.3.8](#)), representing ~11% of mouse tRNA genes. The mice ranged in age from 0.67-35 months.

Three of the 51 tRNAs showed Bonferroni significant DNA methylation changes with age ( $p\text{-value} < 1.08 \times 10^{-4}$ ) and all were in the hypermethylation direction. These three are tRNA-Asp-GTC-1-12, tRNA-Ile-AAT-1-4, tRNA-Glu-TTC-1-3 (Figure 4.30). Full Age model Results are available in [Supplementary File S4](#).

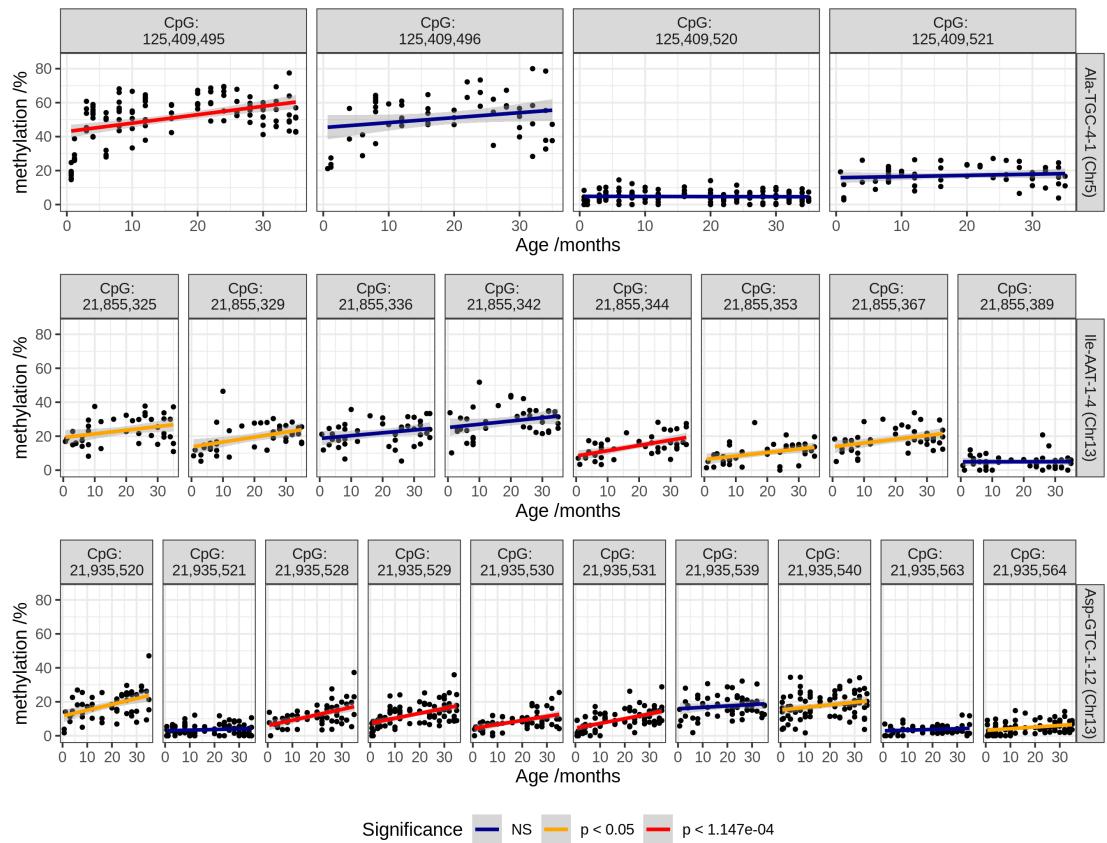


Figure 4.30: DNA methylation of CpGs in 3 tRNA which significantly hypermethylate with age in mice, in data from Petkovich et al. [334]. 6 CpGs reach Bonferroni significance and 7 show nominally significant increases.

In order to investigate the mouse results further data from Thornlow et al. [338], which had established tRNA ortholog sets for 29 mammalian species, were used. They identified 197 mouse tRNAs as having direct human orthologs with 44 of these included in the mouse results from Petkovich et al. [334]. However, unfortunately, although 2 of the top 3 tRNAs (tRNA-Ser-AGA-2-6 & tRNA-Ile-AAT-4-1) have mouse orthologs (tRNA-Ser-AGA-2-2 & tRNA-Ile-AAT-1-1), they are not covered in these mouse data. Furthermore, none of the tRNAs showing significant hypermethylation in mice (at nominal  $p < 0.05$ ) had human orthologs in our MeDIP-seq study-wide significant hypermethylating set. Therefore, whilst specific tRNA loci cannot be directly compared due to this lack of coverage, it is interesting that the small number of significant tRNA genes in the mouse data also hypermethylate with age.

## 4.6 Discussion

This work has identified a previously unknown enrichment for age-related epigenetic changes within the tRNA genes of the human genome. This observation was strongly directional with increasing DNA methylation with age [339]. The MeDIP-seq dataset employed brought advantages in exploring this undefined terrain of the tRNAome. Firstly, being genome-wide it provides much increased access, as these regions are poorly covered by current arrays. Secondly, being a fragment-based regional assessment of DNA methylation, the individual but highly similar small tRNA genes can be surrounded by unique sequence.

It was determined by genome-wide permutation that this strong hypermethylation signal was specific to the tRNAome, and not merely driven by the underlying CpG density of these loci. A targeted BiS-seq experiment validated the defined nature of the tRNA change in an independent dataset, with a pooled sample approach, which may also be useful for other ageing-related targeted DNA methylome evaluations. Additionally, support was gained for these results from limited DNA methylation array data.

Whilst the changes in DNA methylation observed here are relatively small (i.e. 3.7% between 4-year to 78-year pools in tRNA-iMet-CAT-1-4), this is consistent with the effect size seen in the majority of positions differentially methylated with age in many other studies [120,340], except for the noted extreme outliers, such as *EVOLV2* [341]. Furthermore, effect size cut-offs have distinct limitations, they do not generalise well between studies and they have been observed to remove a large fraction of true age-DMPs [342].

Additional exploration of what was driving this age-related phenomenon and its possible biological implications was undertaken. This result was observed in DNA from peripheral blood, a source with a heterogeneous cell-type population [188]. Moreover, there are well known age-related proportional changes in peripheral blood cell composition [163]. The TwinsUK MeDIP-seq and 450k array DNA methylation data included measured haematological values. Therefore, major cell type effects such as a myeloid skew were adjusted for, and distinct tRNAs were still significant. Although, a caveat to this is that one cannot exclude changes in minor specific sub-cell fraction types. However, that these age-related effects were strong enough to be observed in both a regional MeDIP-seq assessment, DNA methylation arrays with minimal coverage of the tRNA genes, and a pooled sequencing approach, implies that they not extremely subtle. Age-related tRNA gene DNA methylation changes were observed in the limited subset of mouse tRNA genes covered in publicly available RRBS data (~13%) and were able to identify tRNAs exhibiting DNA hypermethylation with age in this set. This suggests that age-related tRNA gene hypermethylation may not be unique to humans, and may extend at least across mammals.

The high number of hypermethylating tRNAs prior to cell-type correction pointed to the possibility that the epigenetic state of this small tRNAome fraction of genome could capture cell-type.

tRNA gene DNA methylation was found to be capable of separating myeloid from lymphoid lineages. There also was some suggestion of more fine-grained blood cell-type signatures in tRNA DNAm, such as the separation of CD19+ B cells from CD4/8+ T cells. Ageing is also known to lead to an increase in senescent cells (*e.g.* CD8+ CD28- cells). Whether these epigenetic changes in the tRNAome uniquely represent these cell-types will require technical advances to enable future single cell DNA methylome analysis to accurately assess these regions. If further supported, the epigenetic state of these loci may aid the taxonomy of cell-type definition.

This signal within the tRNA families was observed to occur at specific Isodecoders. After correcting for major cell types, two tRNA genes, tRNA-iMet-CAT-1-4 and tRNA-Ser-AGA-2-6, were identified which had the most consistent hypermethylation across 3 different assays. Isodecoders expand in number with organismal complexity and the high prevalence in mammals has been suggested to be due to their additional regulatory functionality [293,343]. They also have distinct translational efficiency [344], which can also have consequences in human disease [286]. Furthermore, there is great complexity to the fragmentation of tRNA [299], with physiological processes such as stress shown to induce fragment production [345]. These resultant tsRNAs can feedback on protein synthesis by regulating ribosome biogenesis [297] and others have diverse regulatory functions such as targeting transposable element transcripts [298]. They are also observed to circulate in the blood in a cell-free fashion, and fragment levels can be modulated by ageing and calorie restriction [305]. The isodecoder specific nature of our findings frame a possible hypothesis for regulatory change with age and future work will be required to unravel this potential.

There was an inconsistent result for tRNA-Ile-AAT-4-1, it is covered by a MeDIP-seq 500bp window which exhibited genome-wide significant hypermethylation, but also partially overlaps tRNA-Ser-AGA-2-6 (Figure 3B). Whilst the 450k array probe overlapping tRNA-Ile-AAT-4-1 (cg06382303) appears to replicate this hypermethylation, it is a borderline case for exclusion flagged by Zhou et al. [186] due to non-unique 30bp 3' subsequence. In the targeted Bisulfite sequencing data, tRNA-Ile-AAT-4-1 exhibited a loss of methylation. Therefore, this may suggest that the hypermethylation signal observed at this locus in the MeDIP-seq data could have been ‘pulled up’ by the neighbouring tRNA-Ser-AGA-2-6 hypermethylation signal. Of the 16 study-wide significant tRNA genes, only these two have a shared significant window, furthermore, in the expanded set of 44 only tRNA-Thr-AGT-1-2 could be similarly affected.

Codon usage frequencies have been claimed to be mostly invariant in the transcriptomes of a wide range of tissues, as well as across developmental time [336]. Although, others have found the majority of isodecoders are transcribed in different cell types [262].

However, differences in the codon usage of genes with highly tissue-specific expression patterns have been observed [346]. Transcriptome codon usage frequency and corresponding tRNA gene expression levels have been claimed to vary with several other factors; notably, the replicative

state of cells, as well as across developmental time [347]. Experimental stress-related states have revealed changes with an over-representation of codons that are translated by rare tRNAs [282]. Differences in the codon usage of genes with highly tissue-specifically expressed genes have been observed [348]. It has however been argued that these differences are substantially explained by variation in GC content [349], and the extent to which codon usage plays a role in regulation of translation remains contested. Changes in the epigenetic state of specific tRNAs could be modulating transcription efficiency or even codon availability in the ageing cell. tRNA gene dosage is quite closely matched to amino acid usage frequency in the human exome.

The location of tRNA-Ser-AGA-2-6 and tRNA-Ile-AAT-4-1 immediately downstream of *CTC1* and of tRNA-Ile-AAT-4-1 within the 3'UTR of an alternate isoform of *CTC1*, which undergoes nonsense mediated decay raises the possibility that the gene body epigenetic regulation of *CTC1* may impact on the state of these tRNA genes. *CTC1*'s function in telomere maintenance [350], DNA replication licensing [351], and it's role in a rare progeroid condition [352] indicate that it has ageing-relevant biology. The possible relationship between the regulation of *CTC1* and that of the tRNA genes downstream of it warrants further study.

Several transcription factors acting via TFIIIB [276] have a negative (the tumour suppressors p53 [274] and Rb [275]) or positive (the proto-oncogene c-Myc) influence [276]. Regulatory sequence in the flanking or the internal regions of tRNA genes do not explain tRNA expression variation [347]. Whilst DNAm is able to repress the expression of tRNA genes [279] in a plasmid expression system, the broader chromatin environment also affects tRNA transcription. Due to the co-ordinated nature of epigenomic modifications, it may also be revealing to evaluate ageing-related histone modification in these tRNA loci.

tRNA sequences themselves are under strong structural (both secondary and tertiary) [343] as well as functional constraint, which leads to an order of magnitude reduction in variation compared the background genomic mutation rate [262]. However, polymorphic tRNA could be another potential caveat to this work. Although, there is no significant population variation in, for example, tRNA iMet sequences in 1,000 Genomes data. Indeed, there are only 11 new isodecoder sequences with high confidence (tRNAscan scores  $\geq 50$ ) at  $>1\%$  population frequency [262]. Despite strong purifying selection maintaining very low variation in tRNA gene bodies, tRNA genes are subject to high levels of transcription-associated mutagenesis (TAM) leading to elevated mutation rates over evolutionary time in their immediate flanking sequences [353]. There is also some evidence for tRNA copy number variation at specific loci, although this remains under-characterised [354,355]. Another potential cause considered was whether age-related somatic copy number increases could be occurring in these loci. Population or somatic copy number expansions could lead to increased methylated reads in MeDIP-seq without any epigenetic state change. However, this would not be consistent with the targeted and array Bis conversion methodologies, where the proportion of methylated to unmethylated reads would still

be constant.

It is worth noting the parallels with known cancer and ageing epigenetic changes, and that tRNAs are also dysregulated in cancer [278], with proposed utility as prognostic markers [277]. Furthermore, the early replicating state of tRNA loci, potentially associated with high expression [356], may make them prone to hypermethylate, as is observed in early replicating loci in both cancer [357] and senescent cells [358]. Interestingly, tRNA gene loci may also play a role in local as well as large scale genome organisation [313]. tRNA gene clusters act as insulators [312] and have extensive long-range chromatin interactions with other tRNA gene loci [313]. The coordinated transcription of tRNAs at subnuclear foci and the B-box sequence elements bound by TFIIIC and not polIII may represent an organising principle for 3D-chromatin by providing spatial constraints [314]. Therefore, these tRNA epigenetic changes could contribute to the structural changes that are also observed in ageing [359].

A predominantly unmethylated state across fetal (Figure 4.26) and adult tissues (Figure 4.25) was observed at tRNA gene loci, consistent with the high rate of transcription at many tRNA gene loci. We suspect that the tRNA genes largely remain unmethylated through development and that the moderate increases in DNAm that are observed with age at these loci are being driven by changes arising primarily in older individuals. Distinct biological changes have been observed recently in aged individuals [360,361]. This would also be consistent with the lack of significant differences in the tRNA loci detected between the neonate and the 26-year-old adult in the Heyn et al. [132] data. This low baseline DNA methylation of the tRNA genes may also explain why the observed age related changes are predominantly hypermethylation. Whether this is driven by mechanisms, such as aberrant DNA methylation targeting of the tRNA loci or specific sub-celltype effects with age, will require further experimental investigation.

The attempt to estimate tRNA transcription by identifying fragments which may be derived from pre-tRNA sequences has serious limitations. It assumes that pre-tRNA levels reflect the amount of tRNA transcription, however, this may not be the case as pre-tRNAs are rapidly processed to mature tRNAs. The requirement for a read length of at least 40nt only has only a limited ability to distinguish full length tRNAs from tRNA derived small-RNAs, and the short read length of the data prevented the use of a higher threshold. Furthermore the limited number of reads unambiguously mapping to a specific tRNA locus limits the utility of this approach for inferring the effects of epigenetic changes at a specific locus on the expression at that locus. tsRNA abundance has been associated with locus specific tRNA gene expression, in some cases independent of mature tRNA levels [290]. This has important implications for the interpretation of these results given the multi-copy nature of genes like tRNA-iMet-CAT-1-4, as even if expression levels of mature iMet tRNAs are unaffected by changes in one copy's DNA methylation, these changes could still influence the levels of particular tsRNAs derived from specific tRNA loci. The implications of these changes in DNAm levels at tRNA genes for

biological ageing warrant further exploration.

In conclusion, due to the unique challenges that make the tRNAome difficult to examine it has remained epigenetically under-characterised despite its critical importance for cell function. This work directly interrogates the epigenetic DNA methylation state of the functionally important tRNAome, across the age spectrum in a range of datasets as well as methodologies and identified an enrichment for age-related DNA hypermethylation in the human tRNA genes.

# Chapter 5

## DNA methylation clocks in Alu Elements

### 5.1 Abstract

Age predictors trained using elastic net regression on MeDIP-seq data for Alu repeat elements ( $n = 774$ ) were able to predict the age of an unrelated set of ( $n = 664$ ) individuals with an  $R$  of 0.65 and a median absolute error of 8.1 years. The Alu age estimators were prone to overestimating the age of older samples and underestimating the age of younger samples. The difference between predicted age and chronological age, the Alu age acceleration was used to perform GWAS which identified some SNPs with associations with phenotypes plausibly related to biological ageing. However, Alu Age acceleration is strongly correlated with Chronological age rendering it difficult to separate genetic associations with Alu age acceleration from those with chronological age.

### 5.2 Introduction

As we age biological systems from the molecular to the macroscopic undergo functional decline. This increasingly impaired function with time results in increased risk for a wide variety of chronic disease states and indeed for co-morbidities with these conditions. Ranging from the cardiovascular, ophthalmic, and musculoskeletal through to neurological conditions and obesity, all have age as a primary risk factor [248]. Among the molecular changes that occur with age are alterations to the epigenome [13]. The epigenome here refers to: The chemical modifications to the genome and the molecules which package it. These modifications do not affect the genomic sequence but do impact the accessibility or expression of the genome in different contexts, such as different cell-types [17].

Approximately 45% of the human genome is comprised of repetitive or transposable elements [362,363]. The repression of these elements to prevent their expansion or the impact of their latent regulatory potential from effecting genomic stability and regulation is a core function of epigenetic mechanisms [364]. The most common DNA modification, 5-methylcytosine (5mC) at CpG loci, exhibits the most well documented age-related changes, and plays a central role in the repression of repetitive and transposable elements [11]. DNA methylation also classically functions to repress gene expression at promoters [365], similar to the repressive role that is plays at repetitive elements. Beyond solely repressive effects, DNA methylation in gene bodies can modulate expression [366] and influence splicing patterns [367]. Furthermore, DNA methylation levels at transcription factor binding sites (TFBS) modulate the binding affinity of transcriptions factors for these sites [65,368].

Significant decreases in bulk 5mC with age [109] and in cancer tissues [369] are among the earliest pioneering measure of global DNAm changes. Demethylation of repetitive elements was thought to be a major contributor to the global decreases in 5mC [370]. More recent studies using methods able to resolve the genomic position of DNA methylation changes have revealed that some regions of the genome show increasing DNA methylation with age. Notably, the promoters of polycomb target genes [127], bivalent chromatin domains [125] and specific tRNA genes [1].

Alu elements are the single most abundant transposon in the human genome with ~1.19 million copies and comprising ~10.7% of the total genomic sequence [371,372]. Alu elements are members of the short interspersed nuclear elements (SINEs) family of retrotransposons, are ~300bp in length, and are derived from the 7SL RNA gene [373]. The Alu elements arose from the 7SL sequence via the FAM sequences, which gave rise to truncated variants, FRAM and FLAM the fusion of which yielded the current dimeric Alu sequence [374].

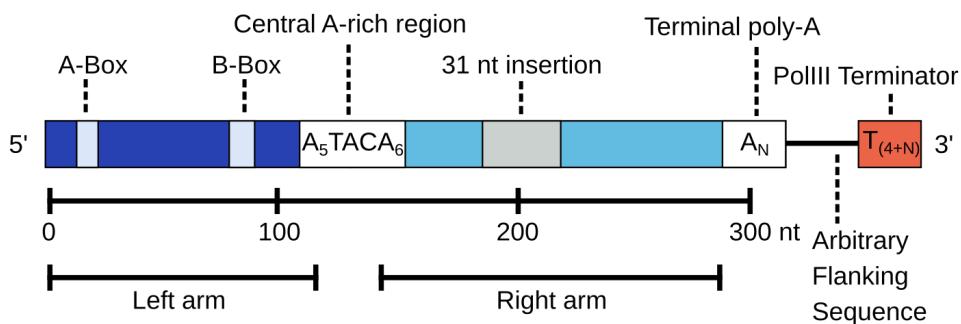


Figure 5.1: **Structure of an Alu element** A complete Alu element contains the A and B box sequence elements from an RNA polIII promoter, a central A rich region, a short insertion in the right arm, and a poly A region involved in retrotransposition and insertion. There is not an explicit polIII terminator in the Alu sequence so transcription continues through arbitrary flanking sequence until a run of at least 4 Ts, is encountered [375,376]. FLAM sequences are the origin of the Left arm and FRAMs the origin of the Right [374].

Alus are primate specific and arose ~65 million years ago (MYA), with a peak in amplification around ~40 MYA (Figure 5.2) [377]. Non-LTR retrotransposons, including Alu, L1 and SVA elements, are active in humans. However, Alu elements are non-autonomous and require the presence of an L1 element to provide the molecular machinery needed for retrotransposition [378]. Mutagenesis arising from Alu element insertion and Alu element mediated recombination is responsible for an estimated 0.4% of all human genetic disorders [363,379]

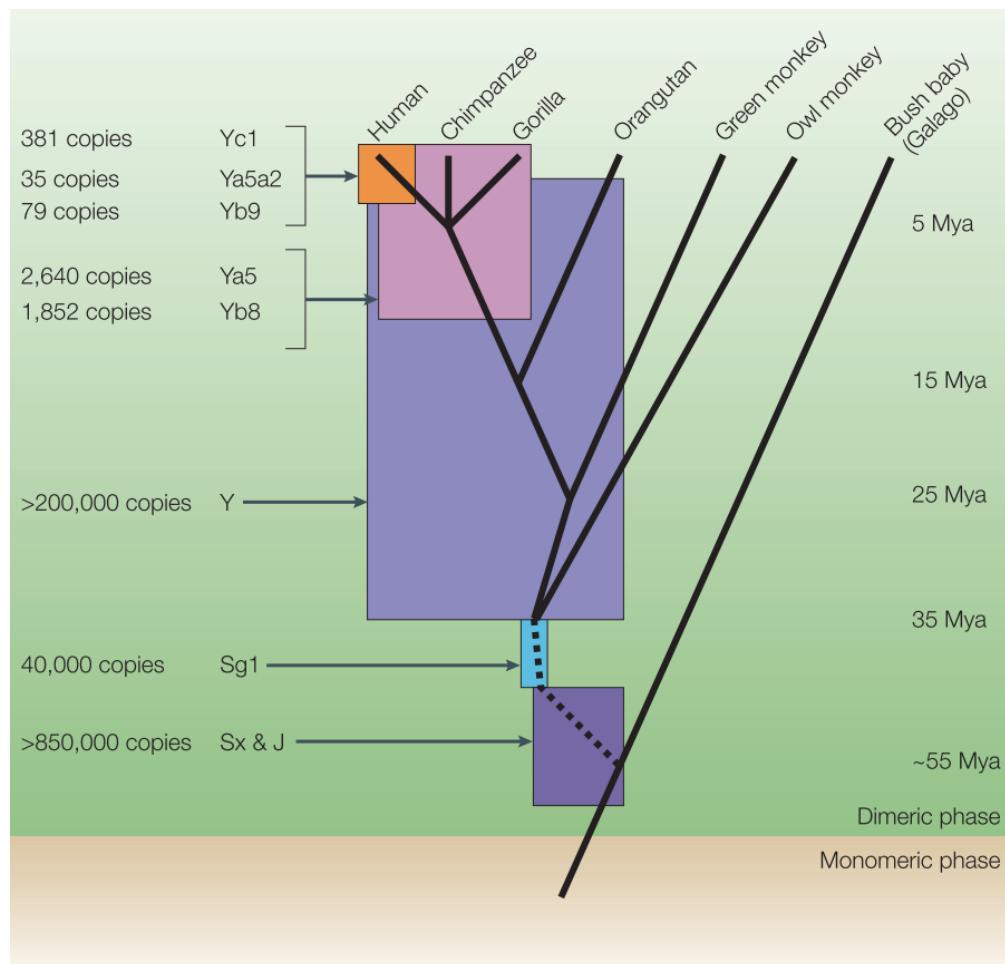


Figure 5.2: **Alu element Age and history of expansion in primates** Mya (million years ago). Reproduced from Batzer & Deininger 2002 [376].

Alu elements are CpG rich (Figure 5.3), accounting for ~7.5 million CpGs, or ~25% of all human CpG sites [372]. High levels of DNA methylation predominate at Alu elements across normal tissues. Furthermore Alu elements are more resistant to hypomethylation than other repetitive elements in cancer cells. This is indicative of strong pressure to maintain their repressed state [380].

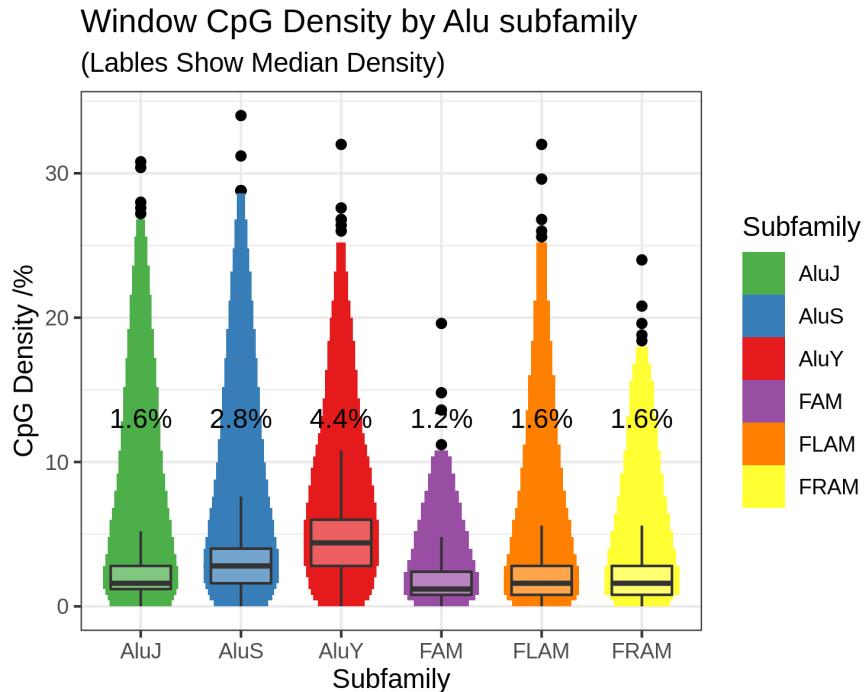


Figure 5.3: **Alu elements are CpG dense**

Alu elements contain many sequences with the potential to effect gene regulation, and indeed do effect gene expression [381,382]. There are numerous TFBSS in the canonical Alu sequence [383], and some may possess enhancer activity [384]. Thus deterioration in the effectiveness of the repression of Alu elements has the potential to disruptively expose this latent regulatory function [250,385,386]. Epigenetic reactivation of cryptic cis-regulatory elements (CREs) within transposons has been shown to act as an oncogenic driver [387], including specifically from Alu elements [388]. Alu element methylation has also previously been correlated with age-related phenotypes. Hypomethylation of Alu elements in blood cells has been associated with increasing age and with lower bone mass [389]. Alu and LINE-1 elements were more methylated in the circulating cell-free DNA of women with an older age at menopause [390]. De-repression of Alu elements also has the potential to influence common non-malignant diseases [250,391].

Quantitative assays of DNA methylation with high positional resolution which cover a large number of loci, combined with statistical methods for sparse regression analysis have permitted the construction of ‘DNA methylation clocks’ which predict chronological age with a high degree of accuracy [168,169]. The difference between chronological and predicted age known as the ‘age acceleration’ has been found to capture aspects of biological ageing in a number of these clocks. For example, Age acceleration is a strong predictor of all-cause mortality [154]. Whilst DNA hypermethylation and hypomethylation are driven by distinct mechanisms [78,392] CpGs undergoing these processes are approximately equally represented in DNA methylation clocks. Interestingly however, hypomethylating CpGs were identified by Zhang et al. to be the

most indicative of biological ageing-related all-cause mortality effects in blood [159]. Most DNA methylation clocks constructed in humans have been constructed using the Illumina bead chip array based DNA methylation assays. The coverage of these assays is strongly focused on promoters [126], CpG islands [133], and more recently enhancers [173]. However, array coverage of repetitive elements is lacking due in part to the technical challenges of accurate mapping in these spaces [169]. Thus, much of the genomic hypomethylation goes uncovered by the assays most commonly used to construct human epigenetic clocks.

In this study elastic net regression models were used to construct age predictors comprised of the DNA methylation state of Alu elements. DNA methylation data at Alu elemetns was drawn from a large ( $>3,000$ ) immunoprecipitation-derived sequencing-based dataset (MeDIP-seq) from peripheral blood DNA. Alu-specific DNA methylation ageing clocks were constructed and compared to other clocks' measures. The best performing Alu age predictor had a Pearson correlation coefficient of 0.65 and a median absolute error of 8.1 years.

## 5.3 Methods

### 5.3.1 DNA methylome data

The Methylated DNA Immunoprecipitation sequencing (MeDIP-seq) data was processed as previously described [84,170] and detailed in [Methods 2.2](#). These processed data are available from the European Genome-phenome Archive (EGA) (<https://www.ebi.ac.uk/ega>) under study number EGAS00001001910 and dataset EGAD00010000983 and were generated by BGI Shenzhen for TwinsUK.

MEDIPS (v1.0) was used for the MeDIP-seq analysis [99]. This produced reads per million base pairs (RPM) values binned into 500bp windows with a 250bp slide in the BED format, resulting in ~12.8 million windows on the genome.

### 5.3.2 Alu element annotation and pre-filtering

Alu element annotations were taken from Repeat Masker [393], only Alus on the autosomes were considered. 2,758,588 500bp windows with a 250bp slide overlapped the 1,137,653 Alu elements annotated.

Technically problematic and genetically confounded Alu element loci were removed. Alu elements located in [blacklisted](#) 834 regions of the hg19 genome, as identified by Amemiya et al. [316] were excluded as candidate loci when constructing DNAm clocks. Alu elements overlapping [known structural variants](#) identified by the 1000 genomes project [391,394] were excluded as candidate clock loci. Genetic variants which impact on the rate of biological ageing could confound DNA methylation ageing signal. Haplotype specific methylation peaks, in which a common disease associated genetic variant is predictive of DNA methylation state, were also excluded from the clock; Bell et al. identified 7,173 such regions [84]. To exclude loci with low coverage which are likely to be uninformative and noisy a zero filter was used. In quantile normalised data loci in which 0.1% or more of samples have an RPM value of 0 are excluded.

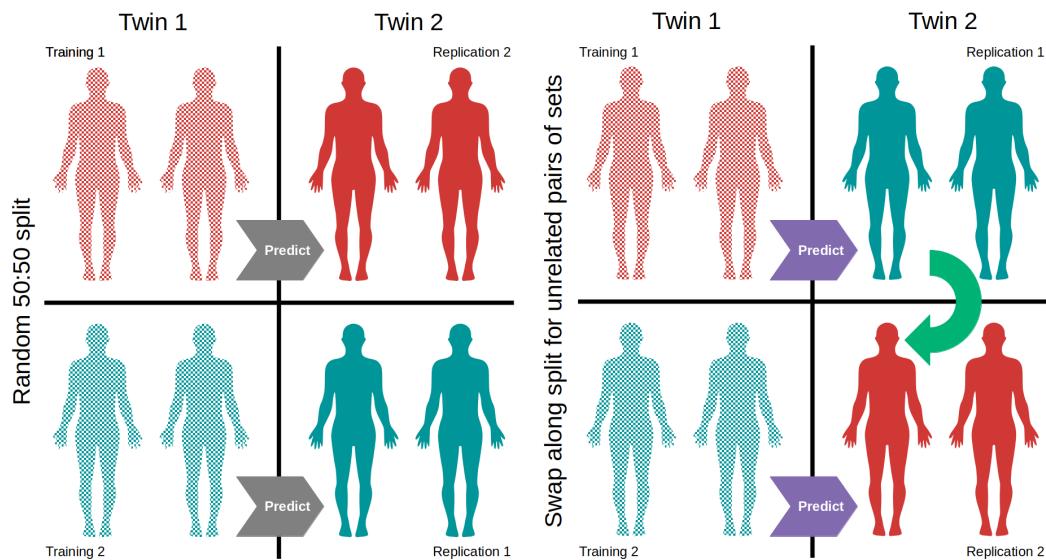
### 5.3.3 Elastic net regression

[Bigstatsr](#) was used to fit penalised regression models [395]. Elastic net regression with an  $\alpha = 0.5$ , and cross-model selection and averaging (CMSA) with 10 fold cross-validation was used. This procedure obviates the necessity of selecting a  $\lambda$  regularisation hyperparameter. CMSA generates model coefficient by averaging across coefficients generated by cross-validation [139,395,396].

### 5.3.4 Sample Selection

Two training and two replication sets were created from the twin data so as to permit the training and replication sets to be unrelated. The training sets, which were permitted to contain

singletons, contained 1548 individuals and the replication sets contained the 1308 individuals who were twins of those in the replication sets. The 1548 individuals in the training sets were randomly split into 2 groups of 774 the replication set for each training set was comprised of the twins of the other training set such that two training sets with unrelated replication sets could be constructed (Figure 5.4). Giving a total of 2856 individuals in all sets, 2617 were female and 239 male, sex was approximately equally distributed between sets with all sets being ~8% male.



**Figure 5.4: Construction of training and replication sets** Twin pairs were split one twin in each group. The first group was twins was split into two groups with a random number generator. The twins of the individuals in one random half of the group were assigned to be the replication set for the other half, producing two sets of unrelated training and replication groups. Models were also used to predict the age of related individuals to gauge the effect of genetic factors on the predictive power of the models.

The age distribution of these training and replication sets are highly similar (Figure 5.5).



Figure 5.5: **Age distribution of training and replication sets** Training 1 ( $n = 774$ , 64M/710F); Training 2 ( $n = 774$ , 65M/709F); Replication 1 ( $n = 664$ , 55M/609F), unrelated to Training 1; Replication 2 ( $n = 644$ , 55M/589F), unrelated to Training 2. M = Male, F = Female

All sets are comprised of  $\sim 2/3$  Monozygotic (MZ) twins and  $\sim 1/3$  Dizygotic (DZ) twins, exact numbers are detailed in (Figure 5.6).

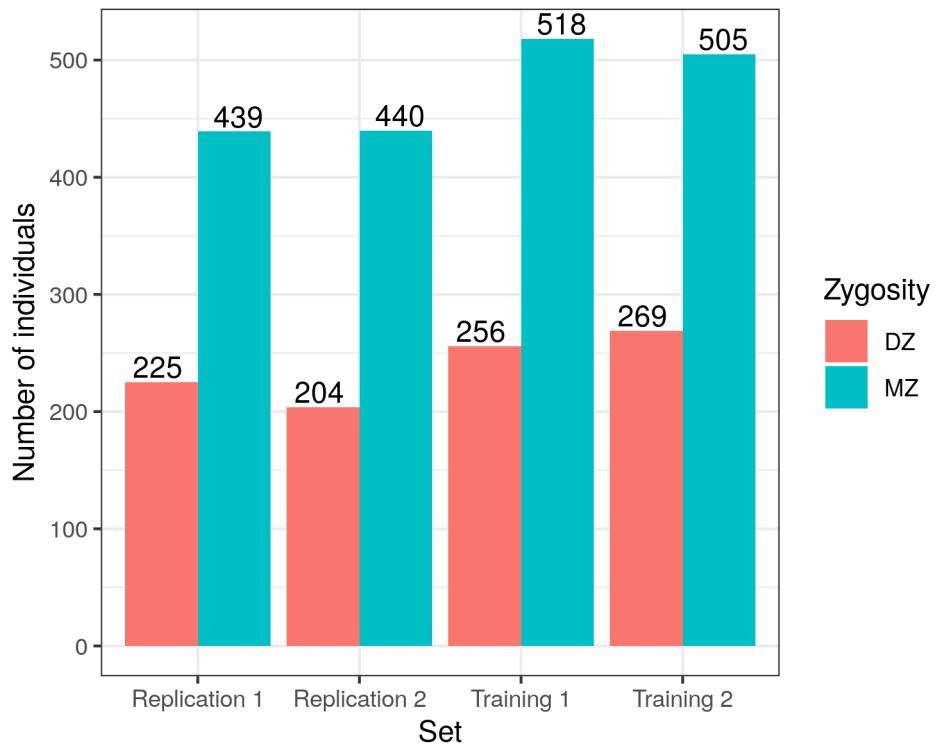


Figure 5.6: **Zygosity Make-up of sets** Replication 1 ( $n = 664$ , 225 DZ, 439 MZ); Replication 2 ( $n = 644$ , 204 DZ, 440 MZ), Training 1 ( $n = 774$ , 256 DZ, 518 MZ); Training 2 ( $n = 774$ , 269 DZ, 505 MZ).

### 5.3.5 Binarisation

Binarisation of DNA methylation data was achieved by a simple procedure: assigning windows in a given sample a DNA methylation value of 1 if their DNA methylation value was greater than the median RPM value for this window across samples and a 0 if it was less.

Binarisation of gene expression data yielded increased predictive accuracy in gene expression based ageing clocks [397]. The reasoning behind this approach is that the quantitation of changes in gene expression is noisy and binarising the data may reduce the noise whilst retaining the information important for age prediction.

### 5.3.6 Comparisons to other clocks in TwinsUK Illumina 450k array data

Illumina Infinium DNA methylation 450k arrays ((C) Illumina) were also performed on TwinsUK participants (Methods 2.1). 574 Blood-derived DNA samples had matched MeDIP-seq data not used in the training or replication sets in the Alu clock analysis. These data were available for this analysis in a pre-processed form, Methylation ‘beta’ values subject to beta-mixture quantile normalisation (BMIQ) as previously described [84,170].

DNA methylation ages were computed for the Horvath [138], Hannum et al. [137] and Levine et al. ‘PhenoAge’ [157] clocks, using the coefficients, intercept values and transformations specified in the original publications, and a model of the form:  $DNA\text{Am}Age = a\beta_{probe\ 1}...b\beta_{probe\ n} + intercept$ . Where  $a$ ,  $b$  and *intercept* are model parameters provided by the authors of the respective DNA methylation age predictors. In the absence of an intercept value 0 was used, in the case of the Hannum clock. Age acceleration values were computed as the predicted age minus the chronological age.

### 5.3.7 Genome Wide Association studies of Alu Age Acceleration

In collaboration with Pirro Hysi and Colleagues at KCL Genome wide association studies (GWAS) were performed for Alu age acceleration. Genotypic and phenotypeic data from the TwinsUK adult twin registry (see [methods 5.3](#)), with 1108 individuals being analysed in the GWAS. Genotype data was generated with a combination of Illumina HumanHap300 and HumanHap610Q chips [398]. Intensity data for each of the arrays were pooled separately. Genotypes were called with the Illuminus32 calling algorithm, thresholding on a maximum posterior probability of 0.95. Pre-phasing was performed using Eagle v2.4 [399], and imputation was performed using Minimach 3.0 [400], using haplotype information from the HRC r1.1 2016 reference panel [401]. linear mixed models were fitted to test for associations between the phenotypes and genetic variants using GEMMA v0.94 [402]. Family relatedness arising from the presence of siblings in the cohort was taken into account in constructing linear mixed models of each variant. Associations were considered suggestive at a p-value threshold of  $p < 1 \times 10^{-5}$  and significant at  $p < 5 \times 10^{-8}$ .

## 5.4 Results

A variety of different initial filtering criteria were applied to the Alu elements included in the models, as well as a number of data transformations. Pre-filtering steps were performed partially for quality control and partially to see if subsets of the Alu space with biologically interesting differences could yield clocks of equivalent accuracy but which might capture different aspects of biological ageing related to repeat de-repression.

Preliminary analysis were carried out using RPM values, Quantile normalised RPM values, Absolute methylation value estimates from MEDIPS and Binarised DNA methylation data. Quantile normalised RPM values yielded the best and most consistent results and were used in subsequent analyses.

Whilst Human epigenetic age is logarithmic with time across the lifespan the effects of this are more relevant at younger ages [403]. We fitted Alu Clock models on the natural logarithm of age in years this did not improve the quality of predicted to actual age values, likely due in part to the older age distribution of our data.

An Alu element length filter of at least 180bp or approximately two-thirds the length of a full length Alu element ( $\sim 280\text{bp}$ ) was imposed. This minimises the number of reads with the potential to misalign to other Alu elements, reads which, in this context, could lead to reduced accuracy of positionally specific Alu DNA methylation estimates [404].

CpG density filters were employed to prioritise regions likely to have dynamic and functionally interesting differences in DNA methylation. Filters were used for CpG densities above the genomic background of 2%, as well as for  $>5\%$  the approximate threshold for CpG island shores and low methylation regions [61,62].

To limit the selection of predictive sites to those likely to have consistent directions of change across samples a clock was constructed with only windows which show nominally significant changes in both members of an MZ twin pair in a consistent direction. Linear models to predict age based on DNA methylation were fitted for 687 pairs of monozygotic (MZ) twins. Batch information and blood cell counts (eosinophils, lymphocytes, monocytes, neutrophils) were included as covariates.

To produce a model that specifically captures age as a function of the de-repression of Alu elements clocks were also constructed with only windows which exhibited nominally significant ( $p < 0.05$ ) decreases in DNAm in both MZ twins. This filtering was based on the model used in selecting windows with a consistent direction of change.

An age predictor was also constructed using only AluY elements. AluY is the youngest (Figure 5.2), and most CpG dense 5.3 family of Alu elements. Having a CpG density closer to that of typically de-repressed functional elements may make such elements more readily susceptible to

stable de-repression. In addition the host genome has had less time to adapt to the presence of these elements so their de-repression may be more disruptive than older elements.

Figure 5.7 summarises the details of Alu clocks trained with these various pre-filtering criteria applied to their starting sets. The correlation with the predicted age and the chronological age of samples unrelated to the individuals on which the clock was trained indicate model quality (Figure 5.9 a). The Quality of the Alu clock predictions is quite poor in comparison to other DNA methylation clocks with the highest R of 0.65 and a median absolute error of 8.1 years, compared to 0.96 and 3.6 years for Horvath's multi-tissue clock [138].

A	Training Set	Prediction Set	R	R <sup>2</sup>	RMSE	MAE	Avail. Windows	Clock Windows	Avail. Alu Elms.	Clock Alu Elms.	Description	% Coef. -ve	Alu Type Proportions
Training 1	Replication 1		<b>0.651</b>	<b>0.424</b>	11.0	8.1	2,656,388	6,942	1,056,610	6,513	No Filters		
			<b>0.619</b>	<b>0.383</b>	11.2	8.3	1,858,147	6,616	788,953	6,196	>180bp, 99.9% non-zero		
			<b>0.609</b>	<b>0.371</b>	11.2	8.2	981,452	6,293	561,509	5,894	+ CpG Density >2%		
			<b>0.591</b>	<b>0.349</b>	11.2	8.4	1,858,147	2,131	788,953	2,066	+ ln(Age)		
			<b>0.566</b>	<b>0.321</b>	11.3	8.0	6,000	1,269	5,283	1,040	+ Consistent Direction		
			<b>0.524</b>	<b>0.275</b>	11.5	8.5	192,299	4,884	138,506	4,573	+ CpG Density >5%		
			<b>0.432</b>	<b>0.186</b>	12.2	9.0	4,289	1,187	3,708	924	+ Hypomethylation		
			<b>0.292</b>	<b>0.085</b>	12.9	9.6	92,473	1,880	62,731	1,769	+ Only AluY Elms.		
B	Training Set	Prediction Set	R	R <sup>2</sup>	RMSE	MAE	Avail. Windows	Clock Windows	Avail. Alu Elms.	Clock Alu Elms.	Description	% Coef. -ve	Alu Type Proportions
Training 2	Replication 2		<b>0.615</b>	<b>0.378</b>	12.0	8.8	2,656,388	7,091	1,056,610	6,610	No Filters		
			<b>0.582</b>	<b>0.339</b>	11.8	8.4	192,299	5,007	138,506	4,650	+ CpG Density >5%		
			<b>0.576</b>	<b>0.332</b>	12.3	8.9	1,858,147	6,435	788,953	6,021	>180bp, 99.9% non-zero		
			<b>0.572</b>	<b>0.328</b>	11.8	8.9	1,858,147	2,621	788,953	2,546	+ ln(Age)		
			<b>0.572</b>	<b>0.328</b>	12.3	8.6	981,452	5,518	561,509	5,200	+ CpG Density >2%		
			<b>0.520</b>	<b>0.271</b>	12.2	8.2	6,000	1,238	5,283	1,026	+ Consistent Direction		
			<b>0.404</b>	<b>0.163</b>	13.0	9.3	4,289	1,045	3,708	853	+ Hypomethylation		
			<b>0.290</b>	<b>0.084</b>	13.7	9.5	92,473	2,656	62,731	2,478	+ Only AluY Elms.		

Figure 5.7: **Alu Clock Model Summaries** A) Models fitted with Training 1 set predicting Replication 1 ages. B) Models fitted with Training 2 set predicting Replication 2 ages. Alu Type Composition Colour code: **AluY**, **AluS**, **AluJ**, **FAM**, **FLAM**, **FRAM**. R is Pearson's correlation coefficient. RMSE = root mean squared Error, MAE = median absolute error. '+' signifies this filter is in addition to element length >180bp and 99.9% non-zero RPM values at this locus.

Figure 5.8 shows the model summaries when correlations are computed for the related sets comprised of the twins of individuals in the training set. The correlations are mostly slightly higher for this prediction than for the prediction of unrelated individuals (Figure 5.9 b). This suggests that there are not large genetic effects on the models as their predictive efficacy is only marginally better when employed to predict the age of individuals closely related to the training set compared to unrelated individuals.

		Training Set	Prediction Set	R	R <sup>2</sup>	RMSE	MAE	Avail. Windows	Clock Windows	Avail. Alu Elms.	Clock Alu Elms.	Description	% Coef. -ve	Alu Type Proportions
Training 1	Replication 2			0.677	0.458	11.4	8.2	2,656,388	6,311	1,056,610	5,930	No Filters	1	Red
				0.644	0.415	11.6	8.4	1,858,147	6,327	788,953	5,948	>180bp, 99.9% non-zero	1	Red
				0.631	0.398	11.6	8.2	981,452	6,610	561,509	6,215	+ CpG Density >2%	1	Red
				0.600	0.360	11.6	7.9	192,299	5,052	138,506	4,707	+ CpG Density >5%	1	Red
				0.575	0.330	11.9	8.3	6,000	1,208	5,283	1,005	+ Consistent Direction	1	Red
				0.569	0.324	11.7	8.5	1,858,147	2,401	788,953	2,334	+ ln(Age)	1	Red
				0.492	0.242	12.4	9.2	4,289	1,297	3,708	1,002	+ Hypomethylation	1	Red
				0.302	0.091	13.4	9.3	92,473	1,951	62,731	1,839	+ Only AluY Elms.	1	Red
		Training Set	Prediction Set	R	R <sup>2</sup>	RMSE	MAE	Avail. Windows	Clock Windows	Avail. Alu Elms.	Clock Alu Elms.	Description	% Coef. -ve	Alu Type Proportions
Training 2	Replication 1			0.634	0.402	11.2	8.0	2,656,388	6,857	1,056,610	6,406	No Filters	1	Red
				0.586	0.343	11.2	8.2	1,858,147	2,052	788,953	1,986	+ ln(Age)	1	Red
				0.582	0.339	11.6	8.6	981,452	5,226	561,509	4,893	+ CpG Density >2%	1	Red
				0.573	0.328	11.6	8.7	1,858,147	5,907	788,953	5,545	>180bp, 99.9% non-zero	1	Red
				0.568	0.323	11.2	8.0	192,299	4,821	138,506	4,449	+ CpG Density >5%	1	Red
				0.513	0.263	11.7	8.4	6,000	933	5,283	815	+ Consistent Direction	1	Red
				0.398	0.159	12.4	8.8	92,473	3,158	62,731	2,947	+ Only AluY Elms.	1	Red
				0.329	0.104	12.6	8.8	4,289	752	3,708	641	+ Hypomethylation	1	Red

Figure 5.8: **Summaries of Alu Clock Models Predicting Twin Groups** for the purpose of comparing predictive of age models in related and unrelated samples. A) Models fitted with Training 1 set predicting Replication 2 ages. B) Models fitted with Training 2 set predicting Replication 1 ages. Alu Type Composition Colour code: **AluY**, **AluS**, **AluJ**, **FAM**, **FLAM**, **FRAM**. R is Pearson’s correlation coefficient. RMSE = root mean squared Error, MAE = median absolute error. ‘+’ signifies this filter is in addition to element length >180bp and 99.9% non-zero RPM values at this locus.

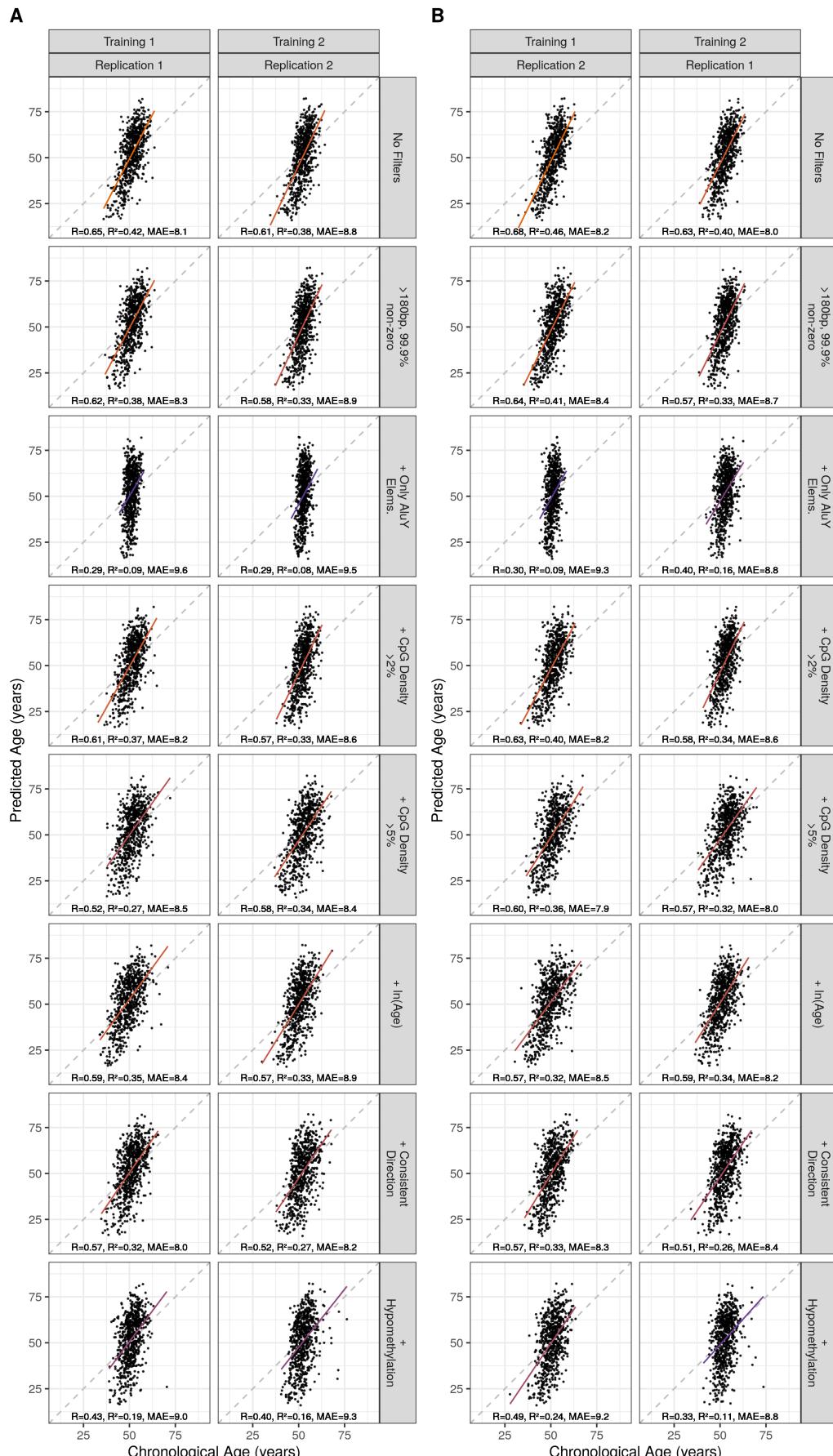


Figure 5.9: Correlation of Predicted with Chronological Age Across Models MAE = median absolute error. Orange = High  $R^2$ , purple = low  $R^2$ . '+' signifies this filter is in addition to element length  $>180\text{bp}$  and 99.9% non-zero RPM values at this locus.

### 5.4.1 Comparison to other DNA methylation clocks

To compare the Alu clock results with those of other DNA methylation age clocks DNAm age was predicted for the: Horvath [138], Hannum et al. [137] and Levine et al. ‘PhenoAge’ [157] clocks (Methods 5.3.6) using samples for which both DNA methylation array data and MeDIP-seq data was available ( $n = 574$ ). The array based clocks all perform similarly with  $R$  values in excess of 0.8, the Alu clock only attains an  $R$  of 0.46, though with a narrower error than PhenoAge. This however, is not entirely surprising given that PhenoAge aims to capture ageing phenotypes and not narrowly chronological age.

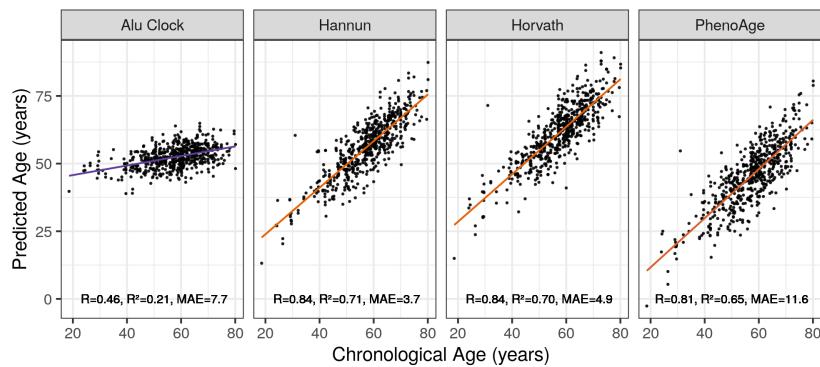


Figure 5.10: **Correlation of Alu Age Acceleration with Chronological Age Across Models** MAE = median absolute error. Orange = High  $R^2$ , purple = low  $R^2$ . Alu Clock here refers to the unfiltered model trained on the Training 1 set.

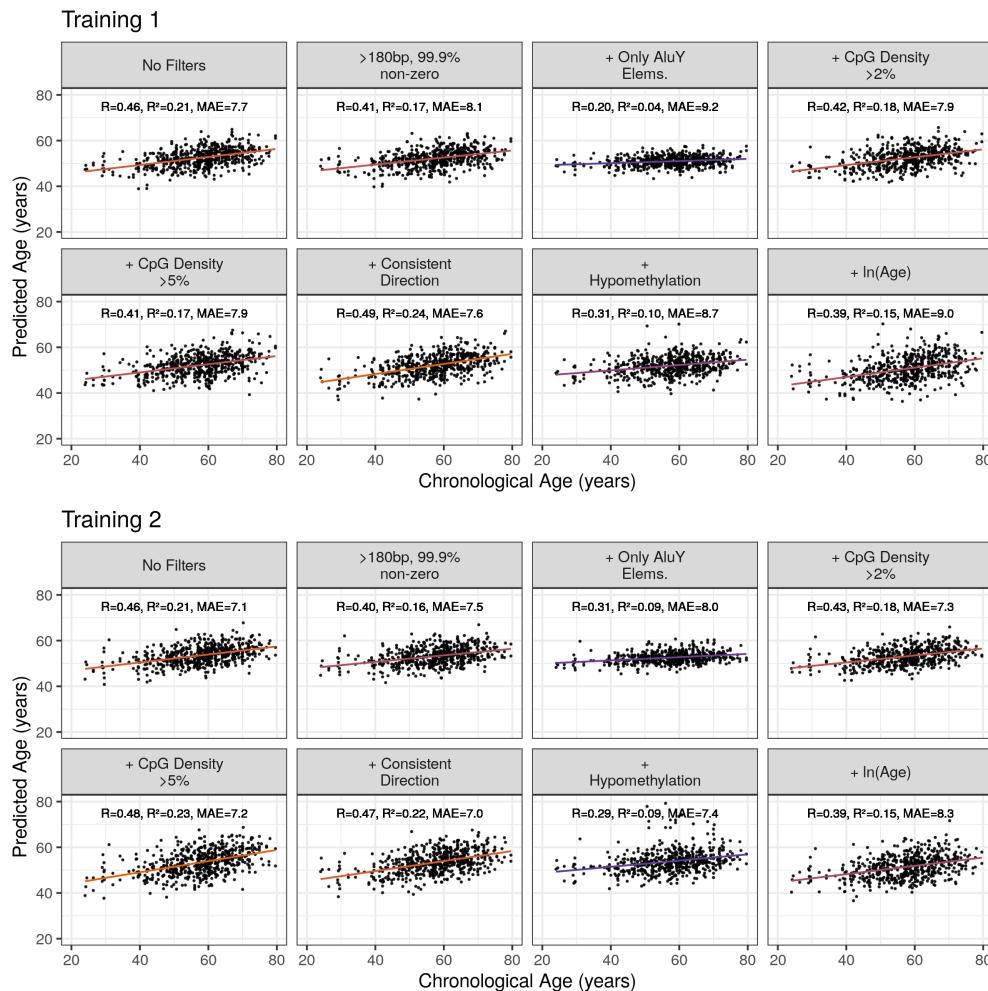


Figure 5.11: **Correlation of Alu Age Acceleration with Chronological Age Across Models** MAE = median absolute error. Orange = High  $R^2$ , purple = low  $R^2$ . Alu Clock here refers to the unfiltered model trained on the Training 1 set.

#### 5.4.2 Alu Age Acceleration

The age acceleration values (Chronological - Predicted Age) for the Alu clocks are strongly correlated with chronological age (Figure 5.12), unlike the age acceleration of the Horvath DNA methylation clock and others (Figure 5.13). The older the individual and more prone their age is to be overestimated and vice versa. This means that Alu clock age acceleration values cannot be used independently of chronological age to make statements about biological age.

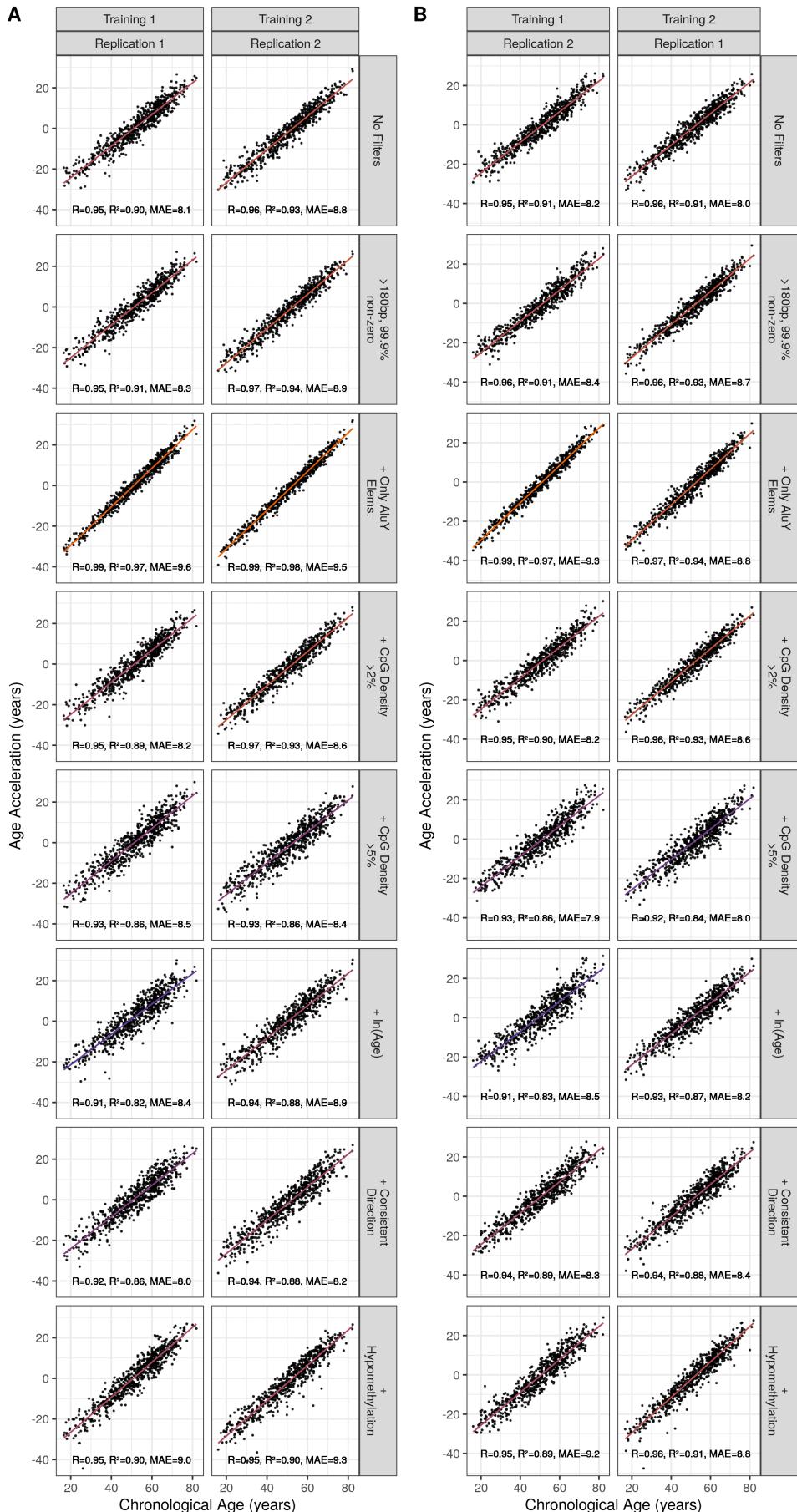


Figure 5.12: Correlation of Alu Age Acceleration with Chronological Age Across Models MAE = median absolute error. Orange = High  $R^2$ , purple = low  $R^2$ . ‘+’ signifies this filter is in addition to element length >180bp and 99.9% non-zero RPM values at this locus.

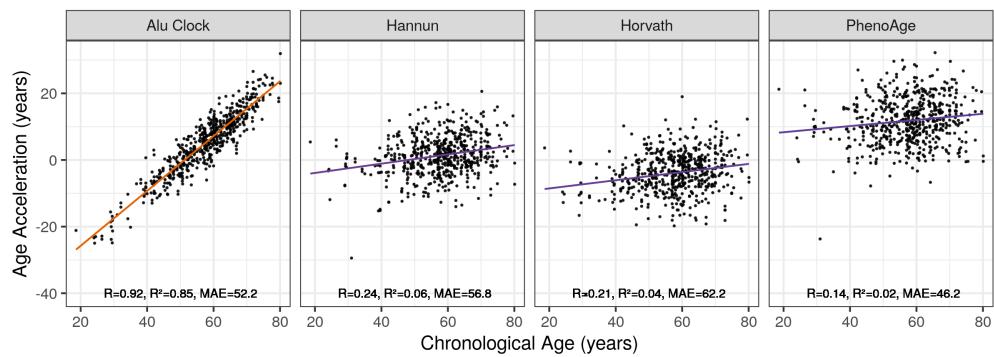


Figure 5.13: Correlation of Age Acceleration with Chronological Age Across Models  
MAE = median absolute error. Orange = High  $R^2$ , purple = low  $R^2$ . Alu Clock here refers to the unfiltered model trained on the Training 1 set. '+' signifies this filter is in addition to element length >180bp and 99.9% non-zero RPM values at this locus.

Age acceleration in the Alu clocks is not strongly correlated with age acceleration in the Horvath, Hannun or PhenoAge clocks.

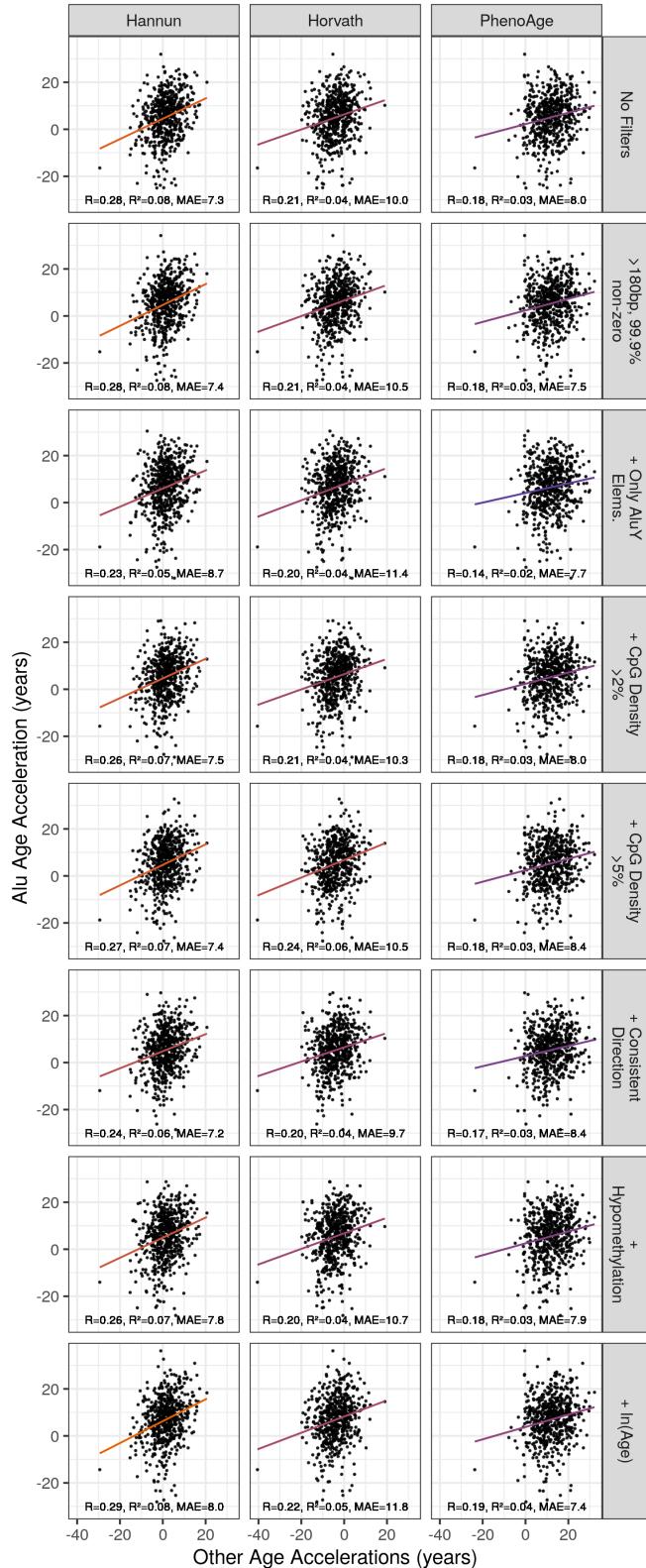


Figure 5.14: **Correlation of Alu Age Accelerations with Age Acceleration from Other Models** MAE = median absolute error. Orange = High  $R^2$ , purple = low  $R^2$ . Alu Clock here refers to the unfiltered model trained on the Training 1 set. ‘+’ signifies this filter is in addition to element length  $>180\text{bp}$  and 99.9% non-zero RPM values at this locus.

### 5.4.3 Top results Genome Wide Association studies of Alu Age Acceleration are shared by Alu Age predictors

Genome wide association studies (GWAS) for Alu age acceleration using acceleration metrics from two models were performed. The models used were the unfiltered (model 1) and the element length and percent zero filtered (model 2) models trained on training set 1. The Manhattan plots in figures 5.15 & 5.16 Illustrate that GWAS for these two models have similar results with the GWAS for model 2 attaining some lower p-values than model 1.

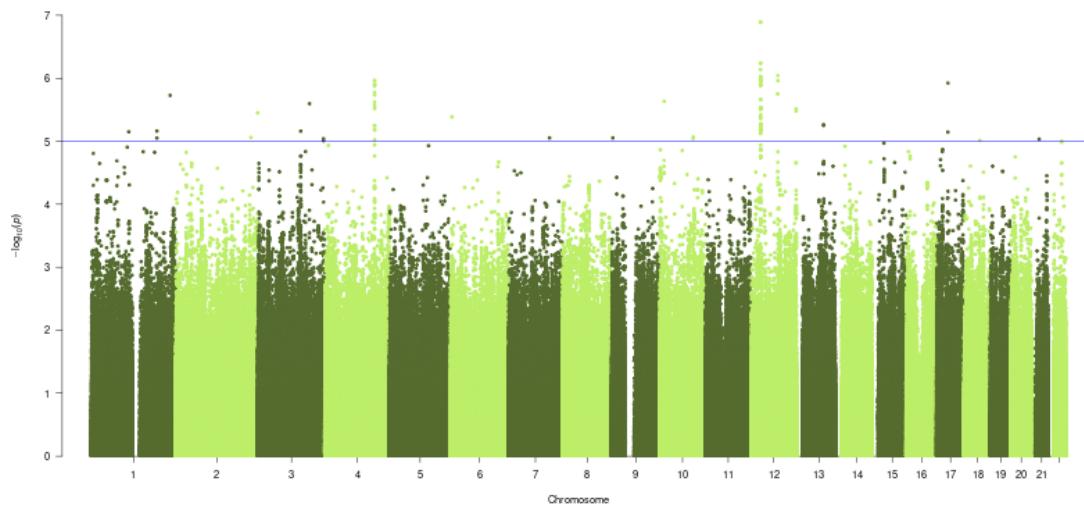


Figure 5.15: Manhattan plot showing the results of a GWAS for Alu Age acceleration computed with the predictor trained using Training set 1 on unfiltered windows.

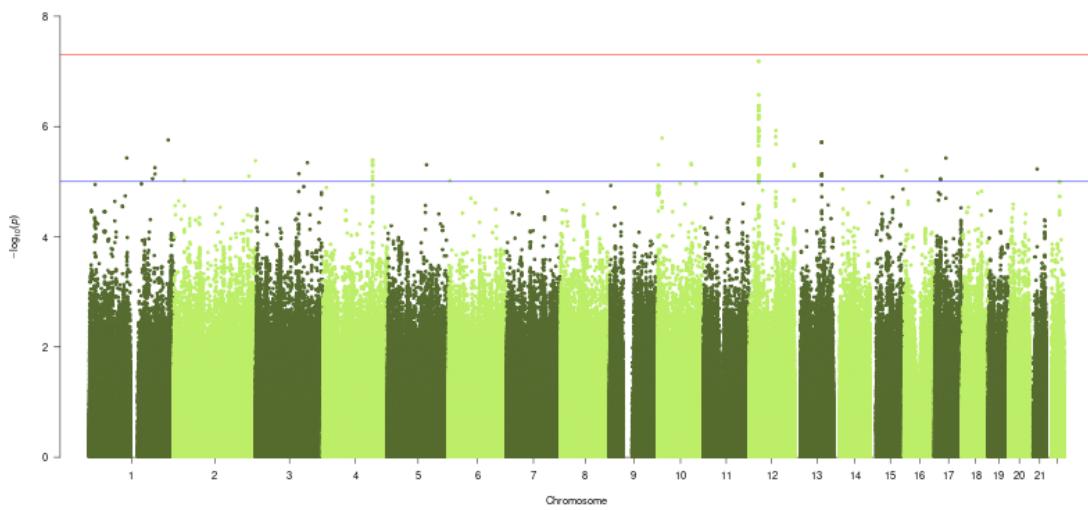


Figure 5.16: Manhattan plot showing the results of a GWAS for Alu Age acceleration computed with the predictor trained using Training set 1 on windows filtered to have an element length  $>180\text{bp}$  and 99.9% non-zero RPM values at each locus.

Neither model had any associations with a genome wide significant p-value of  $p < 5 \times 10^{-8}$ , model 1 had 115 and model 2 130 SNPs below the suggestive threshold of  $p < 1 \times 10^{-5}$ . These sites are listed in Supplementary Files S8 & S9 with the top 10 sites from each GWAS shown in figures 5.17 & 5.18 respectively.

locus	ref	alt	AF	SNP	p-value
Chr12: 28,247,840	C	T	<b>0.013</b>	rs192929352	<b>1.28e-07</b>
Chr12: 28,249,160	G	A	<b>0.013</b>	rs73080027	<b>1.29e-07</b>
Chr12: 28,275,831	C	T	<b>0.014</b>	rs12320589	<b>5.77e-07</b>
Chr12: 28,275,948	A	G	<b>0.014</b>	rs73081931	<b>5.78e-07</b>
Chr12: 28,437,969	G	A	<b>0.013</b>	rs73083941	<b>7.3e-07</b>
Chr12: 78,820,816	T	C	<b>0.526</b>	rs2060031	<b>9.1e-07</b>
Chr12: 28,329,758	G	C	<b>0.013</b>	rs73081984	<b>9.4e-07</b>
Chr12: 28,334,179	T	G	<b>0.013</b>	rs115527307	<b>9.42e-07</b>
Chr12: 28,319,636	T	C	<b>0.013</b>	rs73081974	<b>9.77e-07</b>
Chr12: 28,311,049	C	T	<b>0.013</b>	rs115986349	<b>9.95e-07</b>

Figure 5.17: Top 10 SNPs from model 1.

locus	ref	alt	AF	SNP	p-value
Chr12: 28,247,840	C	T	<b>0.013</b>	rs192929352	<b>6.55e-08</b>
Chr12: 28,249,160	G	A	<b>0.013</b>	rs73080027	<b>6.57e-08</b>
Chr12: 28,275,831	C	T	<b>0.014</b>	rs12320589	<b>2.66e-07</b>
Chr12: 28,275,948	A	G	<b>0.014</b>	rs73081931	<b>2.66e-07</b>
Chr12: 28,258,339	C	T	<b>0.016</b>	rs35792578	<b>4.18e-07</b>
Chr12: 28,249,612	C	T	<b>0.016</b>	rs11830823	<b>4.25e-07</b>
Chr12: 28,249,623	A	C	<b>0.016</b>	rs11835433	<b>4.25e-07</b>
Chr12: 28,250,701	C	T	<b>0.016</b>	rs7956083	<b>4.25e-07</b>
Chr12: 28,250,549	A	G	<b>0.016</b>	rs7955851	<b>4.25e-07</b>
Chr12: 28,249,196	A	G	<b>0.016</b>	rs11835353	<b>4.26e-07</b>

Figure 5.18: Top 10 SNPs from model 2.

The top 4 sites in both models are the same: rs192929352, rs73080027, rs12320589, & rs73081931. These sites are located in close proximity to one another on chromosome 12 (Figure 5.19).

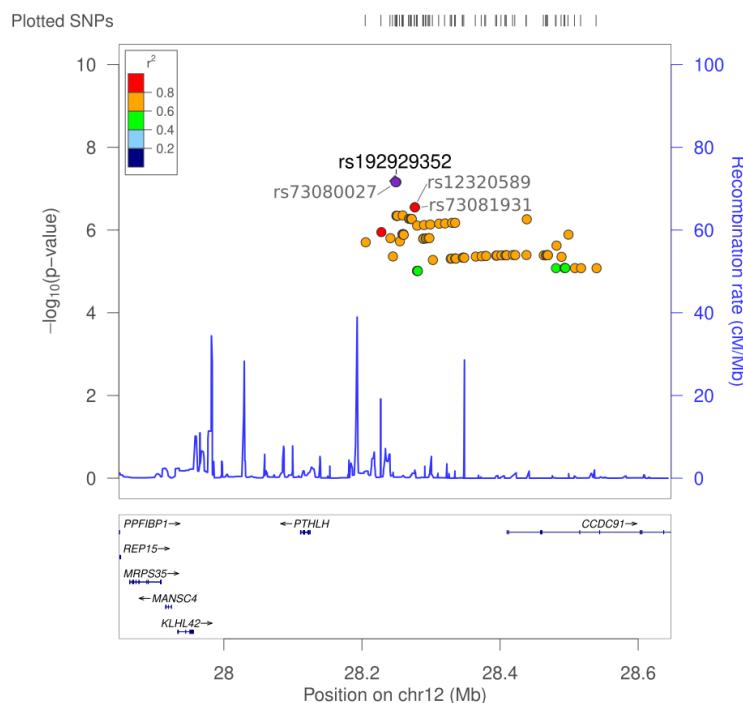


Figure 5.19: Top four shared SNPs are all located just upstream of the the CCDC91 gene. Adapted from the output of LocusZoom [405].

Further investigation of these top 4 SNPs with PhenoScanner [406,407] yielded an association with an increased risk of Treatment with ‘oestrogel’ topical estradiol ( $p\text{-value } 3.79 \times 10^{-7}$ , UK biobank,  $N=337,159$ ) for SNPs rs192929352 & rs73080027. The top phenotypes associated with these SNPs in PheWeb [408] also show some overlap (Figure 5.20)

N	Term
4	Other perinatal conditions of fetus or newborn
4	Spermatocele
3	Polyp of corpus uteri
3	Pernicious anemia
2	Acute pharyngitis
2	Sinoatrial node dysfunction (Bradycardia)
2	Pain in limb
2	Sensorineural hearing loss
2	Chronic bronchitis
2	Decreased libido
2	Malignant neoplasm of gallbladder and extrahepatic bile ducts
2	Obstructive chronic bronchitis
2	Congenital anomalies of the integument
2	Cornea replaced by transplant
2	Diabetes or abnormal glucose tolerance complicating pregnancy
1	Polyp of female genital organs
1	Acquired deformities of knee
1	Cramp of limb
1	Abscess or ulceration of vulva

Figure 5.20: *Phenotypes Common to the top 4 shared SNPs.* N = number of the top 4 SNPs in which a term appears in the top 10 Phenotypes Top four that SNP in PheWeb [408].

The risk of having many of these disease phenotypes increases with age. Risk of cancer of the gallbladder for instance increases significantly with age with more than half of new cases in the over 75s and peak incidence at 85-89 years according to cancer research UK [409]. A common feature of ageing progressive sensorineural hearing loss, or presbycusis [410]. Uterine polyps have a peak prevalence in perimenopausal women but have an increased risk of being malignant with age [411]. It is well known that incidence of cardio-vascular disease increases with age and indeed there are more specific changes in cardiac characteristics with age which could relate to sinoatrial node dysfunction. Such as reduced vagal control and altered heart rate variability [412] and, a reduced velocity long-axis systolic shortening in the left ventricle suggestive of impaired ventricular relaxation [413]. These SNP associations have biologically plausible associations with Alu age acceleration, however they are plausibly associated with chronological age and the strong correlation between Alu age acceleration and chronological age means these results cannot be interpreted as due purely to Alu age acceleration.

## 5.5 Discussion

The criteria for being called an epigenetic clock laid out by Horvath & Raj [168] are that an  $r > 0.8$  should be achieved by the clock. By this metric the Alu age predictors generated here should not be considered epigenetic clocks.

The age acceleration of the Alu age predictors is strongly correlated with age. This may be due to the relatively poor ability of MeDIP-seq to provide absolute quantitation of DNA methylation. There is substantial technical noise in the value representing DNA methylation at a given locus. The observed systematic bias towards underestimating the ages of the young and overestimating the ages of the old may be explained by a broadly consistent direction with an inconsistent magnitude. Simply scaling the value of given locus by a fixed model coefficient when there is a large amount of variation leads to outsized influence of high magnitude outliers on the age estimates. This strong correlation between Alu age acceleration and chronological age means that age is a substantial confounding factor in the GWAS for Alu age acceleration. Consequently the effect of chronological age cannot be readily disentangled from any signature of biological ageing captured by the Alu repeat based age predictors constructed here. Our sample is ~92% female which may explain the association with SNPs associated with increased likelihood of treatment with topical estradiol, and phenotypes of the female reproductive organs in the PheWeb results.

It was an aim of this work to investigate the genetic influences on the Alu clock based measures of age-acceleration, as has been shown previously with an association in the TERT gene with the Horvath clock [414,415]. Interestingly none of the SNPs identified at  $p < 1 \times 10^{-5}$  in either of the two models was present in the list of those identified in a recent GWAS conducted for the age acceleration of several other DNA methylation age clocks [416]. This could suggest that the Alu Age acceleration signature is associating with distinct genetic effects to those of other DNA methylation base age predictors. These results have the potential to reveal mechanisms involved in controlling the epigenetic state of this large portion of the genome. However the strong relationship with Age acceleration and Chronological age in the Alu clock measures of age acceleration renders this difficult as age is a potent confounder of age acceleration in this instance. The results of the genome wide association study whilst they are suggestive of some interesting age related genetic associations do not necessarily represent associations with age acceleration signatures in the Alu repeat regions of the human genome.

## **Chapter 6**

### **Discussion**

...



# References



# Appendices

## 6.1 *BioRxiv* manuscript: The Genomic Loci of Specific Human tRNA Genes Exhibit Ageing-Related DNA Hypermethylation

*Richard J. Acton<sup>1,2,3</sup>, Wei Yuan<sup>4,5</sup>, Fei Gao<sup>6</sup>, Yudong Xia<sup>6</sup>, Emma Bourne<sup>7</sup>, Eva Wozniak<sup>7</sup>, Jordana Bell<sup>4</sup>, Karen Lillycrop<sup>3</sup>, Jun Wang<sup>6</sup>, Elaine Dennison<sup>2</sup>, Nicholas Harvey<sup>2</sup>, Charles A. Mein<sup>7</sup>, Tim D. Spector<sup>4</sup>, Pirro G. Hysi<sup>4</sup>, Cyrus Cooper<sup>2</sup>, Christopher G. Bell<sup>1</sup>*

1. William Harvey Research Institute, Barts & The London School of Medicine and Dentistry, Charterhouse Square, Queen Mary University of London, London, U.K.
2. MRC Lifecourse Epidemiology Unit, University of Southampton, Southampton, U.K.
3. Human Development and Health, Institute of Developmental Sciences, University of Southampton, Southampton, U.K.
4. Department of Twin Research & Genetic Epidemiology, St Thomas Hospital, King's College London, London, U.K.
5. Institute of Cancer Research, Sutton, U.K.
6. BGI-Shenzhen, Shenzhen, China
7. Barts & The London Genome Centre, Blizard Institute, Barts & The London School of Medicine and Dentistry, Queen Mary University of London, London, U.K.

2020-09-17



HOME | ABOUT | SUBMIT | ALERTS / RSS | CHANNELS

Search

Advanced Search

New Results Comment on this paper

## The Genomic Loci of Specific Human tRNA Genes Exhibit Ageing-Related DNA Hypermethylation

Richard J. Acton, Wei Yuan, Fei Gao, Yudong Xia, Emma Bourne, Eva Wozniak, Jordana Bell, Karen Lillycrop, Jun Wang, Elaine Dennison, Nicholas Harvey, Charles A. Mein, Tim D. Spector, Pirro G. Hysi, Cyrus Cooper, Christopher G. Bell  
doi: <https://doi.org/10.1101/870352>

This article is a preprint and has not been certified by peer review [what does this mean?].

Figure 6.1: <https://doi.org/10.1101/870352>

*bioRxiv* [1]

Under Consideration at Nature Communications

1. Acton RJ, Bourne E, Bell J, Lillycrop K, Wang J, Dennison E, Harvey N, Spector TD, Cooper C, Bell CG: **The Genomic Loci of Specific Human tRNA Genes Exhibit Ageing-Related DNA Hypermethylation.** *bioRxiv* 2019;:1–4810.1101/870352 Available: <https://www.biorxiv.org/content/10.1101/870352v1>.
2. Kontis V, Bennett JE, Mathers CD, Li G, Foreman K, Ezzati M: **Future life expectancy in 35 industrialised countries: projections with a Bayesian model ensemble.** *The Lancet* 2017, **389**:1323–133510.1016/S0140-6736(16)32381-9 Available: [http://dx.doi.org/10.1016/S0140-6736\(16\)32381-9](http://dx.doi.org/10.1016/S0140-6736(16)32381-9).
3. Niccoli T, Partridge L: **Ageing as a risk factor for disease.** *Current Biology* 2012, **22**:R741–R75210.1016/j.cub.2012.07.024 Available: <http://dx.doi.org/10.1016/j.cub.2012.07.024>.
4. Belikov AV: **Age-related diseases as vicious cycles.** *Ageing Research Reviews* 2019, **49**:11–2610.1016/j.arr.2018.11.002 Available: <https://doi.org/10.1016/j.arr.2018.11.002>.

5. Fleming DM, Elliot AJ: **The impact of influenza on the health and health care utilisation of elderly people.** *Vaccine* 2005, **23**:S110.1016/j.vaccine.2005.04.018.
6. Barnett K, Mercer SW, Norbury M, Watt G, Wyke S, Guthrie B: **Epidemiology of multimorbidity and implications for health care, research, and medical education: a cross-sectional study.** *The Lancet* 2012, **380**:37–4310.1016/S0140-6736(12)60240-2Available: [http://dx.doi.org/10.1016/S0140-6736\(12\)60240-2](http://dx.doi.org/10.1016/S0140-6736(12)60240-2).
7. United Nations, Department of Economic and Social Affairs PD: **PopulationPyramid.net.** Available: [PopulationPyramid.net](http://PopulationPyramid.net).
8. Goldman D: **The Economic Promise of Delayed Aging.** *Cold Spring Harbor Perspectives in Medicine* 2016, **6**:a02507210.1101/cshperspect.a025072Available: <http://perspectivesinmedicine.cshlp.org/lookup/doi/10.1101/cshperspect.a025072>.
9. Kenyon C, Chang J, Gensch E, Rudner A, Tabtiang R: **A C. elegans mutant that lives twice as long as wild type.** *Nature* 1993, **366**:461–46410.1038/366461a0Available: <http://www.ncbi.nlm.nih.gov/pubmed/8247153>.
10. Medvedev ZA: **An Attempt at a Rational Classification of Theories of Ageing.** *Biological Reviews* 1990, **65**:375–39810.1111/j.1469-185X.1990.tb01428.xAvailable: <http://doi.wiley.com/10.1111/j.1469-185X.1990.tb01428.x>.
11. López-Otín C, Blasco MA, Partridge L, Serrano M, Kroemer G: **The hallmarks of aging.** *Cell* 2013, **153**:1194–21710.1016/j.cell.2013.05.039Available: <http://www.ncbi.nlm.nih.gov/pubmed/23746838>.
12. Avelar RA, Ortega JG, Tacutu R, Tyler EJ, Bennett D, Binetti P, Budovsky A, Chat-sirisupachai K, Johnson E, Murray A, Shields S, Tejada-Martinez D, Thornton D, Fraifeld VE, Bishop CL, Magalhães JP de: **A multidimensional systems biology analysis of cellular senescence in aging and disease.** *Genome Biology* 2020, **21**:9110.1186/s13059-020-01990-9Available: <https://genomebiology.biomedcentral.com/articles/10.1186/s13059-020-01990-9>.
13. Booth LN, Brunet A: **The Aging Epigenome.** *Molecular Cell* 2016, **62**:728–74410.1016/j.molcel.2016.05.013Available: <http://dx.doi.org/10.1016/j.molcel.2016.05.013>.
14. Sinclair DA, Oberdoerffer P: **The ageing epigenome: Damaged beyond repair?** *Ageing Research Reviews* 2009, **8**:189–19810.1016/j.arr.2009.04.004.
15. Hayano M, Yang J-H, Bonkowski MS, Amorim JA, Ross JM, Coppotelli G, Griffin P, Chew YC, Guo W, Yang X, Vera DL, Salfati EL, Das A, Thakur S, Kane AE, Mitchell SJ, Mohri Y, Nishimura EK, Schaevitz L, Garg N, Balta A-M, Rego MA, Gregory-Ksander M, Jakobs TC, Zhong L, Wakimoto H, Mostoslavsky R, Wagers AJ, Tsubota K, Bonasera SJ, et al.: **DNA Break-Induced Epigenetic Drift as a Cause of Mammalian Aging.** *bioRxiv* 2019, 10.1101/808659Available: <https://doi.org/10.1101/808659>.

16. Kane AE, Sinclair DA: **Epigenetic changes during aging and their reprogramming potential.** *Critical Reviews in Biochemistry and Molecular Biology* 2019, **54**:61–8310.1080/10409238.2019.1570075 Available: <https://doi.org/10.1080/10409238.2019.1570075>.
17. Bird A: **Perceptions of epigenetics.** *Nature* 2007, **447**:396–810.1038/nature05913 Available: <http://www.ncbi.nlm.nih.gov/pubmed/17522671>.
18. Waddington CH: **The epigenotype. 1942.** *International journal of epidemiology* 2012, **41**:10–310.1093/ije/dyr184 Available: <http://www.ncbi.nlm.nih.gov/pubmed/22186258>.
19. Russo V, Martienssen R, Riggs AD: *Epigenetic Mechanisms of Gene Regulation.* Cold Spring Harbor Laboratory Press; 1996 Available: <https://cshmonographs.org/index.php/monographs/issue/view/087969490.32>.
20. Tabansky I, Stern JNH, Pfaff DW: **Implications of Epigenetic Variability within a Cell Population for “Cell Type” Classification.** *Frontiers in Behavioral Neuroscience* 2015, **9**:1–1110.3389/fnbeh.2015.00342 Available: <http://journal.frontiersin.org/Article/10.3389/fnbeh.2015.00342/abstract>.
21. Bianconi E, Piovesan A, Facchini F, Beraudi A, Casadei R, Frabetti F, Vitale L, Pelleri MC, Tassani S, Piva F, Perez-Amadio S, Strippoli P, Canaider S: **An estimation of the number of cells in the human body.** *Annals of human biology* 2013, **40**:463–7110.3109/03014460.2013.807878 Available: <http://www.ncbi.nlm.nih.gov/pubmed/23829164>.
22. Regev A, Teichmann SA, Lander ES, Amit I, Benoist C, Birney E, Bodenmiller B, Campbell P, Carninci P, Clatworthy M, Clevers H, Deplancke B, Dunham I, Eberwine J, Eils R, Enard W, Farmer A, Fugger L, Göttgens B, Hacohen N, Haniffa M, Hemberg M, Kim S, Klenerman P, Kriegstein A, Lein E, Linnarsson S, Lundberg E, Lundeberg J, Majumder P, et al.: **The Human Cell Atlas.** *eLife* 2017, **6**:1–3010.7554/eLife.27041 Available: <https://elifesciences.org/articles/27041>.
23. CellSystemsVoices: **What Is Your Conceptual Definition of “Cell Type” in the Context of a Mature Organism?** *Cell Systems* 2017, **4**:255–25910.1016/j.cels.2017.03.006 Available: <https://linkinghub.elsevier.com/retrieve/pii/S2405471217300911>.
24. CellOntology: **Cell Ontology.** 2017, 10.5281/zenodo.168254 Available: <http://obofoundry.org/ontology/cl.html>. Accessed 5 April 2017.
25. Guo JU, Su Y, Shin JH, Shin J, Li H, Xie B, Zhong C, Hu S, Le T, Fan G, Zhu H, Chang Q, Gao Y, Ming G-l, Song H: **Distribution, recognition and regulation of non-CpG methylation in the adult mammalian brain.** *Nature neuroscience* 2014, **17**:215–2210.1038/nm.3607 Available: <http://www.ncbi.nlm.nih.gov/pubmed/24362762>.
26. Lister R, Pelizzola M, Downen RH, Hawkins RD, Hon G, Tonti-Filippini J, Nery JR, Lee L, Ye Z, Ngo Q-m, Edsall L, Antosiewicz-Bourget J, Stewart R, Ruotti V, Millar AH, Thomson

- JA, Ren B, Ecker JR: **Human DNA methylomes at base resolution show widespread epigenomic differences.** *Nature* 2009, **462**:315–2210.1038/nature08514Available: <http://www.ncbi.nlm.nih.gov/pubmed/19829295>.
27. Schultz MD, He Y, Whitaker JW, Hariharan M, Mukamel EA, Leung D, Rajagopal N, Nery JR, Urich MA, Chen H, Lin S, Lin Y, Jung I, Schmitt AD, Selvaraj S, Ren B, Sejnowski TJ, Wang W, Ecker JR: **Human body epigenome maps reveal noncanonical DNA methylation variation.** *Nature* 2015, **523**:212–610.1038/nature14465Available: <http://www.ncbi.nlm.nih.gov/pubmed/26030523>.
28. Penn NW, Suwalski R, O’Riley C, Bojanowski K, Yura R: **The presence of 5-hydroxymethylcytosine in animal deoxyribonucleic acid.** *Biochemical Journal* 1972, **126**:781–79010.1042/bj1260781Available: <https://portlandpress.com/biochemj/article/126/4/781/6593/The-presence-of-5hydroxymethylcytosine-in-animal>.
29. Kriaucionis S, Heintz N: **The Nuclear DNA Base 5-Hydroxymethylcytosine Is Present in Purkinje Neurons and the Brain.** *Science* 2009, **324**:929–93010.1126/science.1169786Available: <https://www.sciencemag.org/lookup/doi/10.1126/science.1169786>.
30. Pfeifer GP, Kadam S, Jin S-G: **5-hydroxymethylcytosine and its potential roles in development and cancer.** *Epigenetics & chromatin* 2013, **6**:1010.1186/1756-8935-6-10Available: <http://www.ncbi.nlm.nih.gov/pubmed/23634848>.
31. Ito S, Shen L, Dai Q, Wu SC, Collins LB, Swenberg JA, He C, Zhang Y: **Tet Proteins Can Convert 5-Methylcytosine to 5-Formylcytosine and 5-Carboxylcytosine.** *Science* 2011, **333**:1300–130310.1126/science.1210597Available: <https://www.sciencemag.org/lookup/doi/10.1126/science.1210597>.
32. Wu TP, Wang T, Seetin MG, Lai Y, Zhu S, Lin K, Liu Y, Byrum SD, Mackintosh SG, Zhong M, Tackett A, Wang G, Hon LS, Fang G, Swenberg Ja, Xiao AZ: **DNA methylation on N(6)-adenine in mammalian embryonic stem cells.** *Nature* 2016, **532**:329–3310.1038/nature17640Available: <http://www.ncbi.nlm.nih.gov/pubmed/27027282>.
33. Xie Q, Wu TP, Gimple RC, Li Z, Prager BC, Wu Q, Yu Y, Wang P, Wang Y, Gorkin DU, Zhang C, Dowiak AV, Lin K, Zeng C, Sui Y, Kim LKY, Miller TE, Jiang L, Lee CH, Huang Z, Fang X, Zhai K, Mack SC, Sander M, Bao S, Kerstetter-Fogle AE, Sloan AE, Xiao AZ, Rich JN: **N-methyladenine DNA Modification in Glioblastoma.** *Cell* 2018, **175**:1228–1243.e2010.1016/j.cell.2018.10.006Available: <https://linkinghub.elsevier.com/retrieve/pii/S009286741831314X>.
34. Losick JDWTABSPBAGMLR: *Molecular Biology of the Gene*. 7th ed. Pearson; 2014.
35. Bannister AJ, Kouzarides T: **Regulation of chromatin by histone modifications.** *Cell Research* 2011, **21**:381–39510.1038/cr.2011.22Available: <http://www.ncbi.nlm.nih.gov/pubmed/>

21321607.

36. Strahl BD, Allis CD: **The language of covalent histone modifications.** *Nature* 2000, **403**:41–4510.1038/47412Available: <http://www.nature.com/articles/47412>.
37. Schreiber SL, Bernstein BE: **Signaling Network Model of Chromatin.** *Cell* 2002, **111**:771–77810.1016/S0092-8674(02)01196-0Available: <https://linkinghub.elsevier.com/retrieve/pii/S0092867402011960>.
38. Voigt P, LeRoy G, Drury WJ, Zee BM, Son J, Beck DB, Young NL, Garcia BA, Reinberg D: **Asymmetrically modified nucleosomes.** *Cell* 2012, **151**:181–9310.1016/j.cell.2012.09.002Available: <http://www.ncbi.nlm.nih.gov/pubmed/23021224>.
39. Hoffman MM, Buske OJ, Wang J, Weng Z, Bilmes J a, Noble WS: **Unsupervised pattern discovery in human chromatin structure through genomic segmentation.** *Nature methods* 2012, **9**:473–610.1038/nmeth.1937Available: <http://www.ncbi.nlm.nih.gov/pubmed/22426492>.
40. Carrillo-de-Santa-Pau E, Juan D, Pancaldi V, Were F, Martin-Subero I, Rico D, Valencia A: **Automatic identification of informative regions with epigenomic changes associated to hematopoiesis.** *Nucleic Acids Research* 2017, **45**:9244–925910.1093/nar/gkx618Available: <https://academic.oup.com/nar/article/45/16/9244/3976483>.
41. Talbert PB, Ahmad K, Almouzni G, Ausió J, Berger F, Bhalla PL, Bonner WM, Cande W, Chadwick BP, Chan SWL, Cross GAM, Cui L, Dimitrov SI, Doenecke D, Eirin-López JM, Gorovsky MA, Hake SB, Hamkalo BA, Holec S, Jacobsen SE, Kamieniarz K, Khochbin S, Ladurner AG, Landsman D, Latham JA, Loppin B, Malik HS, Marzluff WF, Pehrson JR, Postberg J, et al.: **A unified phylogeny-based nomenclature for histone variants.** *Epigenetics & Chromatin* 2012, **5**:710.1186/1756-8935-5-7Available: <http://epigeneticsandchromatin.biomedcentral.com/articles/10.1186/1756-8935-5-7>.
42. Weber CM, Henikoff S: **Histone variants: dynamic punctuation in transcription.** *Genes & development* 2014, **28**:672–8210.1101/gad.238873.114Available: <http://www.ncbi.nlm.nih.gov/pubmed/24696452>.
43. Kaikkonen MU, Lam MTY, Glass CK: **Non-coding RNAs as regulators of gene expression and epigenetics.** *Cardiovascular research* 2011, **90**:430–4010.1093/cvr/cvr097Available: <http://www.ncbi.nlm.nih.gov/pubmed/21558279>.
44. Cerase A, Pintacuda G, Tattermusch A, Avner P: **Xist localization and function: new insights from multiple levels.** *Genome Biology* 2015, **16**:16610.1186/s13059-015-0733-yAvailable: <http://genomebiology.com/2015/16/1/166>.
45. Jenuwein T: **Translating the Histone Code.** *Science* 2001, **293**:1074–108010.1126/science.1063127Available: <http://www.ncbi.nlm.nih.gov/pubmed/20653993>.

46. Smith RWA, Monroe C, Bolnick DA: **Detection of Cytosine methylation in ancient DNA from five native american populations using bisulfite sequencing.** *PloS one* 2015, **10**:e012534410.1371/journal.pone.0125344Available: <http://www.ncbi.nlm.nih.gov/pubmed/26016479>.
47. Briggs AW, Stenzel U, Meyer M, Krause J, Kircher M, Pääbo S: **Removal of deaminated cytosines and detection of in vivo methylation in ancient DNA.** *Nucleic acids research* 2010, **38**:e8710.1093/nar/gkp1163Available: <http://www.ncbi.nlm.nih.gov/pubmed/20028723>.
48. Bauden M, Kristl T, Andersson R, Marko-Varga G, Ansari D: **Characterization of histone-related chemical modifications in formalin-fixed paraffin-embedded and fresh-frozen human pancreatic cancer xenografts using LC-MS/MS.** *Laboratory Investigation* 2017, **97**:279–28810.1038/labinvest.2016.134Available: <http://www.nature.com/dofinder/10.1038/labinvest.2016.134>.
49. Hashimshony T, Zhang J, Keshet I, Bustin M, Cedar H: **The role of DNA methylation in setting up chromatin structure during development.** *Nature Genetics* 2003, **34**:187–19210.1038/ng1158Available: <http://www.nature.com/articles/ng1158>.
50. Estève P-O, Chin HG, Smallwood A, Feehery GR, Gangisetty O, Karpf AR, Carey MF, Pradhan S: **Direct interaction between DNMT1 and G9a coordinates DNA and histone methylation during replication.** *Genes & development* 2006, **20**:3089–10310.1101/gad.1463706Available: <http://www.ncbi.nlm.nih.gov/pubmed/17085482>.
51. Rose NR, Klose RJ: **Understanding the relationship between DNA methylation and histone lysine methylation.** *Biochimica et Biophysica Acta (BBA) - Gene Regulatory Mechanisms* 2014, **1839**:1362–137210.1016/j.bbagr.2014.02.007Available: <http://dx.doi.org/10.1016/j.bbagr.2014.02.007>.
52. Luo Y, Lu X, Xie H: **Dynamic Alu Methylation during Normal Development, Aging, and Tumorigenesis.** *BioMed Research International* 2014, **2014**:1–1210.1155/2014/784706Available: <http://www.hindawi.com/journals/bmri/2014/784706/>.
53. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, Devon K, Dewar K, Doyle M, FitzHugh W, Funke R, Gage D, Harris K, Heaford A, Howland J, Kann L, Lehoczky J, LeVine R, McEwan P, McKernan K, Meldrim J, Mesirov JP, Miranda C, Morris W, Naylor J, Raymond C, Rosetti M, Santos R, Sheridan A, Sougnez C, et al.: **Initial sequencing and analysis of the human genome.** *Nature* 2001, **409**:860–92110.1038/35057062Available: <http://www.ncbi.nlm.nih.gov/pubmed/11237011>.
54. Duncan BK, Miller JH: **Mutagenic deamination of cytosine residues in DNA.** *Nature* 1980, **287**:560–1Available: <http://www.ncbi.nlm.nih.gov/pubmed/6999365>.
55. Saxonov S, Berg P, Brutlag DL: **A genome-wide analysis of CpG dinucleotides in**

**the human genome distinguishes two distinct classes of promoters.** *Proceedings of the National Academy of Sciences* 2006, **103**:1412–141710.1073/pnas.0510310103 Available: <http://www.pnas.org/cgi/doi/10.1073/pnas.0510310103>.

56. Bell JSK, Vertino PM: **Orphan CpG islands define a novel class of highly active enhancers.** *Epigenetics* 2017, **12**:449–46410.1080/15592294.2017.1297910 Available: <http://www.ncbi.nlm.nih.gov/pubmed/28448736>.

57. Illingworth RS, Gruenwald-Schneider U, Webb S, Kerr ARW, James KD, Turner DJ, Smith C, Harrison DJ, Andrews R, Bird AP: **Orphan CpG islands identify numerous conserved promoters in the mammalian genome.** *PLoS genetics* 2010, **6**:e100113410.1371/journal.pgen.1001134 Available: <http://www.ncbi.nlm.nih.gov/pubmed/20885785>.

58. Maunakea AK, Nagarajan RP, Bilenky M, Ballinger TJ, D’Souza C, Fouse SD, Johnson BE, Hong C, Nielsen C, Zhao Y, Turecki G, Delaney A, Varhol R, Thiessen N, Shchors K, Heine VM, Rowitch DH, Xing X, Fiore C, Schillebeeckx M, Jones SJM, Haussler D, Marra MA, Hirst M, Wang T, Costello JF: **Conserved role of intragenic DNA methylation in regulating alternative promoters.** *Nature* 2010, **466**:253–710.1038/nature09165 Available: <http://www.ncbi.nlm.nih.gov/pubmed/20613842>.

59. Ehrlich M, Gama-Sosa MA, Huang L-H, Midgett RM, Kuo KC, McCune RA, Gehrke C: **Amount and distribution of 5-methylcytosine in human DNA from different types of tissues or cells.** *Nucleic Acids Research* 1982, **10**:2709–272110.1093/nar/10.8.2709 Available: <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/10.8.2709>.

60. Bird A: **DNA methylation patterns and epigenetic memory.** *Genes & development* 2002, **16**:6–2110.1101/gad.947102 Available: <http://www.ncbi.nlm.nih.gov/pubmed/11782440>.

61. Stadler MB, Murr R, Burger L, Ivanek R, Lienert F, Schöler A, Nimwegen E van, Wirbelauer C, Oakeley EJ, Gaidatzis D, Tiwari VK, Schübeler D: **DNA-binding factors shape the mouse methylome at distal regulatory regions.** *Nature* 2011, **480**:490–510.1038/nature10716 Available: <http://www.ncbi.nlm.nih.gov/pubmed/22170606>.

62. Irizarry RA, Ladd-Acosta C, Wen B, Wu Z, Montano C, Onyango P, Cui H, Gabo K, Rongione M, Webster M, Ji H, Potash J, Sabunciyan S, Feinberg AP: **The human colon cancer methylome shows similar hypo- and hypermethylation at conserved tissue-specific CpG island shores.** *Nature genetics* 2009, **41**:178–18610.1038/ng.298 Available: <http://www.ncbi.nlm.nih.gov/pubmed/19151715>.

63. Ziller MJ, Gu H, Müller F, Donaghey J, Tsai LT-Y, Kohlbacher O, De Jager PL, Rosen ED, Bennett DA, Bernstein BE, Gnirke A, Meissner A: **Charting a dynamic DNA methylation landscape of the human genome.** *Nature* 2013, **500**:477–8110.1038/nature12433 Available: <http://www.ncbi.nlm.nih.gov/pubmed/23533700>.

<http://www.ncbi.nlm.nih.gov/pubmed/23925113>.

64. Baubec T, Schübeler D: **Genomic patterns and context specific interpretation of DNA methylation.** *Current Opinion in Genetics & Development* 2014, **25**:85–9210.1016/j.gde.2013.11.015Available: <http://www.ncbi.nlm.nih.gov/pubmed/24614011>.
65. Yin Y, Morganova E, Jolma A, Kaasinen E, Sahu B, Khund-Sayeed S, Das PK, Kivioja T, Dave K, Zhong F, Nitta KR, Taipale M, Popov A, Ginno PA, Domcke S, Yan J, Schübeler D, Vinson C, Taipale J: **Impact of cytosine methylation on DNA binding specificities of human transcription factors.** *Science (New York, N.Y.)* 2017, **356**:eaaj223910.1126/science.aaaj2239Available: <http://www.ncbi.nlm.nih.gov/pubmed/28473536>.
66. Zuo Z, Roy B, Chang YK, Granas D, Stormo GD: **Measuring quantitative effects of methylation on transcription factor–DNA binding affinity.** *Science Advances* 2017, **3**:eaao179910.1126/sciadv.aaoo1799Available: <http://advances.sciencemag.org/lookup/doi/10.1126/sciadv.aaoo1799>.
67. Kelly TK, Liu Y, Lay FD, Liang G, Berman BP, Jones P a: **Genome-wide mapping of nucleosome positioning and DNA methylation within individual DNA molecules.** *Genome Research* 2012, **22**:2497–250610.1101/gr.143008.112Available: <http://www.ncbi.nlm.nih.gov/pubmed/22960375>.
68. Huff JT, Zilberman D: **Dnmt1-Independent CG Methylation Contributes to Nucleosome Positioning in Diverse Eukaryotes.** *Cell* 2014, **156**:1286–129710.1016/j.cell.2014.01.029Available: <http://www.ncbi.nlm.nih.gov/pubmed/24630728>.
69. Raiber E-a, Portella G, Martínez Cuesta S, Hardisty R, Murat P, Li Z, Iurlaro M, Dean W, Spindel J, Beraldí D, Liu Z, Dawson MA, Reik W, Balasubramanian S: **5-Formylcytosine organizes nucleosomes and forms Schiff base interactions with histones in mouse embryonic stem cells.** *Nature Chemistry* 2018, **10**:1258–126610.1038/s41557-018-0149-xAvailable: <http://www.nature.com/articles/s41557-018-0149-x>.
70. Lyko F: **The DNA methyltransferase family: a versatile toolkit for epigenetic regulation.** *Nature Reviews Genetics* 2017, **19**:81–9210.1038/nrg.2017.80Available: <http://www.nature.com/doifinder/10.1038/nrg.2017.80>.
71. Vertino PM, Sekowski JA, Coll JM, Applegren N, Han S, Hickey RJ, Malkas LH: **DNMT1 is a component of a multiprotein DNA replication complex.** *Cell cycle (Georgetown, Tex.)* 2002, **1**:416–2310.4161/cc.1.6.270Available: <http://www.ncbi.nlm.nih.gov/pubmed/12548018>.
72. Bostick M, Kim JK, Esteve P-O, Clark A, Pradhan S, Jacobsen SE: **UHRF1 Plays a Role in Maintaining DNA Methylation in Mammalian Cells.** *Science* 2007, **317**:1760–176410.1126/science.1147939Available: <http://www.sciencemag.org/cgi/doi/10.1126/science.1147939>.

147939.

73. Jia D, Jurkowska RZ, Zhang X, Jeltsch A, Cheng X: **Structure of Dnmt3a bound to Dnmt3L suggests a model for de novo DNA methylation.** *Nature* 2007, **449**:248–25110.1038/nature06146 Available: <http://www.nature.com/articles/nature06146>.
74. Kaiser S, Jurkowski TP, Kellner S, Schneider D, Jeltsch A, Helm M: **The RNA methyl-transferase Dnmt2 methylates DNA in the structural context of a tRNA.** *RNA Biology* 2017, **14**:1241–125110.1080/15476286.2016.1236170 Available: <https://www.tandfonline.com/doi/full/10.1080/15476286.2016.1236170>.
75. Tahiliani M, Koh KP, Shen Y, Pastor W a, Bandukwala H, Brudno Y, Agarwal S, Iyer LM, Liu DR, Aravind L, Rao A: **Conversion of 5-Methylcytosine to 5-Hydroxymethylcytosine in Mammalian DNA by MLL Partner TET1.** *Science* 2009, **324**:930–93510.1126/science.1170116 Available: <http://www.sciencemag.org/cgi/doi/10.1126/science.1170116>.
76. Kandoth C, McLellan MD, Vandin F, Ye K, Niu B, Lu C, Xie M, Zhang Q, McMichael JF, Wyczalkowski MA, Leiserson MDM, Miller CA, Welch JS, Walter MJ, Wendt MC, Ley TJ, Wilson RK, Raphael BJ, Ding L: **Mutational landscape and significance across 12 major cancer types.** *Nature* 2013, **502**:333–33910.1038/nature12634.
77. Ko M, Huang Y, Jankowska AM, Pape UJ, Tahiliani M, Bandukwala HS, An J, Lamperti ED, Koh KP, Ganetzky R, Liu XS, Aravind L, Agarwal S, Maciejewski JP, Rao A: **Impaired hydroxylation of 5-methylcytosine in myeloid cancers with mutant TET2.** *Nature* 2010, **468**:839–84310.1038/nature09586 Available: <http://www.nature.com/articles/nature09586>.
78. Wu X, Zhang Y: **TET-mediated active DNA demethylation: mechanism, function and beyond.** *Nature Reviews Genetics* 2017, **18**:517–53410.1038/nrg.2017.33 Available: <http://www.nature.com/doifinder/10.1038/nrg.2017.33>.
79. Lienert F, Wirbelauer C, Som I, Dean A, Mohn F, Schübeler D: **Identification of genetic elements that autonomously determine DNA methylation states.** *Nature genetics* 2011, **43**:1091–710.1038/ng.946 Available: <http://www.ncbi.nlm.nih.gov/pubmed/21964573>.
80. Boks MP, Derkx EM, Weisenberger DJ, Strengman E, Janson E, Sommer IE, Kahn RS, Ophoff RA: **The relationship of DNA methylation with age, gender and genotype in twins and healthy controls.** *PloS one* 2009, **4**:e676710.1371/journal.pone.0006767 Available: <http://www.ncbi.nlm.nih.gov/pubmed/19774229>.
81. Kerkel K, Spadola A, Yuan E, Kosek J, Jiang L, Hod E, Li K, Murty VV, Schupf N, Vilain E, Morris M, Haghghi F, Tycko B: **Genomic surveys by methylation-sensitive SNP analysis identify sequence-dependent allele-specific DNA methylation.** *Nature genetics* 2008, **40**:904–810.1038/ng.174 Available: <http://www.ncbi.nlm.nih.gov/pubmed/18568024>.

82. Schalkwyk LC, Meaburn EL, Smith R, Dempster EL, Jeffries AR, Davies MN, Plomin R, Mill J: **Allelic skewing of DNA methylation is widespread across the genome.** *American journal of human genetics* 2010, **86**:196–21210.1016/j.ajhg.2010.01.014Available: <http://www.ncbi.nlm.nih.gov/pubmed/20159110>.
83. Shoemaker R, Deng J, Wang W, Zhang K: **Allele-specific methylation is prevalent and is contributed by CpG-SNPs in the human genome.** *Genome research* 2010, **20**:883–910.1101/gr.104695.109Available: <http://www.ncbi.nlm.nih.gov/pubmed/20418490>.
84. Bell CG, Gao F, Yuan W, Roos L, Acton RJ, Xia Y, Bell J, Ward K, Mangino M, Hysi PG, Wang J, Spector TD: **Obligatory and facilitative allelic variation in the DNA methylome within common disease-associated loci.** *Nature Communications* 2018, **9**:810.1038/s41467-017-01586-1Available: <http://www.nature.com/articles/s41467-017-01586-1>.
85. Martin-Trujillo A, Vidal E, Monteagudo-Sánchez A, Sanchez-Delgado M, Moran S, Hernandez Mora JR, Heyn H, Guitart M, Esteller M, Monk D: **Copy number rather than epigenetic alterations are the major dictator of imprinted methylation in tumors.** *Nature Communications* 2017, **8**:46710.1038/s41467-017-00639-9Available: <http://www.ncbi.nlm.nih.gov/pubmed/28883545>.
86. Ciernia AV, LaSalle J: **The landscape of DNA methylation amid a perfect storm of autism aetiologies.** *Nature Reviews Neuroscience* 2016, **17**:411–42310.1038/nrn.2016.41Available: <http://dx.doi.org/10.1038/nrn.2016.41>.
87. Ushijima T, Watanabe N, Shimizu K, Miyamoto K, Sugimura T, Kaneda A: **Decreased fidelity in replicating CpG methylation patterns in cancer cells.** *Cancer research* 2005, **65**:11–7Available: <http://www.ncbi.nlm.nih.gov/pubmed/15665274>.
88. Kunkel TA: **DNA Replication Fidelity.** *Journal of Biological Chemistry* 2004, **279**:16895–1689810.1074/jbc.R400006200Available: <http://www.jbc.org/lookup/doi/10.1074/jbc.R400006200>.
89. Laird CD, Pleasant ND, Clark AD, Sneeden JL, Hassan KMA, Manley NC, Vary JC, Morgan T, Hansen RS, Stoger R: **Hairpin-bisulfite PCR: Assessing epigenetic methylation patterns on complementary strands of individual DNA molecules.** *Proceedings of the National Academy of Sciences* 2004, **101**:204–20910.1073/pnas.2536758100Available: <http://www.pnas.org/cgi/doi/10.1073/pnas.2536758100>.
90. Pfeifer GP, Steigerwald SD, Hansen RS, Gartler SM, Riggs AD: **Polymerase chain reaction-aided genomic sequencing of an X chromosome-linked CpG island: methylation patterns suggest clonal inheritance, CpG site autonomy, and an explanation of activity state stability.** *Proceedings of the National Academy of Sciences* 1990, **87**:8252–825610.1073/pnas.87.21.8252Available: <http://www.pnas.org/cgi/doi/10.1073/pnas.87.21.8252>.

91. Riggs AD, Xiong Z: **Methylation and epigenetic fidelity.** *Proceedings of the National Academy of Sciences* 2004, **101**:4–510.1073/pnas.0307781100Available: <http://www.ncbi.nlm.nih.gov/pubmed/14695893>.
92. Jenkinson G, Pujadas E, Goutsias J, Feinberg AP: **Potential energy landscapes identify the information-theoretic nature of the epigenome.** *Nature genetics* 2017, **49**:719–72910.1038/ng.3811Available: <http://www.ncbi.nlm.nih.gov/pubmed/28346445>.
93. Zhao L, Sun M-a, Li Z, Bai X, Yu M, Wang M, Liang L, Shao X, Arnovitz S, Wang Q, He C, Lu X, Chen J, Xie H: **The dynamics of DNA methylation fidelity during mouse embryonic stem cell self-renewal and differentiation.** *Genome Research* 2014, **24**:1296–130710.1101/gr.163147.113Available: <http://genome.cshlp.org/lookup/doi/10.1101/gr.163147.113>.
94. Zhou W, Dinh HQ, Ramjan Z, Weisenberger DJ, Nicolet CM, Shen H, Laird PW, Berman BP: **DNA methylation loss in late-replicating domains is linked to mitotic cell division.** *Nature Genetics* 2018, **50**:591–60210.1038/s41588-018-0073-4Available: <http://dx.doi.org/10.1038/s41588-018-0073-4>.
95. Hatada I, Hayashizaki Y, Hirotsume S, Komatsubara H, Mukai T: **A genomic scanning method for higher organisms using restriction sites as landmarks.** *Proceedings of the National Academy of Sciences of the United States of America* 1991, **88**:9523–710.1073/pnas.88.21.9523Available: <http://www.ncbi.nlm.nih.gov/pubmed/1946366>.
96. 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.** *Nature Genetics* 2005, **37**:853–86210.1038/ng1598Available: <http://www.ncbi.nlm.nih.gov/pubmed/16007088>.
97. Harrison A, Parle-McDermott A: **DNA Methylation: A Timeline of Methods and Applications.** *Frontiers in Genetics* 2011, **2**:632–64910.3389/fgene.2011.00074Available: <http://journal.frontiersin.org/article/10.3389/fgene.2011.00074/abstract>.
98. Down TA, Rakyan VK, Turner DJ, Fliceck P, Li H, Kulesha E, Gräf S, Johnson N, Herrero J, Tomazou EM, Thorne NP, Bäckdahl L, Herberth M, Howe KL, Jackson DK, Miretti MM, Marioni JC, Birney E, Hubbard TJP, Durbin R, Tavaré S, Beck S: **A Bayesian deconvolution strategy for immunoprecipitation-based DNA methylome analysis.** *Nature Biotechnology* 2008, **26**:779–78510.1038/nbt1414Available: <http://www.nature.com/articles/nbt1414>.
99. Lienhard M, Grimm C, Morkel M, Herwig R, Chavez L: **MEDIPS: genome-wide differential coverage analysis of sequencing data derived from DNA enrichment experiments.** *Bioinformatics* 2014, **30**:284–28610.1093/bioinformatics/btt650Available: <http://www.ncbi.nlm.nih.gov/pubmed/24227674>.

100. Serre D, Lee BH, Ting AH: **MBD-isolated Genome Sequencing provides a high-throughput and comprehensive survey of DNA methylation in the human genome.** *Nucleic Acids Research* 2010, **38**:391–39910.1093/nar/gkp992Available: <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkp992>.
101. Li J, Yang J, Zhou P, Le Y, Zhou C, Wang S, Xu D, Lin H-K, Gong Z: **Circular RNAs in cancer: novel insights into origins, properties, functions and implications.** *American journal of cancer research* 2015, **5**:472–80Available: <http://www.ncbi.nlm.nih.gov/pubmed/25973291>.
102. Bock C, Tomazou EM, Brinkman AB, Müller F, Simmer F, Gu H, Jäger N, Gnirke A, Stunnenberg HG, Meissner A: **Quantitative comparison of genome-wide DNA methylation mapping technologies.** *Nature Biotechnology* 2010, **28**:1106–111410.1038/nbt.1681Available: <http://www.ncbi.nlm.nih.gov/pubmed/20852634>.
103. Hayatsu H, Wataya Y, Kai K, Iida S: **Reaction of sodium bisulfite with uracil, cytosine, and their derivatives.** *Biochemistry* 1970, **9**:2858–286510.1021/bi00816a016Available: <http://pubs.acs.org/doi/abs/10.1021/bi00816a016>.
104. Adamowicz M, Maratou K, Aitman TJ: **Multiplexed DNA Methylation Analysis of Target Regions Using Microfluidics (Fluidigm).** In *Methods in molecular biology (clifton, n.j.)* Vol. 1708 2018:349–363. Available: <http://www.ncbi.nlm.nih.gov/pubmed/29224153>.
105. Skvortsova K, Zotenko E, Luu P-L, Gould CM, Nair SS, Clark SJ, Stirzaker C: **Comprehensive evaluation of genome-wide 5-hydroxymethylcytosine profiling approaches in human DNA.** *Epigenetics & Chromatin* 2017, **10**:1610.1186/s13072-017-0123-7Available: <http://epigeneticsandchromatin.biomedcentral.com/articles/10.1186/s13072-017-0123-7>.
106. Williams L, Bei Y, Church HE, Dai N, Dimalanta ET, Ettwiller LM, Evans TC, Langhorst BW, Borgaro JG, Guan S, Marks K, Menin JF, Nichols NM, Chaithanya Ponnaluri VK, Saleh L, Samaranayake M, Sexton BS, Sun Z, Tamanaha E, Vaisvila R, Yigit E, Davis TB: **Enzymatic Methyl-seq: The Next Generation of Methylome Analysis BISULFITE SEQUENCING ALTERNATIVE METHODS FOR DETECTING 5mC AND 5hmC.** 2019,:2–5Available: <https://international.neb.com/tools-and-resources/feature-articles/enzymatic-methyl-seq-the-next-generation-of-methylome-analysis>.
107. Rhoads A, Au KF: **PacBio Sequencing and Its Applications.** *Genomics, Proteomics & Bioinformatics* 2015, **13**:278–28910.1016/j.gpb.2015.08.002Available: <https://linkinghub.elsevier.com/retrieve/pii/S1672022915001345>.
108. Simpson JT, Workman RE, Zuzarte PC, David M, Dursi LJ, Timp W: **Detecting DNA cytosine methylation using nanopore sequencing.** *Nature Methods* 2017, **14**:407–41010.1038/nmeth.4184Available: <http://www.nature.com/articles/nmeth.4184>.

109. Wilson VL, Jones PA: **DNA methylation decreases in aging but not in immortal cells.** *Science (New York, N.Y.)* 1983, **220**:1055–7Available: <http://www.ncbi.nlm.nih.gov/pubmed/6844925>.
110. Romanov GA, Vanyushin BF: **Methylation of reiterated sequences in mammalian DNAs. Effects of the tissue type, age, malignancy and hormonal induction.** *Biochimica et biophysica acta* 1981, **653**:204–18Available: <http://www.ncbi.nlm.nih.gov/pubmed/7225396>.
111. Berdyshev GD, Korotaev GK, Boiarskikh GV, Vaniushin BF: **Nucleotide composition of DNA and RNA from somatic tissues of humpback and its changes during spawning.** *Biokhimiia (Moscow, Russia)* 1967, **32**:988–93Available: <http://www.ncbi.nlm.nih.gov/pubmed/5628601>.
112. Wilson VL, Smith RA, Ma S, Cutler RG: **Genomic 5-methyldeoxycytidine decreases with age.** *The Journal of biological chemistry* 1987, **262**:9948–51Available: <http://www.ncbi.nlm.nih.gov/pubmed/3611071>.
113. Cole JJ, Robertson NA, Rather MI, Thomson JP, McBryan T, Sproul D, Wang T, Brock C, Clark W, Ideker T, Meehan RR, Miller RA, Brown-Borg HM, Adams PD: **Diverse interventions that extend mouse lifespan suppress shared age-associated epigenetic changes at critical gene regulatory regions.** *Genome Biology* 2017, **18**:5810.1186/s13059-017-1185-3Available: <http://genomebiology.biomedcentral.com/articles/10.1186/s13059-017-1185-3>.
114. Meissner A, Gnirke A, Bell GW, Ramsahoye B, Lander ES, Jaenisch R: **Reduced representation bisulfite sequencing for comparative high-resolution DNA methylation analysis.** *Nucleic acids research* 2005, **33**:5868–7710.1093/nar/gki901Available: <http://www.ncbi.nlm.nih.gov/pubmed/16224102>.
115. Feinberg AP, Tycko B: **The history of cancer epigenetics.** *Nature Reviews Cancer* 2004, **4**:143–15310.1038/nrc1279Available: <http://www.nature.com/articles/nrc1279>.
116. Knudson AG: **Mutation and cancer: statistical study of retinoblastoma.** *Proceedings of the National Academy of Sciences of the United States of America* 1971, **68**:820–310.1073/pnas.68.4.820Available: <http://www.ncbi.nlm.nih.gov/pubmed/5279523>.
117. Esteller M, Silva JM, Dominguez G, Bonilla F, Matias-Guiu X, Lerma E, Bussaglia E, Prat J, Harkes IC, Repasky EA, Gabrielson E, Schutte M, Baylin SB, Herman JG: **Promoter hypermethylation and BRCA1 inactivation in sporadic breast and ovarian tumors.** *Journal of the National Cancer Institute* 2000, **92**:564–910.1093/jnci/92.7.564Available: <http://www.ncbi.nlm.nih.gov/pubmed/10749912>.
118. Feinberg AP: **The Key Role of Epigenetics in Human Disease Prevention and Mitigation.** *New England Journal of Medicine* 2018, **378**:1323–133410.1056/NEJMra1402513Available:

<http://www.ncbi.nlm.nih.gov/pubmed/29617578>.

119. Fraga MF, Ballestar E, Paz MF, Ropero S, Setien F, Ballestar ML, Heine-Suner D, Cigudosa JC, Urioste M, Benitez J, Boix-Chornet M, Sanchez-Aguilera A, Ling C, Carlsson E, Poulsen P, Vaag A, Stephan Z, Spector TD, Wu Y-Z, Plass C, Esteller M: **From The Cover: Epigenetic differences arise during the lifetime of monozygotic twins.** *Proceedings of the National Academy of Sciences* 2005, **102**:10604–1060910.1073/pnas.0500398102Available: <http://www.ncbi.nlm.nih.gov/pubmed/16009939>.
120. Slieker RC, Iterson M van, Luijk R, Beekman M, Zhernakova DV, Moed MH, Mei H, Galen M van, Deelen P, Bonder MJ, Zhernakova A, Uitterlinden AG, Tigchelaar EF, Stehouwer CDA, Schalkwijk CG, Kallen CJH van der, Hofman A, Heemst D van, Geus EJ de, Dongen J van, Deelen J, Berg LH van den, Meurs J van, Jansen R, 't Hoen PAC, Franke L, Wijmenga C, Veldink JH, Swertz MA, Greevenbroek MMJ van, et al.: **Age-related accrual of methyomic variability is linked to fundamental ageing mechanisms.** *Genome biology* 2016, **17**:19110.1186/s13059-016-1053-6Available: <http://www.ncbi.nlm.nih.gov/pubmed/27654999>.
121. Bibikova M: **High-throughput DNA methylation profiling using universal bead arrays.** *Genome Research* 2006, **16**:383–39310.1101/gr.4410706Available: <http://www.genome.org/cgi/doi/10.1101/gr.4410706>.
122. Eckhardt F, Lewin J, Cortese R, Rakyan VK, Attwood J, Burger M, Burton J, Cox TV, Davies R, Down TA, Haefliger C, Horton R, Howe K, Jackson DK, Kunde J, Koenig C, Liddle J, Niblett D, Otto T, Pettett R, Seemann S, Thompson C, West T, Rogers J, Olek A, Berlin K, Beck S: **DNA methylation profiling of human chromosomes 6, 20 and 22.** *Nature Genetics* 2006, **38**:1378–138510.1038/ng1909Available: <http://www.nature.com/doifinder/10.1038/ng1909>.
123. Bjornsson HT: **Intra-individual Change Over Time in DNA Methylation With Familial Clustering.** *JAMA* 2008, **299**:287710.1001/jama.299.24.2877Available: <http://jama.jamanetwork.com/article.aspx?doi=10.1001/jama.299.24.2877>.
124. Christensen BC, Houseman EA, Marsit CJ, Zheng S, Wrensch MR, Wiemels JL, Nelson HH, Karagas MR, Padbury JF, Bueno R, Sugarbaker DJ, Yeh R-F, Wiencke JK, Kelsey KT: **Aging and Environmental Exposures Alter Tissue-Specific DNA Methylation Dependent upon CpG Island Context.** *PLoS Genetics* 2009, **5**:e100060210.1371/journal.pgen.1000602Available: <https://dx.plos.org/10.1371/journal.pgen.1000602>.
125. Rakyan VK, Down TA, Maslau S, Andrew T, Yang TP, Beyan H, Whittaker P, McCann OT, Finer S, Valdes AM, Leslie RD, Deloukas P, Spector TD: **Human aging-associated DNA hypermethylation occurs preferentially at bivalent chromatin domains.** *Genome Research* 2010, **20**:434–43910.1101/gr.103101.109Available: <http://genome.cshlp.org/cgi/doi/10.1101/gr.103101.109>.

1101/gr.103101.109.

126. Bibikova M, Le J, Barnes B, Saedinia-Melnyk S, Zhou L, Shen R, Gunderson KL: **Genome-wide DNA methylation profiling using Infinium assay.** *Epigenomics* 2009, **1**:177–200<sup>10.2217/epi.09.14</sup> Available: <https://www.futuremedicine.com/doi/10.2217/epi.09.14>.
127. Teschendorff AE, Menon U, Gentry-Maharaj A, Ramus SJ, Weisenberger DJ, Shen H, Campan M, Noushmehr H, Bell CG, Maxwell AP, Savage DA, Mueller-Holzner E, Marth C, Kocjan G, Gayther SA, Jones A, Beck S, Wagner W, Laird PW, Jacobs IJ, Widschwendter M: **Age-dependent DNA methylation of genes that are suppressed in stem cells is a hallmark of cancer.** *Genome research* 2010, **20**:440–6<sup>10.1101/gr.103606.109</sup> Available: <http://www.ncbi.nlm.nih.gov/pubmed/20219944>.
128. Jung M, Pfeifer GP: **Aging and DNA methylation.** *BMC Biology* 2015, **13**:<sup>710.1186/s12915-015-0118-4</sup> Available: <http://www.biomedcentral.com/1741-7007/13/7>.
129. Bocklandt S, Lin W, Sehl ME, Sánchez FJ, Sinsheimer JS, Horvath S, Vilain E: **Epigenetic Predictor of Age.** *PLoS ONE* 2011, **6**:e14821<sup>10.1371/journal.pone.0014821</sup> Available: <https://dx.plos.org/10.1371/journal.pone.0014821>.
130. Koch CM, Wagner W: **Epigenetic-aging-signature to determine age in different tissues.** *Aging* 2011, **3**:1018–27<sup>10.18632/aging.100395</sup> Available: <http://www.ncbi.nlm.nih.gov/pubmed/22067257>.
131. Bell JT, Tsai P-C, Yang T-P, Pidsley R, Nisbet J, Glass D, Mangino M, Zhai G, Zhang F, Valdes A, Shin S-Y, Dempster EL, Murray RM, Grundberg E, Hedman AK, Nica A, Small KS, MuTHER Consortium, Dermitzakis ET, McCarthy MI, Mill J, Spector TD, Deloukas P: **Epigenome-wide scans identify differentially methylated regions for age and age-related phenotypes in a healthy aging population.** *PLoS genetics* 2012, **8**:e1002629<sup>10.1371/journal.pgen.1002629</sup> Available: <http://www.ncbi.nlm.nih.gov/pubmed/22532803>.
132. Heyn H, Li N, Ferreira HJ, Moran S, Pisano DG, Gomez A, Diez J, Sanchez-Mut JV, Setien F, Carmona FJ, Puca AA, Sayols S, Pujana MA, Serra-Musach J, Iglesias-Platas I, Formiga F, Fernandez AF, Fraga MF, Heath SC, Valencia A, Gut IG, Wang J, Esteller M: **Distinct DNA methylomes of newborns and centenarians.** *Proceedings of the National Academy of Sciences* 2012, **109**:10522–10527<sup>10.1073/pnas.1120658109</sup> Available: <http://www.pnas.org/cgi/doi/10.1073/pnas.1120658109>.
133. Bibikova M, Barnes B, Tsan C, Ho V, Klotzle B, Le JM, Delano D, Zhang L, Schroth GP, Gunderson KL, Fan J-B, Shen R: **High density DNA methylation array with single CpG site resolution.** *Genomics* 2011, **98**:288–95<sup>10.1016/j.ygeno.2011.07.007</sup> Available: <http://www.ncbi.nlm.nih.gov/pubmed/21839163>.

134. Affinito O, Scala G, Palumbo D, Florio E, Monticelli A, Miele G, Avvedimento VE, Usiello A, Chiariotti L, Cocozza S: **Modeling DNA methylation by analyzing the individual configurations of single molecules.** *Epigenetics* 2016, **11**:881–88810.1080/15592294.2016.1246108 Available: <http://www.ncbi.nlm.nih.gov/pubmed/27748645>.
135. Haerter JO, Lökvist C, Dodd IB, Sneppen K: **Collaboration between CpG sites is needed for stable somatic inheritance of DNA methylation states.** *Nucleic Acids Research* 2014, **42**:2235–224410.1093/nar/gkt1235 Available: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3936770>.
136. Garagnani P, Bacalini MG, Pirazzini C, Gori D, Giuliani C, Mari D, Di Blasio AM, Gentilini D, Vitale G, Collino S, Rezzi S, Castellani G, Capri M, Salvioli S, Franceschi C: **Methylation of ELOVL2 gene as a new epigenetic marker of age.** *Aging cell* 2012, **11**:1132–410.1111/acel.12005 Available: <http://www.ncbi.nlm.nih.gov/pubmed/23061750>.
137. Hannum G, Guinney J, Zhao L, Zhang L, Hughes G, Sadda S, Klotzle B, Bibikova M, Fan J-B, Gao Y, Deconde R, Chen M, Rajapakse I, Friend S, Ideker T, Zhang K: **Genome-wide methylation profiles reveal quantitative views of human aging rates.** *Molecular cell* 2013, **49**:359–36710.1016/j.molcel.2012.10.016 Available: <http://www.ncbi.nlm.nih.gov/pubmed/23177740>.
138. Horvath S: **DNA methylation age of human tissues and cell types.** *Genome biology* 2013, **14**:R11510.1186/gb-2013-14-10-r115 Available: <http://www.ncbi.nlm.nih.gov/pubmed/24138928>.
139. Friedman J, Hastie T, Tibshirani R: **Regularization Paths for Generalized Linear Models via Coordinate Descent.** *Journal of Statistical Software* 2010, **33**:7–1010.18637/jss.v033.i01 Available: <http://www.jstatsoft.org/v33/i01/>.
140. Florath I, Butterbach K, Muller H, Bewerunge-Hudler M, Brenner H: **Cross-sectional and longitudinal changes in DNA methylation with age: an epigenome-wide analysis revealing over 60 novel age-associated CpG sites.** *Human Molecular Genetics* 2014, **23**:1186–120110.1093/hmg/ddt531 Available: <https://academic.oup.com/hmg/article-lookup/doi/10.1093/hmg/ddt531>.
141. Bacalini MG, Boattini A, Gentilini D, Giampieri E, Pirazzini C, Giuliani C, Fontanesi E, Remondini D, Capri M, Del Rio A, Luiselli D, Vitale G, Mari D, Castellani G, Di Blasio AM, Salvioli S, Franceschi C, Garagnani P: **A meta-analysis on age-associated changes in blood DNA methylation: results from an original analysis pipeline for Infinium 450k data.** *Aging* 2015, **7**:97–10910.18632/aging.100718 Available: <http://www.aging-us.com/article/100718>.
142. Zaghlool SB, Al-Shafai M, Al Muftah WA, Kumar P, Falchi M, Suhre K: **Association of**

**DNA methylation with age, gender, and smoking in an Arab population.** *Clinical Epigenetics* 2015, **7**:610.1186/s13148-014-0040-6 Available: <http://www.ncbi.nlm.nih.gov/pubmed/25663950>.

143. Benton MC, Sutherland HG, Macartney-Coxson D, Haupt LM, Lea RA, Griffiths LR: **Methylome-wide association study of whole blood DNA in the Norfolk Island isolate identifies robust loci associated with age.** *Aging* 2017, **9**:753–76810.18632/aging.101187 Available: <http://www.aging-us.com/article/101187/text>.

144. Johnson ND, Wiener HW, Smith AK, Nishitani S, Absher DM, Arnett DK, Aslibekyan S, Conneely KN: **Non-linear patterns in age-related DNA methylation may reflect CD4 + T cell differentiation.** *Epigenetics* 2017, **12**:492–50310.1080/15592294.2017.1314419 Available: <http://www.ncbi.nlm.nih.gov/pubmed/28387568>.

145. Smith AK, Kilaru V, Kocak M, Almli LM, Mercer KB, Ressler KJ, Tylavsky FA, Conneely KN: **Methylation quantitative trait loci (meQTLs) are consistently detected across ancestry, developmental stage, and tissue type.** *BMC Genomics* 2014, **15**:14510.1186/1471-2164-15-145 Available: <http://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2164-15-145>.

146. Jaiswal S, Fontanillas P, Flannick J, Manning A, Grauman PV, Mar BG, Lindsley RC, Mermel CH, Burtt N, Chavez A, Higgins JM, Moltchanov V, Kuo FC, Kluk MJ, Henderson B, Kinnunen L, Koistinen HA, Ladenvall C, Getz G, Correa A, Banahan BF, Gabriel S, Kathiresan S, Stringham HM, McCarthy MI, Boehnke M, Tuomilehto J, Haiman C, Groop L, Atzmon G, et al.: **Age-Related Clonal Hematopoiesis Associated with Adverse Outcomes.** *New England Journal of Medicine* 2014, **371**:2488–249810.1056/NEJMoa1408617 Available: <http://www.nejm.org/doi/10.1056/NEJMoa1408617>.

147. Weidner CI, Lin Q, Koch CM, Eisele L, Beier F, Ziegler P, Bauerschlag DO, Jöckel K-H, Erbel R, Mühlleisen TW, Zenke M, Brümmendorf TH, Wagner W: **Aging of blood can be tracked by DNA methylation changes at just three CpG sites.** *Genome biology* 2014, **15**:R2410.1186/gb-2014-15-2-r24 Available: <http://www.ncbi.nlm.nih.gov/pubmed/24490752>.

148. Johnson TE: **Recent results: Biomarkers of aging.** *Experimental Gerontology* 2006, **41**:1243–124610.1016/j.exger.2006.09.006 Available: <https://linkinghub.elsevier.com/retrieve/pii/S0531556506002865>.

149. Thompson MJ, VonHoldt B, Horvath S, Pellegrini M: **An epigenetic aging clock for dogs and wolves.** *Aging* 2017, **9**:1055–106810.18632/aging.101211 Available: <http://www.ncbi.nlm.nih.gov/pubmed/28373601>.

150. Stubbs TM, Bonder MJ, Stark A-K, Krueger F, Meyenn F von, Stegle O, Reik W: **Multi-tissue DNA methylation age predictor in mouse.** *Genome Biology* 2017,

- 18:6810.1186/s13059-017-1203-5 Available: <http://www.ncbi.nlm.nih.gov/pubmed/28399939>.
151. Lowe R, Danson AF, Rakyan VK, Yildizoglu S, Saldmann F, Viltard M, Friedlander G, Faulkes CG: **DNA methylation clocks as a predictor for ageing and age estimation in naked mole-rats , *Heterocephalus glaber***. 2020, **12**:1–1310.18632/aging.102892.
152. Lowe R, Barton C, Jenkins CA, Ernst C, Forman O, Fernandez-Twinn DS, Bock C, Rossiter SJ, Faulkes CG, Ozanne SE, Walter L, Odom DT, Mellersh C, Rakyan VK: **Ageing-associated DNA methylation dynamics are a molecular readout of lifespan variation among mammalian species**. *Genome Biology* 2018, **19**:2210.1186/s13059-018-1397-1 Available: <http://www.ncbi.nlm.nih.gov/pubmed/29452591>.
153. Wang M, Lemos B: **Ribosomal DNA harbors an evolutionarily conserved clock of biological aging**. *Genome research* 2019, **29**:325–33310.1101/gr.241745.118 Available: <http://www.ncbi.nlm.nih.gov/pubmed/30765617>.
154. Marioni RE, Shah S, McRae AF, Chen BH, Colicino E, Harris SE, Gibson J, Henders AK, Redmond P, Cox SR, Pattie A, Corley J, Murphy L, Martin NG, Montgomery GW, Feinberg AP, Fallin MD, Multhaup ML, Jaffe AE, Joehanes R, Schwartz J, Just AC, Lunetta KL, Murabito JM, Starr JM, Horvath S, Baccarelli AA, Levy D, Visscher PM, Wray NR, et al.: **DNA methylation age of blood predicts all-cause mortality in later life**. *Genome biology* 2015, **16**:2510.1186/s13059-015-0584-6 Available: <http://www.ncbi.nlm.nih.gov/pubmed/25633388>.
155. Christiansen L, Lenart A, Tan Q, Vaupel JW, Aviv A, McGue M, Christensen K: **DNA methylation age is associated with mortality in a longitudinal Danish twin study**. *Aging Cell* 2016, **15**:149–15410.1111/acel.12421 Available: <http://doi.wiley.com/10.1111/acel.12421>.
156. Chen BH, Marioni RE, Colicino E, Peters MJ, Ward-Caviness CK, Tsai P-C, Roetker NS, Just AC, Demerath EW, Guan W, Bressler J, Fornage M, Studenski S, Vandiver AR, Moore AZ, Tanaka T, Kiel DP, Liang L, Vokonas P, Schwartz J, Lunetta KL, Murabito JM, Bandinelli S, Hernandez DG, Melzer D, Nalls M, Pilling LC, Price TR, Singleton AB, Gieger C, et al.: **DNA methylation-based measures of biological age: meta-analysis predicting time to death**. *Aging* 2016, **8**:1844–186510.18632/aging.101020 Available: <http://www.ncbi.nlm.nih.gov/pubmed/27690265>.
157. Levine ME, Lu AT, Quach A, Chen BH, Assimes TL, Bandinelli S, Hou L, Baccarelli AA, Stewart JD, Li Y, Whitsel EA, Wilson JG, Reiner AP, Aviv A, Lohman K, Liu Y, Ferrucci L, Horvath S: **An epigenetic biomarker of aging for lifespan and healthspan**. *Aging* 2018, **10**:573–59110.18632/aging.101414 Available: <http://www.ncbi.nlm.nih.gov/pubmed/29676998>.
158. Lu AT, Quach A, Wilson JG, Reiner AP, Aviv A, Raj K, Hou L, Baccarelli AA, Li Y, Stewart JD, Whitsel EA, Assimes TL, Ferrucci L, Horvath S: **DNA methylation GrimAge strongly**

- predicts lifespan and healthspan.** *Aging* 2019, **11**:303–32710.18632/aging.101684Available: <http://www.aging-us.com/article/101684/text>.
159. Zhang Y, Hapala J, Brenner H, Wagner W: **Individual CpG sites that are associated with age and life expectancy become hypomethylated upon aging.** *Clinical Epigenetics* 2017, **9**:910.1186/s13148-017-0315-9Available: <http://clincalepigeneticsjournal.biomedcentral.com/articles/10.1186/s13148-017-0315-9>.
160. Jones MJ, Goodman SJ, Kobor MS: **DNA methylation and healthy human aging.** *Aging Cell* 2015, **14**:924–93210.1111/acel.12349Available: <http://doi.wiley.com/10.1111/acel.12349>.
161. Rimmelé P, Bigarella CL, Liang R, Izac B, Dieguez-Gonzalez R, Barbet G, Donovan M, Brugnara C, Blander JM, Sinclair DA, Ghaffari S: **Aging-like phenotype and defective lineage specification in SIRT1-deleted hematopoietic stem and progenitor cells.** *Stem cell reports* 2014, **3**:44–5910.1016/j.stemcr.2014.04.015Available: <http://www.ncbi.nlm.nih.gov/pubmed/25068121>.
162. Izzo F, Lee SC, Poran A, Chaligne R, Gaiti F, Gross B, Murali RR, Deochand SD, Ang C, Jones PW, Nam AS, Kim K-T, Kothen-Hill S, Schulman RC, Ki M, Lhoumaud P, Skok JA, Viny AD, Levine RL, Kenigsberg E, Abdel-Wahab O, Landau DA: **DNA methylation disruption reshapes the hematopoietic differentiation landscape.** *Nature Genetics* 2020, **52**10.1038/s41588-020-0595-4Available: <http://dx.doi.org/10.1038/s41588-020-0595-4>.
163. Geiger H, Haan G de, Florian MC: **The ageing haematopoietic stem cell compartment.** *Nature Reviews Immunology* 2013, **13**:376–38910.1038/nri3433Available: <http://dx.doi.org/10.1038/nri3433>.
164. Jaffe AE, Irizarry RA: **Accounting for cellular heterogeneity is critical in epigenome-wide association studies.** *Genome biology* 2014, **15**:R3110.1186/gb-2014-15-2-r31Available: <http://www.ncbi.nlm.nih.gov/pubmed/2449553>.
165. 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**:8610.1186/1471-2105-13-86Available: <http://www.biomedcentral.com/1471-2105/13/86>.
166. Schuettengruber B, Chourrout D, Vervoort M, Leblanc B, Cavalli G: **Genome Regulation by Polycomb and Trithorax Proteins.** *Cell* 2007, **128**:735–74510.1016/j.cell.2007.02.009Available: <http://www.ncbi.nlm.nih.gov/pubmed/17320510>.
167. Hahn O, Grönke S, Stubbs TM, Ficz G, Hendrich O, Krueger F, Andrews S, Zhang Q, Wakelam MJ, Beyer A, Reik W, Partridge L: **Dietary restriction protects from age-associated DNA methylation and induces epigenetic reprogramming**

- of lipid metabolism. *Genome Biology* 2017, **18**:5610.1186/s13059-017-1187-1 Available: <http://genomebiology.biomedcentral.com/articles/10.1186/s13059-017-1187-1>.
168. Horvath S, Raj K: **DNA methylation-based biomarkers and the epigenetic clock theory of ageing.** *Nature Reviews Genetics* 2018, **19**:371–38410.1038/s41576-018-0004-3 Available: <http://www.nature.com/articles/s41576-018-0004-3>.
169. Bell CG, Lowe R, Adams PD, Baccarelli AA, Beck S, Bell JT, Christensen BC, Gladyshev VN, Heijmans BT, Horvath S, Ideker T, Issa J-PJ, Kelsey KT, Marioni RE, Reik W, Relton CL, Schalkwyk LC, Teschendorff AE, Wagner W, Zhang K, Rakyan VK: **DNA methylation aging clocks: challenges and recommendations.** *Genome biology* 2019, **20**:24910.1186/s13059-019-1824-y Available: <http://www.ncbi.nlm.nih.gov/pubmed/31767039>.
170. Bell CG, Xia Y, Yuan W, Gao F, Ward K, Roos L, Mangino M, Hysi PG, Bell J, Wang J, Spector TD: **Novel regional age-associated DNA methylation changes within human common disease-associated loci.** *Genome Biology* 2016, **17**:19310.1186/s13059-016-1051-8 Available: <http://www.ncbi.nlm.nih.gov/pubmed/27663977>.
171. Harvey NC, Javaid K, Bishop N, Kennedy S, Papageorghiou AT, Fraser R, Gandhi SV, Schoenmakers I, Prentice A, Cooper C: **MAVIDOS Maternal Vitamin D Osteoporosis Study: study protocol for a randomized controlled trial.** The MAVIDOS Study Group. *Trials* 2012, **13**:1310.1186/1745-6215-13-13 Available: <http://trialsjournal.biomedcentral.com/articles/10.1186/1745-6215-13-13>.
172. Cooper C, Harvey NC, Bishop NJ, Kennedy S, Papageorghiou AT, Schoenmakers I, Fraser R, Gandhi SV, Carr A, D'Angelo S, Crozier SR, Moon RJ, Arden NK, Dennison EM, Godfrey KM, Inskip HM, Prentice A, Mughal MZ, Eastell R, Reid DM, Javaid MK: **Maternal gestational vitamin D supplementation and offspring bone health (MAVIDOS): a multicentre, double-blind, randomised placebo-controlled trial.** *The Lancet Diabetes & Endocrinology* 2016, **4**:393–40210.1016/S2213-8587(16)00044-9 Available: [http://dx.doi.org/10.1016/S2213-8587\(16\)00044-9](http://dx.doi.org/10.1016/S2213-8587(16)00044-9).
173. Moran S, Arribas C, Esteller M: **Validation of a DNA methylation microarray for 850,000 CpG sites of the human genome enriched in enhancer sequences.** *Epigenomics* 2016, **8**:389–9910.2217/epi.15.114 Available: <http://www.ncbi.nlm.nih.gov/pubmed/26673039>.
174. Gunderson KL: **Decoding Randomly Ordered DNA Arrays.** *Genome Research* 2004, **14**:870–87710.1101/gr.2255804 Available: <http://www.genome.org/cgi/doi/10.1101/gr.2255804>.
175. Michael KL, Taylor LC, Schultz SL, Walt DR: **Randomly Ordered Addressable High-Density Optical Sensor Arrays.** *Analytical Chemistry* 1998, **70**:1242–124810.1021/ac971343r Available: <https://pubs.acs.org/doi/10.1021/ac971343r>.
176. Clark C, Palta P, Joyce CJ, Scott C, Grundberg E, Deloukas P, Palotie A, Coffey AJ:

**A Comparison of the Whole Genome Approach of MeDIP-Seq to the Targeted Approach of the Infinium HumanMethylation450 BeadChip for Methylome Profiling.** *PLoS ONE* 2012, **7**:e5023310.1371/journal.pone.0050233 Available: <https://dx.plos.org/10.1371/journal.pone.0050233>.

177. Smith ML, Baggerly KA, Bengtsson H, Ritchie ME, Hansen KD: **illuminaio: An open source IDAT parsing tool for Illumina microarrays.** *F1000Research* 2013, **2**:1–810.12688/f1000research.2-264.v1.

178. Dedeurwaerder S, Defrance M, Calonne E, Denis H, Sotiriou C, Fuks F: **Evaluation of the Infinium Methylation 450K technology.** *Epigenomics* 2011, **3**:771–78410.2217/epi.11.105 Available: <https://www.futuremedicine.com/doi/10.2217/epi.11.105>.

179. Fortin J-P, Triche TJ, Hansen KD: **Preprocessing, normalization and integration of the Illumina HumanMethylationEPIC array with minfi.** *Bioinformatics* 2016,:btw69110.1093/bioinformatics/btw691 Available: <http://biorkiv.org/lookup/doi/10.1101/065490>.

180. Teschendorff AE, Marabita F, Lechner M, Bartlett T, Tegner J, Gomez-Cabrero D, Beck S: **A beta-mixture quantile normalization method for correcting probe design bias in Illumina Infinium 450 k DNA methylation data.** *Bioinformatics* 2013, **29**:189–19610.1093/bioinformatics/bts680 Available: <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/bts680>.

181. Leek JT, Scharpf RB, Bravo HC, Simcha D, Langmead B, Johnson WE, Geman D, Baggerly K, Irizarry RA: **Tackling the widespread and critical impact of batch effects in high-throughput data.** *Nature Reviews Genetics* 2010, **11**:733–73910.1038/nrg2825 Available: <http://www.nature.com/articles/nrg2825>.

182. Fortin J-P, Labbe A, Lemire M, Zanke BW, Hudson TJ, Fertig EJ, Greenwood CMT, Hansen KD: **Functional normalization of 450k methylation array data improves replication in large cancer studies.** *Genome Biology* 2014, **15**:50310.1186/s13059-014-0503-2 Available: <http://genomebiology.biomedcentral.com/articles/10.1186/s13059-014-0503-2>.

183. Ehrlich M, Zoll S, Sur S, Boom D van den: **A new method for accurate assessment of DNA quality after bisulfite treatment.** *Nucleic acids research* 2007, **35**:e2910.1093/nar/gkl1134 Available: <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkl1134%20> <http://www.ncbi.nlm.nih.gov/pubmed/17259213%20> <http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC1865059>.

184. Illumina: *Infinium HD Assay Methylation Protocol Guide*. Illumina; 2015 Available: [https://support.illumina.com/content/dam/illumina-support/documents/documentation/chemistry%7B/\\_%7Ddocumentation/infinium%7B/\\_%7Dassays/infinium%7B/\\_%7Dhd%7B/\\_%7D](https://support.illumina.com/content/dam/illumina-support/documents/documentation/chemistry%7B/_%7Ddocumentation/infinium%7B/_%7Dassays/infinium%7B/_%7Dhd%7B/_%7D)

methylation/infinium-hd-methylation-guide-15019519-01.pdf.

185. Pidsley R, Zotenko E, Peters TJ, Lawrence MG, Risbridger GP, Molloy P, Van Djik S, Muhlhausler B, Stirzaker C, Clark SJ: **Critical evaluation of the Illumina MethylationEPIC BeadChip microarray for whole-genome DNA methylation profiling.** *Genome biology* 2016, **17**:20810.1186/s13059-016-1066-1 Available: <http://www.ncbi.nlm.nih.gov/pubmed/27717381>.
186. Zhou W, Laird PW, Shen H: **Comprehensive characterization, annotation and innovative use of Infinium DNA methylation BeadChip probes.** *Nucleic Acids Research* 2016, **45**:gkw96710.1093/nar/gkw967 Available: <http://www.ncbi.nlm.nih.gov/pubmed/27924034>.
187. Birney E, Smith GD, Greally JM: **Epigenome-wide Association Studies and the Interpretation of Disease -Omics.** *PLOS Genetics* 2016, **12**:e100610510.1371/journal.pgen.1006105 Available: <http://www.ncbi.nlm.nih.gov/pubmed/27336614>.
188. Lappalainen T, Greally JM: **Associating cellular epigenetic models with human phenotypes.** *Nature Reviews Genetics* 2017, **18**:441–45110.1038/nrg.2017.32 Available: <http://www.nature.com/doifinder/10.1038/nrg.2017.32>.
189. Andrews SV, Ladd-Acosta C, Feinberg AP, Hansen KD, Fallin MD: **“Gap hunting” to characterize clustered probe signals in Illumina methylation array data.** *Epigenetics & Chromatin* 2016, **9**:5610.1186/s13072-016-0107-z Available: <http://epigeneticsandchromatin.biomedcentral.com/articles/10.1186/s13072-016-0107-z>.
190. Taiwo O, Wilson GA, Morris T, Seisenberger S, Reik W, Pearce D, Beck S, Butcher LM: **Methylome analysis using MeDIP-seq with low DNA concentrations.** *Nature Protocols* 2012, **7**:617–63610.1038/nprot.2012.012 Available: <http://www.nature.com/articles/nprot.2012.012>.
191. Chavez L, Jozefczuk J, Grimm C, Dietrich J, Timmermann B, Lehrach H, Herwig R, Adjaye J: **Computational analysis of genome-wide DNA methylation during the differentiation of human embryonic stem cells along the endodermal lineage.** *Genome Research* 2010, **20**:1441–145010.1101/gr.110114.110 Available: <http://genome.cshlp.org/cgi/doi/10.1101/gr.110114.110>.
192. Andrews S: **FastQC.** 2010, Available: <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>.
193. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R: **The Sequence Alignment/Map format and SAMtools.** *Bioinformatics* 2009, **25**:2078–207910.1093/bioinformatics/btp352 Available: <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btp352>.

194. Quinlan AR, Hall IM: **BEDTools: a flexible suite of utilities for comparing genomic features.** *Bioinformatics* 2010, **26**:841–84210.1093/bioinformatics/btq033 Available: <http://bedtools.readthedocs.io/en/latest/index.html>.
195. Moayyeri A, Hammond CJ, Valdes AM, Spector TD: **Cohort Profile: TwinsUK and Healthy Ageing Twin Study.** *International Journal of Epidemiology* 2013, **42**:76–8510.1093/ije/dyr207 Available: <http://www.ncbi.nlm.nih.gov/pubmed/22253318>.
196. Korbie D, Lin E, Wall D, Nair SS, Stirzaker C, Clark SJ, Trau M: **Multiplex bisulfite PCR resequencing of clinical FFPE DNA.** *Clinical Epigenetics* 2015, **7**:2810.1186/s13148-015-0067-3 Available: <https://clincialepigenticsjournal.biomedcentral.com/articles/10.1186/s13148-015-0067-3>.
197. Tran H, Porter J, Sun M-a, Xie H, Zhang L: **Objective and Comprehensive Evaluation of Bisulfite Short Read Mapping Tools.** *Advances in Bioinformatics* 2014, **2014**:1–1110.1155/2014/472045 Available: <http://www.hindawi.com/journals/abi/2014/472045/>.
198. Böttcher R, Amberg R, Ruzius FP, Guryev V, Verhaegh WFJ, Beyerlein P, Zaag PJ van der: **Using a priori knowledge to align sequencing reads to their exact genomic position.** *Nucleic Acids Research* 2012, **40**:e125–e12510.1093/nar/gks393 Available: <https://academic.oup.com/nar/article/40/16/e125/1026881>.
199. Krueger F, Andrews SR: **Bismark: a flexible aligner and methylation caller for Bisulfite-Seq applications.** *Bioinformatics* 2011, **27**:1571–157210.1093/bioinformatics/btr167 Available: <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btr167>.
200. Li L-C, Dahiya R: **MethPrimer: designing primers for methylation PCRs.** *Bioinformatics (Oxford, England)* 2002, **18**:1427–31 Available: <http://www.ncbi.nlm.nih.gov/pubmed/12424112>.
201. Staa TP van, Dennison EM, Leufkens HGM, Cooper C: **Epidemiology of fractures in England and Wales.** *Bone* 2001, **29**:517–52210.1016/S8756-3282(01)00614-7 Available: <https://linkinghub.elsevier.com/retrieve/pii/S8756328201006147>.
202. Streubel PN, Ricci WM, Wong A, Gardner MJ: **Mortality After Distal Femur Fractures in Elderly Patients.** *Clinical Orthopaedics and Related Research* 2011, **469**:1188–119610.1007/s11999-010-1530-2 Available: <http://link.springer.com/10.1007/s11999-010-1530-2>.
203. Thompson E, Greenspan S: *National Osteoporosis Foundation Annual Report (2018).* National Osteoporosis Foundation; 2018 Available: [https://cdn.nof.org/wp-content/uploads/2018%7B/\\_%7DNOF%7B/\\_%7DAnnual%7B/\\_%7Dreport%7B/\\_%7DFINAL.pdf](https://cdn.nof.org/wp-content/uploads/2018%7B/_%7DNOF%7B/_%7DAnnual%7B/_%7Dreport%7B/_%7DFINAL.pdf).
204. Hernandez CJ, Beaupré GS, Carter DR: **A theoretical analysis of the relative influences of peak BMD, age-related bone loss and menopause on the development of osteoporosis.** *Osteoporosis International* 2003, **14**:843–84710.1007/s00198-003-1454-8.

205. Nguyen TV, Center JR, Eisman JA: **Femoral Neck Bone Loss Predicts Fracture Risk Independent of Baseline BMD.** *Journal of Bone and Mineral Research* 2005, **20**:1195–120110.1359/JBMR.050215 Available: <http://doi.wiley.com/10.1359/JBMR.050215>.
206. Finkelstein JS, Brockwell SE, Mehta V, Greendale GA, Sowers MR, Ettinger B, Lo JC, Johnston JM, Cauley JA, Danielson ME, Neer RM: **Bone Mineral Density Changes during the Menopause Transition in a Multiethnic Cohort of Women.** *The Journal of Clinical Endocrinology & Metabolism* 2008, **93**:861–86810.1210/jc.2007-1876 Available: <https://academic.oup.com/jcem/article-lookup/doi/10.1210/jc.2007-1876>.
207. Hui SL, Slemenda CW, Johnston CC: **The contribution of bone loss to post-menopausal osteoporosis.** *Osteoporosis International* 1990, **1**:30–3410.1007/BF01880413.
208. Harvey N, Dennison E, Cooper C: **Osteoporosis: A Lifecourse Approach.** *Journal of Bone and Mineral Research* 2014, **29**:1917–192510.1002/jbmr.2286 Available: <http://doi.wiley.com/10.1002/jbmr.2286>.
209. Bikle DD: **Vitamin D and Bone.** *Current Osteoporosis Reports* 2012, **10**:151–15910.1007/s11914-012-0098-z Available: <http://link.springer.com/10.1007/s11914-012-0098-z>.
210. Christakos S, Dhawan P, Porta A, Mady LJ, Seth T: **Vitamin D and intestinal calcium absorption.** *Molecular and Cellular Endocrinology* 2011, **347**:25–2910.1016/j.mce.2011.05.038 Available: <https://linkinghub.elsevier.com/retrieve/pii/S0303720711002930>.
211. Fukumoto S: **Phosphate metabolism and vitamin D.** *BoneKEy Reports* 2014, **3**:1–510.1038/bonekey.2013.231 Available: <http://dx.doi.org/10.1038/bonekey.2013.231>.
212. Baird J, Kurshid MA, Kim M, Harvey N, Dennison E, Cooper C: **Does birthweight predict bone mass in adulthood? A systematic review and meta-analysis.** *Osteoporosis International* 2011, **22**:1323–133410.1007/s00198-010-1344-9.
213. Cooper C, Eriksson JG, Forsén T, Osmond C, Tuomilehto J, Barker DJP: **Maternal height, childhood growth and risk of hip fracture in later life: A longitudinal study.** *Osteoporosis International* 2001, **12**:623–62910.1007/s001980170061.
214. Inskip HM, Godfrey KM, Robinson SM, Law CM, Jp D: **Europe PMC Funders Group Cohort Profile : The Southampton Women ' s Survey.** *International Journal of Epidemiology* 2006, **35**:42–4810.1093/ije/dyi202.
215. Harvey NC, Mahon PA, Robinson SM, Nisbet CE, Javaid MK, Crozier SR, Inskip HM, Godfrey KM, Arden NK, Dennison EM, Cooper C, Taylor P, Greenaway LJ, Hanson M, Barker DJP, Law CM: **Different indices of fetal growth predict bone size and volumetric density at 4 years of age.** *Journal of Bone and Mineral Research* 2010, **25**:920–92710.1359/jbmr.091022.

216. Harvey NC, Cole ZA, Crozier SR, Ntani G, Mahon PA, Robinson SM, Inskip HM, Godfrey KM, Dennison EM, Cooper C: **Fetal and infant growth predict hip geometry at 6 y old: Findings from the Southampton Women's Survey.** *Pediatric Research* 2013, **74**:450–45610.1038/pr.2013.119.
217. Mahon P, Harvey N, Crozier S, Inskip H, Robinson S, Arden N, Swaminathan R, Cooper C, Godfrey K: **Low maternal vitamin D status and fetal bone development: Cohort study.** *Journal of Bone and Mineral Research* 2010, **25**:14–1910.1359/jbmр.090701 Available: <http://doi.wiley.com/10.1359/jbmр.090701>.
218. Viljakainen HT, Saarnio E, Hytinantti T, Miettinen M, Surcel H, Mäkitie O, Andersson S, Laitinen K, Lamberg-Allardt C: **Maternal Vitamin D Status Determines Bone Variables in the Newborn.** *The Journal of Clinical Endocrinology & Metabolism* 2010, **95**:1749–175710.1210/jc.2009-1391 Available: <https://academic.oup.com/jcem/article-lookup/doi/10.1210/jc.2009-1391>.
219. Viljakainen HT, Korhonen T, Hytinantti T, Laitinen EKA, Andersson S, Mäkitie O, Lamberg-Allardt C: **Maternal vitamin D status affects bone growth in early childhood—a prospective cohort study.** *Osteoporosis International* 2011, **22**:883–89110.1007/s00198-010-1499-4 Available: <http://link.springer.com/10.1007/s00198-010-1499-4>.
220. Zhu K, Whitehouse AJ, Hart PH, Kusel M, Mountain J, Lye S, Pennell C, Walsh JP: **Maternal Vitamin D Status During Pregnancy and Bone Mass in Offspring at 20 Years of Age: A Prospective Cohort Study.** *Journal of Bone and Mineral Research* 2014, **29**:1088–109510.1002/jbmr.2138 Available: <http://doi.wiley.com/10.1002/jbmr.2138>.
221. Kip SN, Strehler EE: **Vitamin D 3 upregulates plasma membrane Ca 2+ -ATPase expression and potentiates apico-basal Ca 2+ flux in MDCK cells.** *American Journal of Physiology-Renal Physiology* 2004, **286**:F363–F36910.1152/ajprenal.00076.2003 Available: <http://www.ncbi.nlm.nih.gov/pubmed/14583431>.
222. Burdge GC, Slater-Jeffries J, Torrens C, Phillips ES, Hanson MA, Lillycrop KA: **Dietary protein restriction of pregnant rats in the F 0 generation induces altered methylation of hepatic gene promoters in the adult male offspring in the F 1 and F 2 generations.** *British Journal of Nutrition* 2007, **97**:435–43910.1017/S0007114507352392 Available: [https://www.cambridge.org/core/product/identifier/S0007114507352392/type/journal%7B/\\_%7Darticle](https://www.cambridge.org/core/product/identifier/S0007114507352392/type/journal%7B/_%7Darticle).
223. Lillycrop KA, Phillips ES, Torrens C, Hanson MA, Jackson AA, Burdge GC: **Feeding pregnant rats a protein-restricted diet persistently alters the methylation of specific cytosines in the hepatic PPAR $\alpha$  promoter of the offspring.** *British Journal of Nutrition* 2008, **100**:278–28210.1017/S0007114507894438 Available: [https://www.cambridge.org/core/product/identifier/S0007114507894438/type/journal%7B/\\_%7Darticle](https://www.cambridge.org/core/product/identifier/S0007114507894438/type/journal%7B/_%7Darticle).

224. Burdge GC, Hanson MA, Slater-Jeffries JL, Lillycrop KA: **Epigenetic regulation of transcription: a mechanism for inducing variations in phenotype (fetal programming) by differences in nutrition during early life?** *British Journal of Nutrition* 2007, **97**:1036–104610.1017/S0007114507682920Available: [https://www.cambridge.org/core/product/identifier/S0007114507682920/type/journal%7B/\\_%7Darticle](https://www.cambridge.org/core/product/identifier/S0007114507682920/type/journal%7B/_%7Darticle).
225. Holroyd C, Harvey N, Dennison E, Cooper C: **Epigenetic influences in the developmental origins of osteoporosis.** *Osteoporosis International* 2012, **23**:401–41010.1007/s00198-011-1671-5.
226. Harvey NC, Sheppard A, Godfrey KM, McLean C, Garratt E, Ntani G, Davies L, Murray R, Inskip HM, Gluckman PD, Hanson MA, Lillycrop KA, Cooper C: **Childhood Bone Mineral Content Is Associated With Methylation Status of the RXRA Promoter at Birth.** *Journal of Bone and Mineral Research* 2014, **29**:600–60710.1002/jbmr.2056Available: <http://doi.wiley.com/10.1002/jbmr.2056>.
227. Curtis EM, Murray R, Titcombe P, Cook E, Clarke-Harris R, Costello P, Garratt E, Holbrook JD, Barton S, Inskip H, Godfrey KM, Bell CG, Cooper C, Lillycrop KA, Harvey NC: **Perinatal DNA Methylation at CDKN2A Is Associated With Offspring Bone Mass: Findings From the Southampton Women’s Survey.** *Journal of Bone and Mineral Research* 2017, **32**:2030–204010.1002/jbmr.3153Available: <http://doi.wiley.com/10.1002/jbmr.3153>.
228. Holland ML, Lowe R, Caton PW, Gemma C, Carabajosa G, Danson AF, Carpenter AAM, Loche E, Ozanne SE, Rakyan VK: **Early-life nutrition modulates the epigenetic state of specific rDNA genetic variants in mice.** *Science* 2016, **353**:495–49810.1126/science.aaf7040Available: <http://www.sciencemag.org/lookup/doi/10.1126/science.aaf7040>.
229. Dogan MV, Beach SRH, Philibert RA: **Genetically contextual effects of smoking on genome wide DNA methylation.** *American Journal of Medical Genetics Part B: Neuropsychiatric Genetics* 2017, **174**:595–60710.1002/ajmg.b.32565Available: <http://doi.wiley.com/10.1002/ajmg.b.32565>.
230. Dai L, Mehta A, Mordukhovich I, Just AC, Shen J, Hou L, Koutrakis P, Sparrow D, Vokonas PS, Baccarelli AA, Schwartz JD: **Differential DNA methylation and PM 2.5 species in a 450K epigenome-wide association study.** *Epigenetics* 2017, **12**:139–14810.1080/15592294.2016.1271853Available: <https://www.tandfonline.com/doi/full/10.1080/15592294.2016.1271853>.
231. Wahl S, Drong A, Lehne B, Loh M, Scott WR, Kunze S, Tsai P-C, Ried JS, Zhang W, Yang Y, Tan S, Fiorito G, Franke L, Guerrera S, Kasela S, Kriebel J, Richmond RC, Adamo M, Afzal U, Ala-Korpela M, Albetti B, Ammerpohl O, Apperley JF, Beekman M, Bertazzi PA, Black SL, Blancher C, Bonder M-J, Brosch M, Carstensen-Kirberg M, et al.: **Epigenome-wide**

- association study of body mass index, and the adverse outcomes of adiposity.** *Nature* 2017, **541**:81–8610.1038/nature20784Available: <http://www.ncbi.nlm.nih.gov/pubmed/28002404>.
232. Min J, Hemani G, Smith GD, Relton CL, Suderman M: **Meffil: efficient normalisation and analysis of very large DNA methylation samples.** *bioRxiv* 2017, **44**:12596310.1101/125963Available: <https://www.biorxiv.org/content/early/2017/04/27/125963>.
233. Cardenas A, Allard C, Doyon M, Houseman EA, Bakulski KM, Perron P, Bouchard L, Hivert M-F: **Validation of a DNA methylation reference panel for the estimation of nucleated cells types in cord blood.** *Epigenetics* 2016, **11**:773–77910.1080/15592294.2016.1233091Available: <http://www.ncbi.nlm.nih.gov/pubmed/27668573>.
234. Goede OM de, Razzaghian HR, Price EM, Jones MJ, Kobor MS, Robinson WP, Lavoie PM: **Nucleated red blood cells impact DNA methylation and expression analyses of cord blood hematopoietic cells.** *Clinical Epigenetics* 2015, **7**:9510.1186/s13148-015-0129-6Available: <http://dx.doi.org/10.1186/s13148-015-0129-6>.
235. Bakulski KM, Feinberg JI, Andrews SV, Yang J, Brown S, L. McKenney S, Witter F, Walston J, Feinberg AP, Fallin MD: **DNA methylation of cord blood cell types: Applications for mixed cell birth studies.** *Epigenetics* 2016, **11**:354–36210.1080/15592294.2016.1161875Available: <http://dx.doi.org/10.1080/15592294.2016.1161875>.
236. Gervin K, Page CM, Aass HCD, Jansen MA, Fjeldstad HE, Andreassen BK, Duijts L, Meurs JB van, Zelm MC van, Jaddoe VW, Nordeng H, Knudsen GP, Magnus P, Nystad W, Staff AC, Felix JF, Lyle R: **Cell type specific DNA methylation in cord blood: A 450K-reference data set and cell count-based validation of estimated cell type composition.** *Epigenetics* 2016, **11**:690–69810.1080/15592294.2016.1214782Available: <http://dx.doi.org/10.1080/15592294.2016.1214782>.
237. Leek JT, Storey JD: **Capturing Heterogeneity in Gene Expression Studies by Surrogate Variable Analysis.** *PLoS Genetics* 2007, **3**:e16110.1371/journal.pgen.0030161Available: <https://dx.plos.org/10.1371/journal.pgen.0030161>.
238. Teschendorff AE, Zhuang J, Widswendter M: **Independent surrogate variable analysis to deconvolve confounding factors in large-scale microarray profiling studies.** *Bioinformatics* 2011, **27**:1496–150510.1093/bioinformatics/btr171Available: <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btr171>.
239. McGregor K, Bernatsky S, Colmegna I, Hudson M, Pastinen T, Labbe A, Greenwood CMT: **An evaluation of methods correcting for cell-type heterogeneity in DNA methylation**

- studies.** *Genome Biology* 2016, **17**:8410.1186/s13059-016-0935-y Available: <http://genomebiology.biomedcentral.com/articles/10.1186/s13059-016-0935-y>.
240. Teschendorff AE, Zheng SC: **Cell-type deconvolution in epigenome-wide association studies: a review and recommendations.** *Epigenomics* 2017, **9**:757–76810.2217/epi-2016-0153 Available: <https://www.futuremedicine.com/doi/10.2217/epi-2016-0153>.
241. Lin X, Tan JYL, Teh AL, Lim IY, Liew SJ, MacIsaac JL, Chong YS, Gluckman PD, Kobor MS, Cheong CY, Karnani N: **Cell type-specific DNA methylation in neonatal cord tissue and cord blood: a 850K-reference panel and comparison of cell types.** *Epigenetics* 2018, **13**:941–95810.1080/15592294.2018.1522929 Available: <https://doi.org/10.1080/15592294.2018.1522929>.
242. Wang S: **Method to detect differentially methylated loci with case-control designs using Illumina arrays.** *Genetic Epidemiology* 2011, **35**:686–69410.1002/gepi.20619 Available: <http://doi.wiley.com/10.1002/gepi.20619>.
243. Tsai PC, Bell JT: **Power and sample size estimation for epigenome-wide association scans to detect differential DNA methylation.** *International Journal of Epidemiology* 2015, **44**:1429–144110.1093/ije/dyv041.
244. Graw S, Henn R, Thompson JA, Koestler DC: **PwrEWAS: A user-friendly tool for comprehensive power estimation for epigenome wide association studies (EWAS).** *BMC Bioinformatics* 2019, **20**:1–1110.1186/s12859-019-2804-7.
245. Cohen J: *Statisticsl Power Analysis for the Behavioural Sciences.* 2nd ed. New York: Lawrence Erlbaum Associates; 1988.
246. Joubert BR, Felix JF, Yousefi P, Bakulski KM, Just AC, Breton C, Reese SE, Markunas CA, Richmond RC, Xu C-J, Küpers LK, Oh SS, Hoyo C, Gruzieva O, Söderhäll C, Salas LA, Baïz N, Zhang H, Lepeule J, Ruiz C, Ligthart S, Wang T, Taylor JA, Duijts L, Sharp GC, Jankipersadsing SA, Nilsen RM, Vaez A, Fallin MD, Hu D, et al.: **DNA Methylation in Newborns and Maternal Smoking in Pregnancy: Genome-wide Consortium Meta-analysis.** *The American Journal of Human Genetics* 2016, **98**:680–69610.1016/j.ajhg.2016.02.019 Available: <http://www.ncbi.nlm.nih.gov/pubmed/27040690>.
247. Breton CV, Marsit CJ, Faustman E, Nadeau K, Goodrich JM, Dolinoy DC, Herbstman J, Holland N, LaSalle JM, Schmidt R, Yousefi P, Perera F, Joubert BR, Wiemels J, Taylor M, Yang IV, Chen R, Hew KM, Freeland DMH, Miller R, Murphy SK: **Small-Magnitude Effect Sizes in Epigenetic End Points are Important in Children’s Environmental Health Studies: The Children’s Environmental Health and Disease Prevention Research Center’s Epigenetics Working Group.** *Environmental Health Perspectives* 2017, **125**:511–52610.1289/EHP595 Available: <https://ehp.niehs.nih.gov/doi/10.1289/EHP595>.

248. Partridge L, Deelen J, Slagboom PE: **Facing up to the global challenges of ageing.** *Nature* 2018, **561**:45–5610.1038/s41586-018-0457-8 Available: <http://www.nature.com/articles/s41586-018-0457-8>.
249. Campisi J, Kapahi P, Lithgow GJ, Melov S, Newman JC, Verdin E: **From discoveries in ageing research to therapeutics for healthy ageing.** *Nature* 2019, **571**:183–19210.1038/s41586-019-1365-2 Available: <http://www.nature.com/articles/s41586-019-1365-2>.
250. Chuong EB, Elde NC, Feschotte C: **Regulatory activities of transposable elements: from conflicts to benefits.** *Nature Reviews Genetics* 2017, **18**:71–8610.1038/nrg.2016.139 Available: <http://www.nature.com/articles/nrg.2016.139>.
251. 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**:692–70210.4161/epi.6.6.16196 Available: <http://www.tandfonline.com/doi/abs/10.4161/epi.6.6.16196>.
252. Horvath S, Oshima J, Martin GM, Lu AT, Quach A, Cohen H, Felton S, Matsuyama M, Lowe D, Kabacik S, Wilson JG, Reiner AP, Maierhofer A, Flunkert J, Aviv A, Hou L, Baccarelli AA, Li Y, Stewart JD, Whitsel EA, Ferrucci L, Matsuyama S, Raj K: **Epigenetic clock for skin and blood cells applied to Hutchinson Gilford Progeria Syndrome and ex vivo studies.** *Aging* 2018, **10**:1758–177510.18632/aging.101508 Available: <http://www.aging-us.com/article/101508/text>.
253. Field AE, Robertson NA, Wang T, Havas A, Ideker T, Adams PD: **DNA Methylation Clocks in Aging: Categories, Causes, and Consequences.** *Molecular Cell* 2018, **71**:882–89510.1016/j.molcel.2018.08.008 Available: <https://doi.org/10.1016/j.molcel.2018.08.008>.
254. Bhattacharyya S, Varshney U: **Evolution of initiator tRNAs and selection of methionine as the initiating amino acid.** *RNA Biology* 2016, **13**:810–81910.1080/15476286.2016.1195943 Available: <http://dx.doi.org/10.1080/15476286.2016.1195943>.
255. Kamhi E, Raitskin O, Sperling R, Sperling J: **A potential role for initiator-tRNA in pre-mRNA splicing regulation.** *Proceedings of the National Academy of Sciences of the United States of America* 2010, **107**:11319–1132410.1073/pnas.0911561107.
256. Birch J, Clarke CJ, Campbell AD, Campbell K, Mitchell L, Liko D, Kalna G, Strathdee D, Sansom OJ, Neilson M, Blyth K, Norman JC: **The initiator methionine tRNA drives cell migration and invasion leading to increased metastatic potential in melanoma.** *Biology Open* 2016, **5**:1371–137910.1242/bio.019075.
257. Rideout EJ, Marshall L, Grewal SS: **Drosophila RNA polymerase III repressor Maf1 controls body size and developmental timing by modulating tRNAiMet synthesis**

- and systemic insulin signaling. *Proceedings of the National Academy of Sciences* 2012, **109**:1139–114410.1073/pnas.1113311109 Available: <http://www.pnas.org/cgi/doi/10.1073/pnas.1113311109>.
258. Kolitz SE, Lorsch JR: **Eukaryotic initiator tRNA: Finely tuned and ready for action.** *FEBS Letters* 2010, **584**:396–40410.1016/j.febslet.2009.11.047 Available: <http://doi.wiley.com/10.1016/j.febslet.2009.11.047>.
259. Pavon-Eternod M, Gomes S, Rosner MR, Pan T: **Overexpression of initiator methionine tRNA leads to global reprogramming of tRNA expression and increased proliferation in human epithelial cells.** *RNA* 2013, **19**:461–46610.1261/rna.037507.112 Available: <http://www.ncbi.nlm.nih.gov/pubmed/23431330>.
260. Eigen M, Lindemann B, Tietze M, Winkler-Oswatitsch R, Dress A, Haeseler A von: **How old is the genetic code? Statistical geometry of tRNA provides an answer.** *Science* 1989, **244**:673–67910.1126/science.2497522 Available: <http://www.sciencemag.org/cgi/doi/10.1126/science.2497522>.
261. Tavernarakis N: **Ageing and the regulation of protein synthesis: a balancing act?** *Trends in Cell Biology* 2008, **18**:228–23510.1016/j.tcb.2008.02.004 Available: <https://linkinghub.elsevier.com/retrieve/pii/S0962892408000792>.
262. Parisien M, Wang X, Pan T: **Diversity of human tRNA genes from the 1000-genomes project.** *RNA Biology* 2013, **10**:1853–186710.4161/rna.27361 Available: <http://www.tandfonline.com/doi/full/10.4161/rna.27361>.
263. Chan PP, Lowe TM: **GtRNAdB: a database of transfer RNA genes detected in genomic sequence.** *Nucleic acids research* 2009, **37**:D93–710.1093/nar/gkn787 Available: <http://www.ncbi.nlm.nih.gov/pubmed/18984615>.
264. Torres AG: **Enjoy the Silence: Nearly Half of Human tRNA Genes Are Silent.** *Bioinformatics and Biology Insights* 2019, **13**:11779322198684510.1177/1177932219868454 Available: <http://journals.sagepub.com/doi/10.1177/1177932219868454>.
265. Lodish H, Berk A, Zipursky SL, Matsudaira P, Baltimore D, Darnell J: *Molecular Cell Biology, 4th edition.* 4th ed. New York: W. H. Freeman; 2000 Available: <https://www.ncbi.nlm.nih.gov/books/NBK21475/>.
266. Schramm L: **Recruitment of RNA polymerase III to its target promoters.** *Genes & Development* 2002, **16**:2593–262010.1101/gad.1018902 Available: <http://www.genesdev.org/cgi/doi/10.1101/gad.1018902>.
267. Canella D, Praz V, Reina JH, Cousin P, Hernandez N: **Defining the RNA polymerase III transcriptome: Genome-wide localization of the RNA polymerase**

- III transcription machinery in human cells.** *Genome research* 2010, **20**:710–2110.1101/gr.101337.109 Available: <http://www.ncbi.nlm.nih.gov/pubmed/20413673>.
268. Dieci G, Sentenac A: **Facilitated Recycling Pathway for RNA Polymerase III.** *Cell* 1996, **84**:245–25210.1016/S0092-8674(00)80979-4 Available: <https://linkinghub.elsevier.com/retrieve/pii/S0092867400809794>.
269. Murawski M, Szczesniak B, Zoladek T, Hopper AK, Martin NC, Boguta M: **maf1 mutation alters the subcellular localization of the Mod5 protein in yeast.** *Acta biochimica Polonica* 1994, **41**:441–8 Available: <http://www.ncbi.nlm.nih.gov/pubmed/7732762>.
270. Pluta K, Lefebvre O, Martin NC, Smagowicz WJ, Stanford DR, Ellis SR, Hopper AK, Sentenac A, Boguta M: **Maf1p, a Negative Effector of RNA Polymerase III in *Saccharomyces cerevisiae*.** *Molecular and Cellular Biology* 2001, **21**:5031–504010.1128/MCB.21.15.5031-5040.2001 Available: <http://www.ncbi.nlm.nih.gov/pubmed/11438659>.
271. Vorländer MK, Baudin F, Moir RD, Wetzel R, Hagen WJH, Willis IM, Müller CW: **Structural basis for RNA polymerase III transcription repression by Maf1.** *Nature Structural & Molecular Biology* 2020, [10.1038/s41594-020-0383-y](https://doi.org/10.1038/s41594-020-0383-y) Available: <http://www.nature.com/articles/s41594-020-0383-y>.
272. Mange F, Praz V, Migliavacca E, Willis IM, Schütz F, Hernandez N: **Diurnal regulation of RNA polymerase III transcription is under the control of both the feeding–fasting response and the circadian clock.** *Genome Research* 2017, **27**:973–98410.1101/gr.217521.116 Available: <http://genome.cshlp.org/lookup/doi/10.1101/gr.217521.116>.
273. Kennedy BK, Lamming DW: **The Mechanistic Target of Rapamycin: The Grand ConducTOR of Metabolism and Aging.** *Cell Metabolism* 2016, **23**:990–100310.1016/j.cmet.2016.05.009 Available: <http://dx.doi.org/10.1016/j.cmet.2016.05.009>.
274. Crighton D: **p53 represses RNA polymerase III transcription by targeting TBP and inhibiting promoter occupancy by TFIIIB.** *The EMBO Journal* 2003, **22**:2810–282010.1093/emboj/cdg265 Available: <http://emboj.embopress.org/cgi/doi/10.1093/emboj/cdg265>.
275. Sutcliffe JE, Brown TRP, Allison SJ, Scott PH, White RJ: **Retinoblastoma Protein Disrupts Interactions Required for RNA Polymerase III Transcription.** *Molecular and Cellular Biology* 2000, **20**:9192–920210.1128/MCB.20.24.9192-9202.2000 Available: <http://mcb.asm.org/cgi/doi/10.1128/MCB.20.24.9192-9202.2000>.
276. Gomez-Roman N, Grandori C, Eisenman RN, White RJ: **Direct activation of RNA polymerase III transcription by c-Myc.** *Nature* 2003, **421**:290–29410.1038/nature01327 Available:

<http://www.nature.com/articles/nature01327>.

277. Krishnan P, Ghosh S, Wang B, Heyns M, Li D, Mackey JR, Kovalchuk O, Damaraju S: **Genome-wide profiling of transfer RNAs and their role as novel prognostic markers for breast cancer.** *Scientific Reports* 2016, **6**:3284310.1038/srep32843Available: <http://dx.doi.org/10.1038/srep32843>.
278. Huang S-q, Sun B, Xiong Z-p, Shu Y, Zhou H-h, Zhang W, Xiong J, Li Q: **The dysregulation of tRNAs and tRNA derivatives in cancer.** *Journal of Experimental & Clinical Cancer Research* 2018, **37**:10110.1186/s13046-018-0745-zAvailable: <https://jeccr.biomedcentral.com/articles/10.1186/s13046-018-0745-z>.
279. Besser D, Götz F, Schulze-Forster K, Wagner H, Kröger H, Simon D: **DNA methylation inhibits transcription by RNA polymerase III of a tRNA gene, but not of a 5S rRNA gene.** *FEBS letters* 1990, **269**:358–6210.1016/0014-5793(90)81193-RAvailable: <http://www.ncbi.nlm.nih.gov/pubmed/2401361>.
280. Varshney D, Vavrova-Anderson J, Oler AJ, Cowling VH, Cairns BR, White RJ: **SINE transcription by RNA polymerase III is suppressed by histone methylation but not by DNA methylation.** *Nature communications* 2015, **6**:656910.1038/ncomms7569Available: <http://www.ncbi.nlm.nih.gov/pubmed/25798578>.
281. Gerber A, Ito K, Chu C-S, Roeder RG: **Gene-Specific Control of tRNA Expression by RNA Polymerase II.** *Molecular Cell* 2020, **78**:765–778.e710.1016/j.molcel.2020.03.023Available: <https://doi.org/10.1016/j.molcel.2020.03.023>.
282. Gingold H, Dahan O, Pilpel Y: **Dynamic changes in translational efficiency are deduced from codon usage of the transcriptome.** *Nucleic acids research* 2012, **40**:10053–6310.1093/nar/gks772Available: <http://www.ncbi.nlm.nih.gov/pubmed/22941644>.
283. Athey J, Alexaki A, Osipova E, Rostovtsev A, Santana-Quintero LV, Katneni U, Simonyan V, Kimchi-Sarfaty C: **A new and updated resource for codon usage tables.** *BMC Bioinformatics* 2017, **18**:39110.1186/s12859-017-1793-7Available: <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-017-1793-7>.
284. Dittmar KA, Goodenbour JM, Pan T: **Tissue-specific differences in human transfer RNA expression.** *PLoS genetics* 2006, **2**:e22110.1371/journal.pgen.0020221Available: <http://www.ncbi.nlm.nih.gov/pubmed/17194224>.
285. Sagi D, Rak R, Gingold H, Adir I, Maayan G, Dahan O, Broday L, Pilpel Y, Rechavi O: **Tissue- and Time-Specific Expression of Otherwise Identical tRNA Genes.** *PLOS Genetics* 2016, **12**:e100626410.1371/journal.pgen.1006264Available: <http://www.ncbi.nlm.nih.gov/pubmed/27560950>.
286. Kirchner S, Cai Z, Rauscher R, Kastelic N, Anding M, Czech A, Kleizen B, Ost-

- edgaard LS, Braakman I, Sheppard DN, Ignatova Z: **Alteration of protein function by a silent polymorphism linked to tRNA abundance.** *PLOS Biology* 2017, **15**:e200077910.1371/journal.pbio.2000779 Available: <http://dx.plos.org/10.1371/journal.pbio.2000779>.
287. Ishimura R, Nagy G, Dotu I, Zhou H, Yang X-L, Schimmel P, Senju S, Nishimura Y, Chuang JH, Ackerman SL: **RNA function. Ribosome stalling induced by mutation of a CNS-specific tRNA causes neurodegeneration.** *Science (New York, N.Y.)* 2014, **345**:455–910.1126/science.1249749 Available: <http://www.ncbi.nlm.nih.gov/pubmed/25061210>.
288. Pliatsika V, Loher P, Magee R, Telonis AG, Londin E, Shigematsu M, Kirino Y, Rigoutsos I: **MINTbase v2.0: a comprehensive database for tRNA-derived fragments that includes nuclear and mitochondrial fragments from all The Cancer Genome Atlas projects.** *Nucleic Acids Research* 2018, **46**:D152–D15910.1093/nar/gkx1075 Available: <http://academic.oup.com/nar/article/46/D1/D152/4653530>.
289. Lee YS, Shibata Y, Malhotra A, Dutta A: **A novel class of small RNAs: tRNA-derived RNA fragments (tRFs).** *Genes & development* 2009, **23**:2639–4910.1101/gad.1837609 Available: <http://www.ncbi.nlm.nih.gov/pubmed/19933153>.
290. Torres AG, Reina O, Stephan-Otto Attolini C, Ribas de Pouplana L: **Differential expression of human tRNA genes drives the abundance of tRNA-derived fragments.** *Proceedings of the National Academy of Sciences* 2019, **116**:20182112010.1073/pnas.1821120116 Available: <http://www.pnas.org/lookup/doi/10.1073/pnas.1821120116>.
291. Li S, Xu Z, Sheng J: **tRNA-Derived Small RNA: A Novel Regulatory Small Non-Coding RNA.** *Genes* 2018, **9**:24610.3390/genes9050246 Available: <http://www.mdpi.com/2073-4425/9/5/246>.
292. Xu W-L, Yang Y, Wang Y-D, Qu L-H, Zheng L-L: **Computational Approaches to tRNA-Derived Small RNAs.** *Non-Coding RNA* 2017, **3**:210.3390/ncrna3010002 Available: <http://www.mdpi.com/2311-553X/3/1/2>.
293. Keam SP, Young PE, McCorkindale AL, Dang THY, Clancy JL, Humphreys DT, Preiss T, Hutvagner G, Martin DIK, Cropley JE, Suter CM: **The human Piwi protein Hiwi2 associates with tRNA-derived piRNAs in somatic cells.** *Nucleic Acids Research* 2014, **42**:8984–899510.1093/nar/gku620 Available: <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gku620>.
294. Honda S, Kawamura T, Loher P, Morichika K, Rigoutsos I, Kirino Y: **The biogenesis pathway of tRNA-derived piRNAs in Bombyx germ cells.** *Nucleic Acids Research* 2017, **45**:9108–912010.1093/nar/gkx537 Available: <http://academic.oup.com/nar/article/45/15/9108/3883741/The-biogenesis-pathway-of-tRNADerived-piRNAs-in>.

295. Tosar JP, Rovira C, Cayota A: **Non-coding RNA fragments account for the majority of annotated piRNAs expressed in somatic non-gonadal tissues.** *Communications Biology* 2018, **1**:210.1038/s42003-017-0001-7 Available: <http://www.nature.com/articles/s42003-017-0001-7>.
296. Loher P, Telonis AG, Rigoutsos I: **MINTmap: fast and exhaustive profiling of nuclear and mitochondrial tRNA fragments from short RNA-seq data.** *Scientific reports* 2017, **7**:4118410.1038/srep41184 Available: <http://www.ncbi.nlm.nih.gov/pubmed/28220888>.
297. Kim HK, Fuchs G, Wang S, Wei W, Zhang Y, Park H, Roy-Chaudhuri B, Li P, Xu J, Chu K, Zhang F, Chua M-S, So S, Zhang QC, Sarnow P, Kay MA: **A transfer-RNA-derived small RNA regulates ribosome biogenesis.** *Nature* 2017, **552**:57–6210.1038/nature25005 Available: <http://www.nature.com/articles/nature25005>.
298. Martinez G, Choudury SG, Slotkin RK: **tRNA-derived small RNAs target transposable element transcripts.** *Nucleic Acids Research* 2017, **45**:5142–515210.1093/nar/gkx103 Available: <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkx103>.
299. Schimmel P: **The emerging complexity of the tRNA world: mammalian tRNAs beyond protein synthesis.** *Nature reviews. Molecular cell biology* 2018, **19**:45–5810.1038/nrm.2017.77 Available: <http://www.ncbi.nlm.nih.gov/pubmed/28875994>.
300. Cristodero M, Polacek N: **The multifaceted regulatory potential of tRNA-derived fragments.** *Non-coding RNA Investigation* 2017, **1**:7–710.21037/ncri.2017.08.07 Available: <http://ncri.amegroups.com/article/view/3820/4459>.
301. Pace DA: **Fixed metabolic costs for highly variable rates of protein synthesis in sea urchin embryos and larvae.** *Journal of Experimental Biology* 2006, **209**:158–17010.1242/jeb.01962 Available: <http://jeb.biologists.org/cgi/doi/10.1242/jeb.01962>.
302. Nwanaji-Enwerem JC, Weisskopf MG, Baccarelli AA: **Multi-tissue DNA methylation age: Molecular relationships and perspectives for advancing biomarker utility.** *Ageing Research Reviews* 2018, **45**:15–2310.1016/j.arr.2018.04.005 Available: <https://doi.org/10.1016/j.arr.2018.04.005>.
303. Hansen M, Taubert S, Crawford D, Libina N, Lee S-J, Kenyon C: **Lifespan extension by conditions that inhibit translation in *Caenorhabditis elegans*.** *Aging Cell* 2007, **6**:95–11010.1111/j.1474-9726.2006.00267.x Available: <http://doi.wiley.com/10.1111/j.1474-9726.2006.00267.x>.
304. Filer D, Thompson MA, Takhayev V, Dobson AJ, Kotronaki I, Green JWM, Heinemann M, Tullet JMA, Alic N: **RNA polymerase III limits longevity downstream of TORC1.**

*Nature* 2017, **552**:263–26710.1038/nature25007Available: <http://www.nature.com/articles/nature25007>.

305. Dhahbi JM, Spindler SR, Atamna H, Yamakawa A, Boffelli D, Mote P, Martin DI: **5' tRNA halves are present as abundant complexes in serum, concentrated in blood cells, and modulated by aging and calorie restriction.** *BMC Genomics* 2013, **14**:29810.1186/1471-2164-14-298Available: <http://www.ncbi.nlm.nih.gov/pubmed/23638709>.

306. Yoshihisa T: **Handling tRNA introns, archaeal way and eukaryotic way.** *Frontiers in genetics* 2014, **5**:21310.3389/fgene.2014.00213Available: <http://www.ncbi.nlm.nih.gov/pubmed/25071838>.

307. Bhargava P: **Epigenetic regulation of transcription by RNA polymerase III.** *Biochimica et biophysica acta* 2013, **1829**:1015–2510.1016/j.bbagr.2013.05.005Available: <http://www.ncbi.nlm.nih.gov/pubmed/23732820>.

308. Wilusz JE: **Controlling translation via modulation of tRNA levels.** *Wiley Interdisciplinary Reviews: RNA* 2015, **6**:453–47010.1002/wrna.1287Available: <http://doi.wiley.com/10.1002/wrna.1287>.

309. Cozen AE, Quartley E, Holmes AD, Hrabetá-Robinson E, Phizicky EM, Lowe TM: **ARM-seq: AlkB-facilitated RNA methylation sequencing reveals a complex landscape of modified tRNA fragments.** *Nature Methods* 2015, **12**:879–88410.1038/nmeth.3508Available: <http://www.nature.com/articles/nmeth.3508>.

310. Shigematsu M, Honda S, Loher P, Telonis AG, Rigoutsos I, Kirino Y: **YAMAT-seq: an efficient method for high-throughput sequencing of mature transfer RNAs.** *Nucleic acids research* 2017, **45**:e7010.1093/nar/gkx005Available: <http://www.ncbi.nlm.nih.gov/pubmed/28108659>.

311. Gogakos T, Brown M, Garzia A, Meyer C, Hafner M, Tuschl T: **Characterizing Expression and Processing of Precursor and Mature Human tRNAs by Hydro-tRNaseq and PAR-CLIP.** *Cell Reports* 2017, **20**:1463–147510.1016/j.celrep.2017.07.029Available: <http://dx.doi.org/10.1016/j.celrep.2017.07.029>.

312. Raab JR, Chiu J, Zhu J, Katzman S, Kurukuti S, Wade PA, Haussler D, Kamakaka RT: **Human tRNA genes function as chromatin insulators.** *The EMBO journal* 2012, **31**:330–5010.1038/embj.2011.406Available: <http://www.ncbi.nlm.nih.gov/pubmed/22085927>.

313. Van Bortle K, Phanstiel DH, Snyder MP: **Topological organization and dynamic regulation of human tRNA genes during macrophage differentiation.** *Genome Biology* 2017, **18**:18010.1186/s13059-017-1310-3Available: <http://genomebiology.biomedcentral.com/articles/10.1186/s13059-017-1310-3>.

314. Noma K-i, Cam HP, Maraia RJ, Grewal SIS: **A role for TFIIIC transcription factor**

- complex in genome organization.** *Cell* 2006, **125**:859–7210.1016/j.cell.2006.04.028 Available: <http://www.ncbi.nlm.nih.gov/pubmed/16751097>.
315. Syddall H, Aihie Sayer A, Dennison E, Martin H, Barker D, Cooper C: **Cohort Profile: The Hertfordshire Cohort Study.** *International Journal of Epidemiology* 2005, **34**:1234–124210.1093/ije/dyi127 Available: <http://academic.oup.com/ije/article/34/6/1234/707357/Cohort-Profile-The-Hertfordshire-Cohort-Study>.
316. Amemiya HM, Kundaje A, Boyle AP: **The ENCODE Blacklist: Identification of Problematic Regions of the Genome.** *Scientific Reports* 2019, **9**:935410.1038/s41598-019-45839-z Available: <http://www.nature.com/articles/s41598-019-45839-z>.
317. Derrien T, Estellé J, Marco Sola S, Knowles DG, Raineri E, Guigó R, Ribeca P: **Fast Computation and Applications of Genome Mappability.** *PLoS ONE* 2012, **7**:e3037710.1371/journal.pone.0030377 Available: <https://dx.plos.org/10.1371/journal.pone.0030377>.
318. Ernst J, Kheradpour P, Mikkelsen TS, Shoresh N, Ward LD, Epstein CB, Zhang X, Wang L, Issner R, Coyne M, Ku M, Durham T, Kellis M, Bernstein BE: **Mapping and analysis of chromatin state dynamics in nine human cell types.** *Nature* 2011, **473**:43–4910.1038/nature09906 Available: <http://www.nature.com/articles/nature09906>.
319. Bell CG, Wilson GA, Butcher LM, Roos C, Walter L, Beck S: **Human-specific CpG "beacons" identify loci associated with human-specific traits and disease.** *Epigenetics* 2012, **7**:1188–9910.4161/epi.22127 Available: <http://www.ncbi.nlm.nih.gov/pubmed/22968434>.
320. North BV, Curtis D, Sham PC: **A Note on the Calculation of Empirical P Values from Monte Carlo Procedures.** *The American Journal of Human Genetics* 2003, **72**:498–49910.1086/346173 Available: <https://linkinghub.elsevier.com/retrieve/pii/S0002929707605618>.
321. Ewels P, Magnusson M, Lundin S, Käller M: **MultiQC: summarize analysis results for multiple tools and samples in a single report.** *Bioinformatics (Oxford, England)* 2016, **32**:3047–810.1093/bioinformatics/btw354 Available: <http://www.ncbi.nlm.nih.gov/pubmed/27312411>.
322. Krueger F: **Trim Galore.** 2015, Available: [https://www.bioinformatics.babraham.ac.uk/projects/trim%7B/\\_%7Dgalore/](https://www.bioinformatics.babraham.ac.uk/projects/trim%7B/_%7Dgalore/).
323. Martin M: **Cutadapt removes adapter sequences from high-throughput sequencing reads.** *EMBnet.journal* 2011, **17**:1010.14806/ej.17.1.200 Available: <http://journal.embnet.org/index.php/embnetjournal/article/view/200>.
324. Langmead B, Salzberg SL: **Fast gapped-read alignment with Bowtie 2.** *Nature Methods* 2012, **9**:357–35910.1038/nmeth.1923 Available: <http://www.ncbi.nlm.nih.gov/pubmed/22388286>.

325. Müller F, Scherer M, Assenov Y, Lutsik P, Walter J, Lengauer T, Bock C: **RnBeads 2.0: comprehensive analysis of DNA methylation data.** *Genome Biology* 2019, **20**:5510.1186/s13059-019-1664-9 Available: <https://genomebiology.biomedcentral.com/articles/10.1186/s13059-019-1664-9>.
326. Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, Smyth GK: **limma powers differential expression analyses for RNA-sequencing and microarray studies.** *Nucleic Acids Research* 2015, **43**:e47–e4710.1093/nar/gkv007 Available: <http://www.ncbi.nlm.nih.gov/pubmed/25605792>.
327. Meuleman W: **Epilogos.** 2019, Available: <https://epilogos.altius.org/> <https://github.com/Altius/epilogos>.
328. Reinius LE, Acevedo N, Joerink M, Pershagen G, Dahlén S-E, Greco D, Söderhäll 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**:e4136110.1371/journal.pone.0041361 Available: <https://dx.plos.org/10.1371/journal.pone.0041361>.
329. Sean D, Meltzer PS: **GEOquery: A bridge between the Gene Expression Omnibus (GEO) and BioConductor.** *Bioinformatics* 2007, **23**:1846–184710.1093/bioinformatics/btm254.
330. Grossman RL, Heath AP, Ferretti V, Varmus HE, Lowy DR, Kibbe WA, Staudt LM: **Toward a Shared Vision for Cancer Genomic Data.** *New England Journal of Medicine* 2016, **375**:1109–111210.1056/NEJMp1607591 Available: <http://www.nejm.org/doi/10.1056/NEJMp1607591>.
331. Yang Z, Wong A, Kuh D, Paul DS, Rakyan VK, Leslie RD, Zheng SC, Widschwendter M, Beck S, Teschendorff AE: **Correlation of an epigenetic mitotic clock with cancer risk.** *Genome Biology* 2016, **17**:20510.1186/s13059-016-1064-3 Available: <http://genomebiology.biomedcentral.com/articles/10.1186/s13059-016-1064-3>.
332. Nazor KL, Altun G, Lynch C, Tran H, Harness JV, Slavin I, Garitaonandia I, Müller F-J, Wang Y-C, Boscolo FS, Fakunle E, Dumevska B, Lee S, Park HS, Olee T, D'Lima DD, Semeckin R, Parast MM, Galat V, Laslett AL, Schmidt U, Keirstead HS, Loring JF, Laurent LC: **Recurrent Variations in DNA Methylation in Human Pluripotent Stem Cells and Their Differentiated Derivatives.** *Cell Stem Cell* 2012, **10**:620–63410.1016/j.stem.2012.02.013 Available: <http://www.ncbi.nlm.nih.gov/pubmed/10449618>.
333. Juzenas S, Venkatesh G, Hüenthal M, Hoeppner MP, Du ZG, Paulsen M, Rosenstiel P, Senger P, Hofmann-Apitius M, Keller A, Kupcinskas L, Franke A, Hemmrich-Stanisak G: **A comprehensive, cell specific microRNA catalogue of human peripheral blood.** *Nucleic Acids Research* 2017, **45**:9290–930110.1093/nar/gkx706 Available: <http://academic.oup.com/nar>

/article/45/16/9290/4080663.

334. Petkovich DA, Podolskiy DI, Lobanov AV, Lee S-G, Miller RA, Gladyshev VN: **Using DNA Methylation Profiling to Evaluate Biological Age and Longevity Interventions.** *Cell Metabolism* 2017, **25**:954–960.e610.1016/j.cmet.2017.03.016 Available: <https://linkinghub.elsevier.com/retrieve/pii/S1550413117301687>.

335. Lowe TM, Chan PP: **tRNAscan-SE On-line: integrating search and context for analysis of transfer RNA genes.** *Nucleic Acids Research* 2016, **44**:W54–W5710.1093/nar/gkw413 Available: <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkw413>.

336. Schmitt BM, Rudolph KLM, Karagianni P, Fonseca NA, White RJ, Talianidis I, Odom DT, Marioni JC, Kutter C: **High-resolution mapping of transcriptional dynamics across tissue development reveals a stable mRNA–tRNA interface.** *Genome Research* 2014, **24**:1797–180710.1101/gr.176784.114 Available: <http://www.ncbi.nlm.nih.gov/pubmed/25122613>.

337. Gu Z, Eils R, Schlesner M: **Complex heatmaps reveal patterns and correlations in multidimensional genomic data.** *Bioinformatics* 2016, **32**:2847–284910.1093/bioinformatics/btw313 Available: <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btw313>.

338. Thornlow BP, Armstrong J, Holmes AD, Howard JM, Corbett-Detig RB, Lowe TM: **Predicting transfer RNA gene activity from sequence and genome context.** *Genome Research* 2020, **30**:85–9410.1101/gr.256164.119.

339. Ehrlich M: **DNA hypermethylation in disease: mechanisms and clinical relevance.** *Epigenetics* 2019, **14**:1141–116310.1080/15592294.2019.1638701 Available: <http://www.ncbi.nlm.nih.gov/pubmed/31284823>.

340. Xu Z, Taylor JA: **Genome-wide age-related DNA methylation changes in blood and other tissues relate to histone modification, expression and cancer.** *Carcinogenesis* 2014, **35**:356–36410.1093/carcin/bgt391 Available: <https://academic.oup.com/carcin/article-lookup/doi/10.1093/carcin/bgt391>.

341. Slieker RC, Relton CL, Gaunt TR, Slagboom PE, Heijmans BT: **Age-related DNA methylation changes are tissue-specific with ELOVL2 promoter methylation as exception.** *Epigenetics and Chromatin* 2018, **11**:1–1110.1186/s13072-018-0191-3 Available: <https://doi.org/10.1186/s13072-018-0191-3>.

342. Zhu T, Zheng SC, Paul DS, Horvath S, Teschendorff AE: **Cell and tissue type independent age-associated DNA methylation changes are not rare but common.** *Aging* 2018, **10**:3541–355710.18632/aging.101666 Available: <http://www.ncbi.nlm.nih.gov/pubmed/30482885>.

343. Goodenbour JM, Pan T: **Diversity of tRNA genes in eukaryotes.** *Nucleic acids research* 2006, **34**:6137–4610.1093/nar/gkl725 Available: <http://www.ncbi.nlm.nih.gov/pubmed/17088292>.
344. Geslain R, Pan T: **Functional Analysis of Human tRNA Isodecoders.** *Journal of Molecular Biology* 2010, **396**:821–83110.1016/j.jmb.2009.12.018 Available: <https://linkinghub.elsvier.com/retrieve/pii/S002228360901523X>.
345. Li S, Shi X, Chen M, Xu N, Sun D, Bai R, Chen H, Ding K, Sheng J, Xu Z: **Angiogenin promotes colorectal cancer metastasis via tiRNA production.** *International Journal of Cancer* 2019, :ijc.3224510.1002/ijc.32245 Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/ijc.32245>.
346. Plotkin JB, Robins H, Levine AJ: **Tissue-specific codon usage and the expression of human genes.** *Proceedings of the National Academy of Sciences of the United States of America* 2004, **101**:12588–9110.1073/pnas.0404957101 Available: <http://www.ncbi.nlm.nih.gov/pubmed/15314228>.
347. Schmitt BM, Rudolph KLM, Karagianni P, Fonseca NA, White RJ, Talianidis I, Odom DT, Marioni JC, Kutter C: **High-resolution mapping of transcriptional dynamics across tissue development reveals a stable mRNA-tRNA interface.** *Genome research* 2014, **24**:1797–80710.1101/gr.176784.114 Available: <http://www.ncbi.nlm.nih.gov/pubmed/25122613>.
348. Powell JR, Moriyama EN: **Evolution of codon usage bias in Drosophila.** *Proceedings of the National Academy of Sciences* 1997, **94**:7784–779010.1073/pnas.94.15.7784 Available: <http://www.pnas.org/cgi/doi/10.1073/pnas.94.15.7784>.
349. Rudolph KLM, Schmitt BM, Villar D, White RJ, Marioni JC, Kutter C, Odom DT: **Codon-Driven Translational Efficiency Is Stable across Diverse Mammalian Cell States.** *PLoS genetics* 2016, **12**:e100602410.1371/journal.pgen.1006024 Available: <http://www.ncbi.nlm.nih.gov/pubmed/27166679>.
350. Gu T, Lin X, Cullen SM, Luo M, Jeong M, Estecio M, Shen J, Hardikar S, Sun D, Su J, Rux D, Guzman A, Lee M, Qi LS, Chen J-J, Kyba M, Huang Y, Chen T, Li W, Goodell MA: **DNMT3A and TET1 cooperate to regulate promoter epigenetic landscapes in mouse embryonic stem cells.** *Genome Biology* 2018, **19**:8810.1186/s13059-018-1464-7 Available: <http://genomebiology.biomedcentral.com/articles/10.1186/s13059-018-1464-7>.
351. Wang Y, Brady KS, Caiello BP, Ackerson SM, Stewart JA: **Human CST suppresses origin licensing and promotes AND-1/Ctf4 chromatin association.** *Life Science Alliance* 2019, **2**:e20180027010.26508/lsc.201800270 Available: <http://www.life-science-alliance.org/lookup/doi/10.26508/lsc.201800270>.
352. Sargolzaeiaval F, Zhang J, Schleit J, Lessel D, Kubisch C, Precioso DR, Sillence D, Hisama

- FM, Dorschner M, Martin GM, Oshima J: **CTC1 mutations in a Brazilian family with progeroid features and recurrent bone fractures.** *Molecular Genetics & Genomic Medicine* 2018, **6**:1148–115610.1002/mgg3.495 Available: <http://doi.wiley.com/10.1002/mgg3.495>.
353. Thornlow BP, Hough J, Roger JM, Gong H, Lowe TM, Corbett-Detig RB: **Transfer RNA genes experience exceptionally elevated mutation rates.** *Proceedings of the National Academy of Sciences* 2018, **115**:8996–900110.1073/pnas.1801240115 Available: <http://www.pnas.org/lookup/doi/10.1073/pnas.1801240115>.
354. Iben JR, Maraia RJ: **tRNA gene copy number variation in humans.** *Gene* 2014, **536**:376–38410.1016/j.gene.2013.11.049 Available: <https://linkinghub.elsevier.com/retrieve/pii/S0378111913015758>.
355. Darrow EM, Chadwick BP: **A novel tRNA variable number tandem repeat at human chromosome 1q23.3 is implicated as a boundary element based on conservation of a CTCF motif in mouse.** *Nucleic Acids Research* 2014, **42**:6421–643510.1093/nar/gku280 Available: <http://academic.oup.com/nar/article/42/10/6421/2435926/A-novel-tRNA-variable-number-tandem-repeat-at>.
356. Müller CA, Nieduszynski CA: **DNA replication timing influences gene expression level.** *The Journal of Cell Biology* 2017, **216**:1907–191410.1083/jcb.201701061 Available: <http://www.jcb.org/lookup/doi/10.1083/jcb.201701061>.
357. Du Q, Bert SA, Armstrong NJ, Caldon CE, Song JZ, Nair SS, Gould CM, Luu P-L, Peters T, Khouri A, Qu W, Zotenko E, Stirzaker C, Clark SJ: **Replication timing and epigenome remodelling are associated with the nature of chromosomal rearrangements in cancer.** *Nature Communications* 2019, **10**:41610.1038/s41467-019-08302-1 Available: <http://dx.doi.org/10.1038/s41467-019-08302-1>.
358. Cruickshanks HA, McBryan T, Nelson DM, VanderKraats ND, Shah PP, Tuyn J van, Singh Rai T, Brock C, Donahue G, Dunican DS, Drotar ME, Meehan RR, Edwards JR, Berger SL, Adams PD: **Senescent cells harbour features of the cancer epigenome.** *Nature Cell Biology* 2013, **15**:1495–150610.1038/ncb2879 Available: <http://www.ncbi.nlm.nih.gov/pubmed/24270890>.
359. Sun L, Yu R, Dang W: **Chromatin Architectural Changes during Cellular Senescence and Aging.** *Genes* 2018, **9**:21110.3390/genes9040211 Available: <http://www.mdpi.com/2073-4425/9/4/211>.
360. Moskowitz DM, Zhang DW, Hu B, Le Saux S, Yanes RE, Ye Z, Buenrostro JD, Weyand CM, Greenleaf WJ, Goronzy JJ: **Epigenomics of human CD8 T cell differentiation and aging.** *Science Immunology* 2017, **2**:eaag019210.1126/sciimmunol.aag0192 Available: <https://immunology.science.org/lookup/doi/10.1126/sciimmunol.aag0192>.

361. Hashimoto K, Kouno T, Ikawa T, Hayatsu N, Miyajima Y, Yabukami H, Terooatea T, Sasaki T, Suzuki T, Valentine M, Pasarella G, Okazaki Y, Suzuki H, Shin JW, Minoda A, Taniguchi I, Okano H, Arai Y, Hirose N, Carninci P: **Single-cell transcriptomics reveals expansion of cytotoxic CD4 T cells in supercentenarians.** *Proceedings of the National Academy of Sciences of the United States of America* 2019, **116**:24242–24251<sup>10.1073/pnas.1907883116</sup>.
362. Deniz Ö, Frost JM, Branco MR: **Regulation of transposable elements by DNA modifications.** *Nature Reviews Genetics* 2019, **20**:417–431<sup>10.1038/s41576-019-0106-6</sup>.
363. Kazazian HH, Moran JV: **Mobile DNA in health and disease.** *New England Journal of Medicine* 2017, **377**:361–370<sup>10.1056/NEJMra1510092</sup>.
364. Gregory TR: **Synergy between sequence and size in large-scale genomics.** *Nature Reviews Genetics* 2005, **6**:699–708<sup>10.1038/nrg1674</sup>.
365. Deaton AM, Bird A: **CpG islands and the regulation of transcription.** *Genes & development* 2011, **25**:1010–22<sup>10.1101/gad.2037511</sup> Available: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3137511/>
366. Hellman A, Chess A: **Gene body-specific methylation on the active X chromosome.** *Science* 2007, **315**:1141–1143<sup>10.1126/science.1136352</sup>.
367. Shukla S, Kavak E, Gregory M, Imashimizu M, Shutinoski B, Kashlev M, Oberdoerffer P, Sandberg R, Oberdoerffer S: **CTCF-promoted RNA polymerase II pausing links DNA methylation to splicing.** *Nature* 2011, **479**:74–79<sup>10.1038/nature10442</sup> Available: <http://dx.doi.org/10.1038/nature10442>.
368. Blattler A, Yao L, Witt H, Guo Y, Nicolet CM, Berman BP, Farnham PJ: **Global loss of DNA methylation uncovers intronic enhancers in genes showing expression changes.** *Genome Biology* 2014, **15**:469<sup>10.1186/s13059-014-0469-0</sup> Available: <http://genomebiology.biomedcentral.com/articles/10.1186/s13059-014-0469-0>.
369. Gama-Sosa MA, Slagel VA, Trewyn RW, Oxenhacker R, Kuo KC, Gehrke CW, Ehrlich M: **The 5-methylcytosine content of DNA from human tumors.** *Nucleic Acids Research* 1983, **11**:6883–6894<sup>10.1093/nar/11.19.6883</sup> Available: <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/11.19.6883>.
370. Bollati V, Schwartz J, Wright R, Litonjua A, Tarantini L, Suh H, Sparrow D, Vokonas P, Baccarelli A: **Decline in genomic DNA methylation through aging in a cohort of elderly subjects.** *Mechanisms of Ageing and Development* 2009, **130**:234–239<sup>10.1016/j.mad.2008.12.003</sup> Available: <https://linkinghub.elsevier.com/retrieve/pii/S0047637408002820>.
371. Dewannieux M, Esnault C, Heidmann T: **LINE-mediated retrotransposition of marked Alu sequences.** *Nature Genetics* 2003, **35**:41–48<sup>10.1038/ng1223</sup> Available:

<http://www.nature.com/articles/ng1223>.

372. Deininger P: **Alu elements: know the SINEs.** *Genome Biology* 2011, **12**:23610.1186/gb-2011-12-12-236Available: <http://genomebiology.biomedcentral.com/articles/10.1186/gb-2011-12-12-236>.
373. Ullu E, Tschudi C: **Alu sequences are processed 7SL RNA genes.** *Nature* 1984, **312**:171–17210.1038/312171a0Available: <http://www.nature.com/articles/312171a0>.
374. Quentin Y: **Origin of the Alu family:** *Nucleic Acids Research* 1992, **20**:3397–3401.
375. Chen L-L, Carmichael GG: **Gene regulation by SINES and inosines: biological consequences of A-to-I editing of Alu element inverted repeats.** *Cell Cycle* 2008, **7**:3294–330110.4161/cc.7.21.6927Available: <http://www.tandfonline.com/doi/abs/10.4161/cc.7.21.6927>.
376. Batzer MA, Deininger PL: **Alu repeats and human genomic diversity.** *Nature Reviews Genetics* 2002, **3**:370–37910.1038/nrg798Available: <http://www.nature.com/articles/nrg798>.
377. Ade C, Roy-Engel AM, Deininger PL: **Alu elements: an intrinsic source of human genome instability.** *Current Opinion in Virology* 2013, **3**:639–64510.1016/j.coviro.2013.09.002Available: <https://linkinghub.elsevier.com/retrieve/pii/S1879625713001533>.
378. Cordaux R, Batzer MA: **The impact of retrotransposons on human genome evolution.** *Nature Reviews Genetics* 2009, **10**:691–70310.1038/nrg2640Available: <http://www.nature.com/articles/nrg2640>.
379. Belancio VP, Hedges DJ, Deininger P: **Mammalian non-LTR retrotransposons: For better or worse, in sickness and in health.** *Genome Research* 2008, **18**:343–35810.1101/gr.5558208Available: <http://www.genome.org/cgi/doi/10.1101/gr.5558208>.
380. Jordà M, Díez-Villanueva A, Mallona I, Martín B, Lois S, Barrera V, Esteller M, Vavouri T, Peinado MA: **The epigenetic landscape of Alu repeats delineates the structural and functional genomic architecture of colon cancer cells.** *Genome Research* 2017, **27**:118–13210.1101/gr.207522.116Available: <http://www.ncbi.nlm.nih.gov/pubmed/27999094>.
381. Chen L-L, Yang L: **ALU ternative Regulation for Gene Expression.** *Trends in Cell Biology* 2017, **27**:480–49010.1016/j.tcb.2017.01.002Available: <https://linkinghub.elsevier.com/retrieve/pii/S0962892417300028>.
382. Ferrari R, de Llobet Cucalon LI, Di Vona C, Le Dilly F, Vidal E, Lioutas A, Oliete JQ, Jochem L, Cutts E, Dieci G, Vannini A, Teichmann M, Luna S de la, Beato M: **TFIIC Binding to Alu Elements Controls Gene Expression via Chromatin Looping and Histone Acetylation.** *Molecular Cell* 2019,:1–1310.1016/j.molcel.2019.10.020Available: <https://linkinghub.elsevier.com/retrieve/pii/S1097276519307981>.
383. Polak P, Domany E: **Alu elements contain many binding sites for transcription**

**factors and may play a role in regulation of developmental processes.** *BMC genomics* 2006, **7**:13310.1186/1471-2164-7-133 Available: <http://www.ncbi.nlm.nih.gov/pubmed/16740159>.

384. Su M, Han D, Boyd-Kirkup J, Yu X, Han J-DJ: **Evolution of Alu Elements toward Enhancers.** *Cell Reports* 2014, **7**:376–38510.1016/j.celrep.2014.03.011 Available: <http://dx.doi.org/10.1016/j.celrep.2014.03.011>.

385. Ward MC, Wilson MD, Barbosa-Morais NL, Schmidt D, Stark R, Pan Q, Schwalie PC, Menon S, Lukk M, Watt S, Thybert D, Kutter C, Kirschner K, Fliceck P, Blencowe BJ, Odom DT: **Latent Regulatory Potential of Human-Specific Repetitive Elements.** *Molecular Cell* 2013, **49**:262–27210.1016/j.molcel.2012.11.013 Available: <http://dx.doi.org/10.1016/j.molcel.2012.11.013>.

386. Xie M, Hong C, Zhang B, Lowdon RF, Xing X, Li D, Zhou X, Lee HJ, Maire CL, Ligon KL, Gascard P, Sigaroudinia M, Tlsty TD, Kadlec T, Weiss A, O'Geen H, Farnham PJ, Madden PAF, Mungall AJ, Tam A, Kamoh B, Cho S, Moore R, Hirst M, Marra MA, Costello JF, Wang T: **DNA hypomethylation within specific transposable element families associates with tissue-specific enhancer landscape.** *Nature Genetics* 2013, **45**:836–84110.1038/ng.2649 Available: <http://www.ncbi.nlm.nih.gov/pubmed/23708189>.

387. Jang HS, Shah NM, Du AY, Dailey ZZ, Pehrsson EC, Godoy PM, Zhang D, Li D, Xing X, Kim S, O'Donnell D, Gordon JI, Wang T: **Transposable elements drive widespread expression of oncogenes in human cancers.** *Nature Genetics* 2019, **51**:611–61710.1038/s41588-019-0373-3 Available: <http://dx.doi.org/10.1038/s41588-019-0373-3>.

388. Rajendiran S, Gibbs LD, Van Treuren T, Klinkebiel DL, Vishwanatha JK: **MIEN1 is tightly regulated by SINE Alu methylation in its promoter.** *Oncotarget* 2016, **7**:65307–6531910.18632/oncotarget.11675 Available: <http://www.oncotarget.com/fulltext/11675>.

389. Jintaridt P, Tungtrongchitr R, Preutthipan S, Mutirangura A: **Hypomethylation of Alu Elements in Post-Menopausal Women with Osteoporosis.** *PLoS ONE* 2013, **8**:e7038610.1371/journal.pone.0070386 Available: <https://dx.plos.org/10.1371/journal.pone.0070386>.

390. Lu S, Niu Z, Chen Y, Tu Q, Zhang Y, Chen W, Tong W, Zhang Z: **Repetitive Element DNA Methylation is Associated with Menopausal Age.** *Aging and Disease* 2018, **9**:43510.14336/AD.2017.0810 Available: <http://www.aginganddisease.org/EN/10.14336/AD.2017.0810>.

391. Payer LM, Steranka JP, Yang WR, Kryatova M, Medabalimi S, Ardeljan D, Liu C, Boeke JD, Avramopoulos D, Burns KH: **Structural variants caused by Alu insertions are associated with risks for many human diseases.** *Proceedings of the National Academy of*

- Sciences* 2017, **114**:E3984–E399210.1073/pnas.1704117114 Available: <http://www.pnas.org/lookup/doi/10.1073/pnas.1704117114>.
392. Mozhui K, Pandey AK: **Conserved effect of aging on DNA methylation and association with EZH2 polycomb protein in mice and humans.** *Mechanisms of Ageing and Development* 2017, **162**:27–3710.1016/j.mad.2017.02.006 Available: <http://dx.doi.org/10.1016/j.mad.2017.02.006>.
393. Smit A, Hubley R, Green: **Repeat Masker.** 2009, Available: <http://www.repeatmasker.org>.
394. Fairley S, Lowy-Gallego E, Perry E, Flicek P: **The International Genome Sample Resource (IGSR) collection of open human genomic variation resources.** *Nucleic Acids Research* 2020, **48**:D941–D94710.1093/nar/gkz836 Available: <https://academic.oup.com/nar/article/48/D1/D941/5580898>.
395. Privé F, Aschard H, Blum MGB: **Efficient Implementation of Penalized Regression for Genetic Risk Prediction.** *Genetics* 2019, **212**:65–7410.1534/genetics.119.302019 Available: [http://www.genetics.org/lookup/doi/10.1534/genetics.119.302019%20https://figshare.com/articles/code/7178750%20https://gsajournals.figshare.com/articles/Supplemental%7B/\\_%7DMaterial%7B/\\_%7Dfor%7B/\\_%7DPriv%7B/\\_%7DAschard%7B/\\_%7Dand%7B/\\_%7DBlum%7B/\\_%7D2019/7851470](http://www.genetics.org/lookup/doi/10.1534/genetics.119.302019%20https://figshare.com/articles/code/7178750%20https://gsajournals.figshare.com/articles/Supplemental%7B/_%7DMaterial%7B/_%7Dfor%7B/_%7DPriv%7B/_%7DAschard%7B/_%7Dand%7B/_%7DBlum%7B/_%7D2019/7851470).
396. Tibshirani R, Bien J, Friedman J, Hastie T, Simon N, Taylor J, Tibshirani RJ: **Strong rules for discarding predictors in lasso-type problems.** *Journal of the Royal Statistical Society. Series B: Statistical Methodology* 2012, **74**:245–26610.1111/j.1467-9868.2011.01004.x Available: <http://arxiv.org/abs/1011.2234>.
397. Meyer DH: **A transcriptome based aging clock near the theoretical limit of accuracy.** 2020, [10.1101/2020.05.29.123430](https://doi.org/10.1101/2020.05.29.123430) Available: <https://doi.org/10.1101/2020.05.29.123430>.
398. Hysi PG, Young TL, Mackey DA, Andrew T, Fernández-Medarde A, Solouki AM, Hewitt AW, Macgregor S, Vingerling JR, Li Y-j, Ikram MK, Fai LY, Sham PC, Manyes L, Porteros A, Lopes MC, Carbonaro F, Fahy SJ, Martin NG, Duijn CM van, Spector TD, Rahi JS, Santos E, Klaver CCW, Hammond CJ: **A genome-wide association study for myopia and refractive error identifies a susceptibility locus at 15q25.** *Nature Genetics* 2010, **42**:902–90510.1038/ng.664 Available: <http://www.nature.com/articles/ng.664>.
399. Loh P-r, Palamara PF, Price AL: **Fast and accurate long-range phasing in a UK Biobank cohort.** *Nature Genetics* 2016, **48**:811–81610.1038/ng.3571 Available: <http://www.nature.com/articles/ng.3571>.
400. Das S, Forer L, Schönherr S, Sidore C, Locke AE, Kwong A, Vrieze SI, Chew EY, Levy S, McGue M, Schlessinger D, Stambolian D, Loh P-R, Iacono WG, Swaroop A, Scott LJ, Cucca F,

- Kronenberg F, Boehnke M, Abecasis GR, Fuchsberger C: **Next-generation genotype imputation service and methods.** *Nature Genetics* 2016, **48**:1284–128710.1038/ng.3656 Available: <http://www.nature.com/articles/ng.3656>.
401. McCarthy S: **A reference panel of 64,976 haplotypes for genotype imputation.** *Nature Genetics* 2016, **48**:1279–128310.1038/ng.3643 Available: <http://www.nature.com/articles/ng.3643>.
402. Zhou X, Stephens M: **Genome-wide efficient mixed-model analysis for association studies.** *Nature Genetics* 2012, **44**:821–82410.1038/ng.2310 Available: <http://www.nature.com/articles/ng.2310>.
403. Snir S, Farrell C, Pellegrini M: **Human epigenetic ageing is logarithmic with time across the entire lifespan.** *Epigenetics* 2019, **0**:1–1510.1080/15592294.2019.1623634 Available: <https://www.tandfonline.com/doi/full/10.1080/15592294.2019.1623634>.
404. Hellen EH, Brookfield JF: **Alu elements in primates are preferentially lost from areas of high GC content.** *PeerJ* 2013, **1**:e7810.7717/peerj.78 Available: <https://peerj.com/articles/78>.
405. Pruim RJ, Welch RP, Sanna S, Teslovich TM, Chines PS, Gliedt TP, Boehnke M, Abecasis GR, Willer CJ: **LocusZoom: regional visualization of genome-wide association scan results.** *Bioinformatics* 2010, **26**:2336–233710.1093/bioinformatics/btq419 Available: <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btq419>.
406. Staley JR, Blackshaw J, Kamat MA, Ellis S, Surendran P, Sun BB, Paul DS, Freitag D, Burgess S, Danesh J, Young R, Butterworth AS: **PhenoScanner: a database of human genotype–phenotype associations.** *Bioinformatics* 2016, **32**:3207–320910.1093/bioinformatics/btw373 Available: <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btw373>.
407. Staley JR, Blackshaw J, Kamat MA, Ellis S, Surendran P, Sun BB, Paul DS, Freitag D, Burgess S, Danesh J, Young R, Butterworth AS: **PhenoScanner: a database of human genotype–phenotype associations.** *Bioinformatics (Oxford, England)* 2016, **32**:3207–320910.1093/bioinformatics/btw373 Available: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC5048068/>.
408. Gagliano Taliun SA, VandeHaar P, Boughton AP, Welch RP, Taliun D, Schmidt EM, Zhou W, Nielsen JB, Willer CJ, Lee S, Fritzsche LG, Boehnke M, Abecasis GR: **Exploring and visualizing large-scale genetic associations by using PheWeb.** *Nature Genetics* 2020, **52**:550–55210.1038/s41588-020-0622-5.
409. Cancer Research UK: **Gallbladder cancer incidence statistics.** 2020, Available: <https://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-cancer-type/gallbladder-cancer/incidence%7B/#%7Dheading-One>. Accessed 17 September 2020.

410. Löhler J, Cebulla M, Shehata-Dieler W, Volkenstein S, Völter C, Walther LE: **Hearing impairment in old age.** *Deutsches Aerzteblatt Online* 2019, **116**:301–31010.3238/arztebl.2019.0301Available: <https://www.aerzteblatt.de/10.3238/arztebl.2019.0301>.
411. Hileeto D, Fadare O, Martel M, Zheng W: **Age dependent association of endometrial polyps with increased risk of cancer involvement.** *World journal of surgical oncology* 2005, **3**:810.1186/1477-7819-3-8Available: <http://www.ncbi.nlm.nih.gov/pubmed/15703068>.
412. Abhishekha HA, Nisarga P, Kisan R, Meghana A, Chandran S, Trichur Raju, Sathyaprabha TN: **Influence of age and gender on autonomic regulation of heart.** *Journal of Clinical Monitoring and Computing* 2013, **27**:259–26410.1007/s10877-012-9424-3Available: <http://link.springer.com/10.1007/s10877-012-9424-3>.
413. Nikitin NP, Witte KKA, Ingle L, Clark AL, Farnsworth TA, Cleland JGF: **Longitudinal myocardial dysfunction in healthy older subjects as a manifestation of cardiac ageing.** *Age and Ageing* 2005, **34**:343–34910.1093/ageing/afi043Available: <http://academic.oup.com/ageing/article/34/4/343/10217/Longitudinal-myocardial-dysfunction-in-healthy>.
414. Lu AT, Hannon E, Levine ME, Crimmins EM, Lunnon K, Mill J, Geschwind DH, Horvath S: **Genetic architecture of epigenetic and neuronal ageing rates in human brain regions.** *Nature Communications* 2017, **8**:1535310.1038/ncomms15353Available: <http://dx.doi.org/10.1038/ncomms15353>.
415. Gibson J, Russ TC, Clarke T-K, Howard DM, Hillary RF, Evans KL, Walker RM, Birmingham ML, Morris SW, Campbell A, Hayward C, Murray AD, Porteous DJ, Horvath S, Lu AT, McIntosh AM, Whalley HC, Marioni RE: **A meta-analysis of genome-wide association studies of epigenetic age acceleration.** *PLOS Genetics* 2019, **15**:e100810410.1371/journal.pgen.1008104Available: <http://www.ncbi.nlm.nih.gov/pubmed/31738745>.
416. Mccartney DL, Min JL, Richmond RC, Lu AT, Maria K, Davies G, Broer L, Guo X, Jeong A, Kasela S, Katrinli S, Kuo P-l, Matias- PR, Mishra PP, Nygaard M, Palviainen T, Soerensen M, Sun D, Tsai P-c, Matthijs D, Xu Z, Yao J, Zhao W, Correa A, Boerwinkle E, Durda P, Elliott HR, Gieger C, Genetics T, Consortium M, et al.: **Genome-wide association studies identify 137 loci for DNA methylation biomarkers of ageing.** *bioRxiv* 2020,:1–5010.1101/2020.06.29.133702Available: <https://doi.org/10.1101/2020.06.29.133702>.