

# Genetic susceptibility to earlier ovarian ageing increases *de novo* mutation rate in offspring

Stasa Stankovic\*<sup>1</sup>, Saleh Shekari\*<sup>2,3</sup>, Qin Qin Huang\*<sup>4</sup>, Eugene J. Gardner\*<sup>1</sup>, Nick D. L. Owens\*<sup>2</sup>, Ajuna Azad<sup>5</sup>, Gareth Hawkes<sup>2</sup>, Katherine A. Kentistou<sup>1</sup>, Robin N. Beaumont<sup>2</sup>, Felix R. Day<sup>1</sup>, Yajie Zhao<sup>1</sup>, The Genomics England Research Consortium<sup>8,9</sup>, Kitale Kennedy<sup>2</sup>, Andrew R. Wood<sup>2</sup>, Michael N. Weedon<sup>2</sup>, Ken K. Ong<sup>1,6</sup>, Caroline F. Wright<sup>2</sup>, Eva R. Hoffmann<sup>5</sup>, Matthew E. Hurles<sup>4</sup>, Katherine S. Ruth<sup>2</sup>, Hilary C. Martin<sup>4</sup>, John R. B. Perry\*<sup>1,7</sup> and Anna Murray\*<sup>2</sup>

\* Denotes equal contribution

Correspondence to John R.B Perry ([john.perry@mrc-epid.cam.ac.uk](mailto:john.perry@mrc-epid.cam.ac.uk)) and Anna Murray ([A.Murray@exeter.ac.uk](mailto:A.Murray@exeter.ac.uk))

## Affiliations

<sup>1</sup>MRC Epidemiology Unit, Wellcome–MRC Institute of Metabolic Science, University of Cambridge, Cambridge, UK

<sup>2</sup>University of Exeter Medical School, University of Exeter, Exeter, UK.

<sup>3</sup>School of Public Health, Faculty of Medicine, University of Queensland, Brisbane, AU.

<sup>4</sup>Wellcome Sanger Institute, Wellcome Genome Campus, Cambridge, Hinxton, UK

<sup>5</sup>DNRF Center for Chromosome Stability, Department of Cellular and Molecular Medicine, Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark.

<sup>6</sup>Department of Paediatrics, University of Cambridge, Cambridge, UK

<sup>7</sup>Metabolic Research Laboratory, Wellcome–MRC Institute of Metabolic Science, University of Cambridge, Cambridge, UK

<sup>8</sup>Genomics England, London, UK

<sup>9</sup>William Harvey Research Institute, Queen Mary University of London, London, EC1M 6BQ, UK

# Abstract

Human genetic studies have provided substantial insight into the biological mechanisms governing ovarian ageing, yet previous approaches have been largely restricted to assessing common genetic variation. Here we report analyses of rare (MAF<0.1%) protein-coding variants in the exomes of 106,973 women from the UK Biobank study, implicating novel genes with effect sizes up to ~5 times larger than previously discovered in analyses of common variants. These include protein truncating variants in *ZNF518A*, which shorten reproductive lifespan by promoting both earlier age at natural menopause (ANM, 5.61 years [4.04-7.18],  $P=2*10^{-12}$ ) and later puberty timing in girls (age at menarche, 0.56 years [0.15-0.97],  $P=9.2*10^{-3}$ ). By integrating ChIP-Seq data, we demonstrate that common variants associated with ANM and menarche are enriched in the binding sites of *ZNF518A*. We also identify further links between ovarian ageing and cancer susceptibility, highlighting damaging germline variants in *SAMHD1* that delay ANM and increase all-cause cancer risk in both males (OR=2.1 [1.7-2.6],  $P=4.7*10^{-13}$ ) and females (OR=1.61 [1.31-1.96],  $P=4*10^{-6}$ ). Finally, we demonstrate that genetic susceptibility to earlier ovarian ageing in women increases *de novo* mutation rate in their offspring. This provides direct evidence that female mutation rate is heritable and highlights an example of a mechanism for the maternal genome influencing child health.

# Introduction

Reproductive longevity in women varies substantially in the general population, and has a profound impact on fertility and health outcomes in later life<sup>1-3</sup>. Women are born with a non-renewable ovarian reserve, which is established during foetal development. This reserve is continuously depleted throughout reproductive life, ultimately leading to menopause<sup>4-6</sup>. Variation in menopause timing is largely dependent on the differences in the size of the initial oocyte pool and the rate of follicle loss<sup>3</sup>. Natural fertility is believed to be closely associated with menopause timing, and it declines on average 10 years before the onset of menopause<sup>4,7</sup>. The effect of early menopause on infertility is becoming increasingly relevant due to the secular trend of delaying parenthood to later maternal age at childbirth, especially in Western populations. In addition, normal variation in reproductive lifespan is causally associated with the risk of a wide range of disease outcomes, such as type 2 diabetes mellitus, cancer and impaired bone health, further highlighting the need for better understanding of the regulators and physiological mechanisms involved in reproductive ageing<sup>1,8</sup>.

The variation in timing of menopause reflects a complex mix of genetic and environmental factors that population-based studies have begun to unravel. Previous genome-wide association studies (GWAS) have successfully identified ~300 distinct common genomic loci associated with the timing of menopause<sup>1</sup>. These reported variants cumulatively explain 10% - 12% of the variance in ANM and 31-38% of the overall estimated SNP heritability<sup>1,9,10</sup>. The majority of these loci implicate genes that regulate DNA damage response (DDR), highlighting the particular sensitivity of oocytes to DNA damage due to the prolonged state of cell cycle arrest across the life-course<sup>1,7,11-20</sup>.

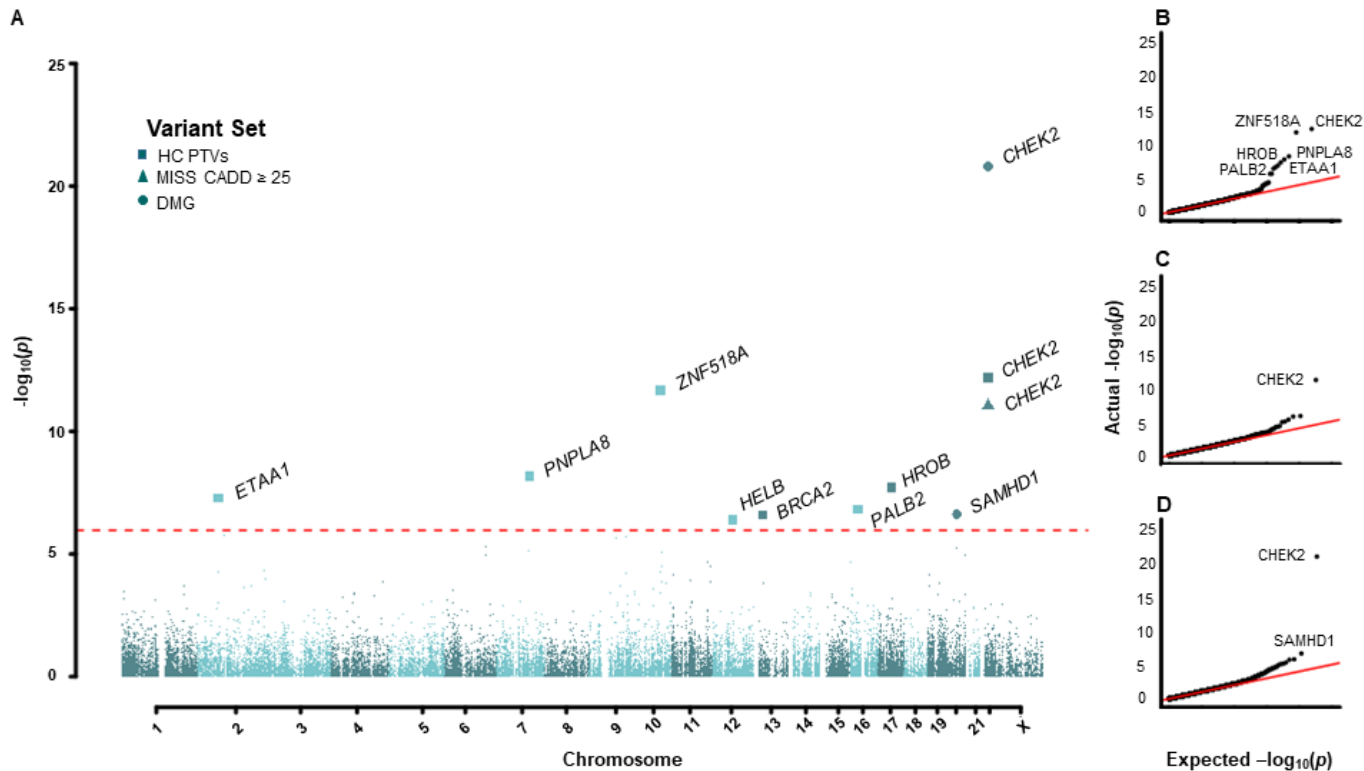
Genetic studies for ANM to date have largely focussed on assessing common genetic variation, with little insight into the role of rarer, protein-coding variants. Initial exome-sequencing (WES) analyses in UK Biobank identified gene-based associations with ANM for *CHEK2*, *DCLRE1A*, *HELB*, *TOP3A*, *BRCA2* and *CLPB*<sup>1,9</sup>. In this study, we aimed to explore the role of rare damaging variants in ovarian ageing in greater detail through a combination of enhanced phenotype curation, better powered statistical tests and assessment of different types of variant classes at lower allele frequency thresholds (**Supplementary Note**). Using these approaches we identify five genes harbouring variants of large effect that have not previously been implicated, highlighting *ZNF518A* as a major transcriptional regulator of ovarian ageing. Furthermore, we extend these observations to show that women at increased genetic risk of earlier menopause have increased rates of *de novo* mutations in their offspring.

## Results

### Exome-wide gene burden associations with ANM

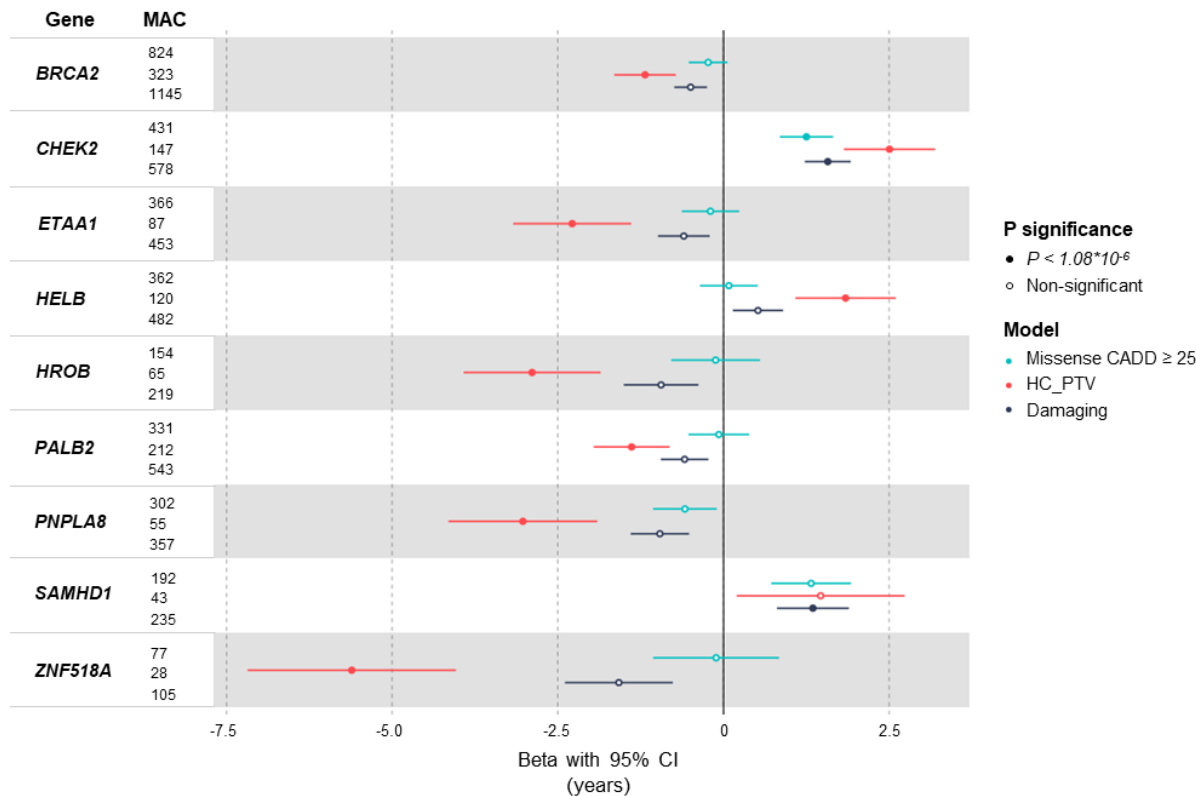
To assess the impact of rare damaging variants on age at natural menopause (ANM), we used whole-exome sequencing (WES) data available in 106,973 post-menopausal UK Biobank female participants of European genetic-ancestry<sup>21</sup>. Individual gene burden association tests were conducted by collapsing genetic variants according to their predicted functional categories. We defined three categories of rare exome variants with minor allele frequency (MAF) < 0.1%: high-confidence Protein Truncating Variants (HC-PTVs), missense variants with CADD score  $\geq 25$ , and 'damaging' variants (defined as combination of HC-PTVs and missense variants with CADD  $\geq 25$ ). We analysed 17,475 protein-coding genes with the minimum of 10 rare allele carriers in at least one of the masks tested. The primary burden association analysis was conducted using BOLT-LMM<sup>22</sup> (**Supplementary Table 1**). The low exome-wide inflation scores (e.g. PTV  $\lambda=1.047$ ) and the absence of significant association with synonymous variant burden for any gene indicate our statistical tests are well calibrated (**Supplementary Figure 1**).

We identified rare variation in nine genes associated with ANM at exome-wide significance ( $P < 1.08 \times 10^{-6}$ , **Figures 1 and 2, Supplementary Figures 2 and 3**). These were confirmed by an independent group of analysts using different QC and analysis pipelines (**Supplementary Tables 1, 2**). Three of these genes have been previously reported in UKBB WES analysis<sup>9</sup> - we confirm the associations of *CHEK2* (beta=1.57 years, 95% CI: 1.23-1.92,  $P=1.60 \times 10^{-21}$ , N=578 damaging allele carriers) and *HELB* (beta=1.84, 95% CI: 1.08-2.60,  $P=4.20 \times 10^{-7}$ , N=120 HC-PTV carriers) with later ANM and a previously borderline association of *HROB* with earlier ANM (beta= -2.89 years, 95% CI: 1.86-3.92,  $P=1.90 \times 10^{-8}$ , N=65 HC-PTV carriers). In addition, our previous ANM GWAS analyses<sup>1</sup> identified an individual low-frequency PTV variant in *BRCA2*, which we now extend to demonstrate that, in aggregate, *BRCA2* HC-PTV carriers exhibit 1.18 years earlier ANM (beta= -1.18, 95% CI: 0.72-1.65,  $P=2.60 \times 10^{-7}$ , N=323). Rare variants in the remaining five genes – *ETAA1*, *ZNF518A*, *PNPLA8*, *PALB2* and *SAMHD1* have not been previously implicated in ovarian ageing. Effect sizes of these associations range from 5.61 years earlier ANM for HC-PTV carriers in *ZNF518A*, to 1.35 years later ANM for women carrying damaging alleles in *SAMHD1*. This contrasts with a maximum effect size of 1.06 years (median 0.12 years) for common variants (MAF>1%) identified by previous ANM GWAS<sup>1</sup>.



**Figure 1: Exome-wide associations with age at natural menopause. (A)** Manhattan plot showing gene burden test results for age at natural menopause. Genes passing exome-wide significance ( $P < 1.08 \times 10^{-6}$ ) are indicated, with point shape signifying the variant class tested. **(B-E)** QQ plots for **(B)** high confidence PTVs **(C)** CADD  $\geq 25$  missense variants **(D)** damaging variants.

We next sought to understand why previous analyses of UKBB WES data missed the associations we report here, and conversely why we did not identify associations with other previously reported genes. Of the seven genes identified by Ward *et al.*<sup>9</sup>, three were also identified by our study (*CHEK2*, *HELB* and *HROB*), three were recovered when we increased our burden test MAF threshold from 0.1% to 1% (*DCLRE1A*, *RAD54L*, *TOP3A*), and an additional gene fell just below our  $P$  value threshold when considering variants with  $< 1\%$  MAF (*CLPB*;  $P = 1.2 \times 10^{-5}$ ). In contrast, our discovery of novel associations that were not reported by Ward *et al.* (*BRCA2*, *ETAA1*, *PALB2*, *PNPLA8*, *SAMHD1* and *ZNF518A*) were likely explained by differences in phenotype preparation, sample size, variant annotation and the statistical model used (see **Supplementary Note and Supplementary Table 3**).



**Figure 2: Forest plot for gene burden associations with age at natural menopause.** Exome-wide significant ( $P < 1.08 \times 10^{-6}$ ) genes are displayed. Points and error bars indicate beta and 95% CI for the variant category indicated. Betas, CIs, Minor Allele Counts (MAC) and P values are derived from BOLT-LMM.

## Exploring common variant associations at identified ANM genes

To explore the overlap between common and rare variant association signals for ANM, we integrated our exome-wide results with data generated from the largest reported common variant GWAS of ANM<sup>1</sup>.

Five of our nine identified WES genes (*CHEK2*, *BRCA2*, *ETAA1*, *HELB* and *ZNF518A*) mapped within 500kb of a common GWAS signal (**Supplementary Table 4**). Notably, we previously reported a common, predicted benign, missense variant (rs35777125-G439R, MAF=11%) in *ETAA1* associated with 0.26 years earlier ANM. In contrast, our WES analysis identified that carriers of rare HC-PTVs in *ETAA1* show a nearly 10-fold earlier ANM (beta= -2.28 years, 95% CI: 1.39-3.17,  $P=5.30 \times 10^{-8}$ , N=87). Furthermore, three independent non-coding common GWAS signals ~150kb apart (MAF: 2.8-47.5%, beta: -0.28-0.28 years per minor allele) were reported proximal to *ZNF518A*, whereas gene burden testing finds that rare HC-PTV carriers show nearly 20-fold earlier ANM than common variant carriers (beta= -5.61 years, 95% CI: 4.04-7.18,  $P=2.10 \times 10^{-12}$ , N=28).

In addition there were two genes within 500kb of GWAS loci (*BRCA1* and *SLCO4A1*) that were associated with ANM by gene burden testing at  $P < 1.7 \times 10^{-5}$ . Effect sizes for common variant associations ranged from 0.07-0.24 years per allele at these loci, whereas gene burden tests for rarer variants at these same loci revealed much larger effect sizes: for *BRCA1*, 2.1 years earlier for PTVs ( $P = 2.4 \times 10^{-6}$ ) and for *SLCO4A1*, 1.13 years earlier ANM for damaging variants ( $P = 1.1 \times 10^{-5}$ ), with non-overlapping 95% confidence intervals between common and rare variant associations for *BRCA1*.

## Common ANM associated variants are enriched in *ZNF518A* binding sites

Heterozygous loss of function of *ZNF518A* had the largest effect on ANM of the genes we identified. *ZNF518A* is a poorly characterised C2H2 zinc finger transcription factor, which has been shown to associate with PRC2 and G9A-GLP repressive complexes along with its paralog *ZNF518B*, suggesting a potential role in transcriptional repression<sup>23</sup>. *ZNF518A* localises robustly to 18,706 sites in the genome, based on ChIP-seq data available from ENCODE<sup>24,25</sup> and binds primarily to gene promoters, with 33.5% (6,263) of *ZNF518A* binding sites within 2kb of a transcription start site (TSS) (**Supplementary Figure 4a-c**). Common variants associated with ANM<sup>1</sup> were enriched in the transcriptional targets of *ZNF518A* ( $P = 1.32 \times 10^{-4}$ ) using fGWAS<sup>26</sup>. We further tested functional enrichment using signed linkage disequilibrium profile (SLDP) regression<sup>27</sup>. This confirmed the enrichment of *ZNF518A* binding sites near to loci associated with ANM and showed that its transcriptional repression is associated with earlier ANM ( $P = 0.02$ ), consistent with evidence from rare variant burden tests. Separating *ZNF518A* sites by those proximal (< 2Kb) and distal (>5kb) from a TSS, demonstrated this association was due to *ZNF518A* binding at regulatory regions distal to the TSS (proximal TSS  $P = 0.3$ , distal *ZNF518A*  $P = 0.002$ ). Notably, these regulatory *ZNF518A* bound loci produce the largest association amongst an SLDP catalogue of 382 transcription factors and regulators (**Supplementary Table 5, Supplementary Figure 4d**). These results suggest a different functional role for *ZNF518A* at TSS and more distal regulatory regions. In order to explore this further we assessed the sequence determinants of *ZNF518A* binding. *De novo* motif discovery identified an AT-rich motif enriched at distal regulatory *ZNF518A* binding sites, but not at TSS bound by *ZNF518A*. This AT-rich motif was centrally enriched within *ZNF518A* ChIP-seq peaks, and matched an unvalidated motif present in the JASPAR transcription factor motif database<sup>28</sup> (**Supplementary Figure 4e**). We found the number of perfect instances of this AT-rich motif to be strongly associated with *ZNF518A* occupancy as assessed by *ZNF518A* ChIP-seq signal at distal regions but not at TSS (**Supplementary Figure 4f,g**). At distal regions, the maximal association between peaks greater than the median height was found at least seven motif instances (Hypergeometric right tail  $P < 10^{-389}$ , Odds Ratio 7.41). These data suggest that *ZNF518A* is recruited by DNA sequence at distal sites, but at TSS may be recruited to gene promoters by interaction with another DNA binding factor.

We next employed public *in vitro* differentiated human primordial germ like-cell data<sup>29,30</sup> to assess the chromatin state at *ZNF518A* bound loci, directly comparing distal regions with TSS. *ZNF518A* bound TSS showed chromatin accessibility<sup>30</sup> and were marked with H3K27ac<sup>29</sup>. In

contrast, distal regions lacked H3K27ac and showed minimal chromatin accessibility (**Supplementary Figure 4h**). Extending this comparison to the Epimap chromatin states<sup>31</sup>, we find that overall *ZNF518A* bound loci are enriched in active TSS and that distal *ZNF518A* regions are variously enriched in active and repressed chromatin (**Supplementary Figure 4i,j**). Consistent with previous data which has found *ZNF518A* in repressive complexes, these data suggest that *ZNF518A* is recruited by DNA sequence to distal regulatory regions where it acts to repress local chromatin.

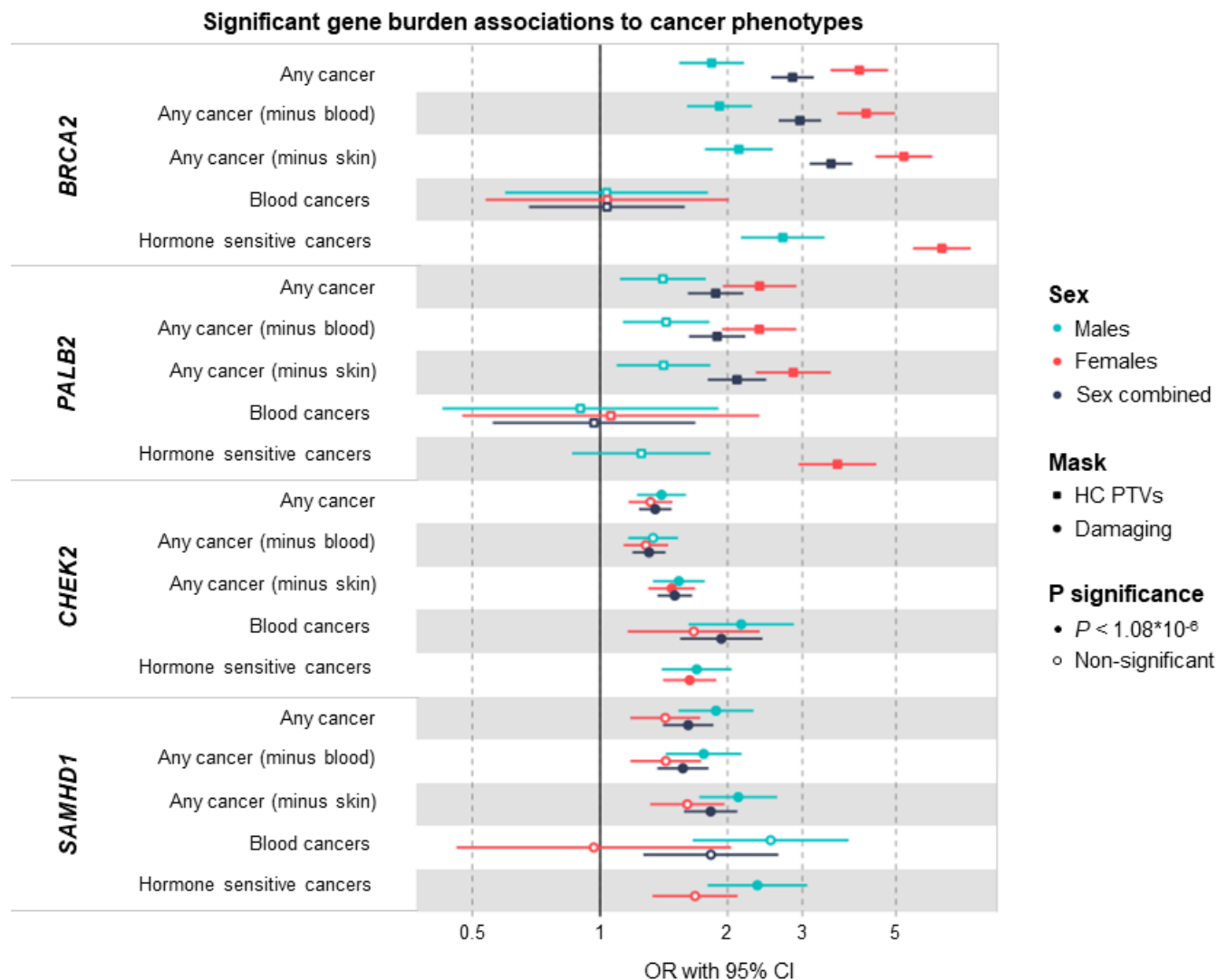
While *ZNF518A* is known to have diverse tissue expression including the ovary, we found that it was highly expressed in fetal germ cells at both the mitotic and meiotic stages (**Supplementary Tables 6 and 7; Supplementary Figures 5 and 6**). The eight other WES genes identified in this study were expressed at varying levels in fetal gonadal cells, oocytes and granulosa cells across different developmental stages (**Supplementary Figures 5 and 6**).

## Identified genes influence other aspects of health and disease

Our genetic studies have previously shown that the genetic mechanisms regulating the end of reproductive life are largely distinct from those determining its beginning<sup>32,33</sup>. However, it is noteworthy that the largest reported GWAS for age at menarche identified a common variant signal at the *ZNF518A* locus for later puberty timing in girls (rs1172955, beta= 0.04 years, 95% CI: 0.03-0.05,  $P=6.6 \times 10^{-12}$ ), which appears nominally associated with earlier ANM (beta=-0.04, 95% CI: 0.01-0.06,  $P=6.6 \times 10^{-3}$ )<sup>32</sup>. To extend this observation, we found that our identified *ZNF518A* PTVs were also associated with later age at menarche (0.56 years, 95% CI: 0.14-0.98,  $P=9.2 \times 10^{-3}$ ). Furthermore, using fGWAS and SLDP, we discovered that, similar to ANM, common variants that influence puberty in girls were enriched in transcriptional targets of *ZNF518A* (**Supplementary Table 5**). These data suggest that loss of *ZNF518A* shortens reproductive lifespan, by delaying puberty and reducing age at menopause.

We next explored what impact ANM-associated genes had on cancer outcomes and found a novel association of *SAMHD1* damaging variants and HC-PTVs with 'All cancer' in both males (OR=2.12, 95% CI: 1.72-2.62,  $P=4.7 \times 10^{-13}$ ) and females (OR=1.61, 95% CI: 1.31-1.96,  $P=4 \times 10^{-6}$ ; **Figure 3, Supplementary Table 8-10**). In addition we replicated previously reported associations with protein truncating variants in *BRCA2*, *CHEK2* and *PALB2* and cancer outcomes in males and females (**Supplementary Tables 8-10**).



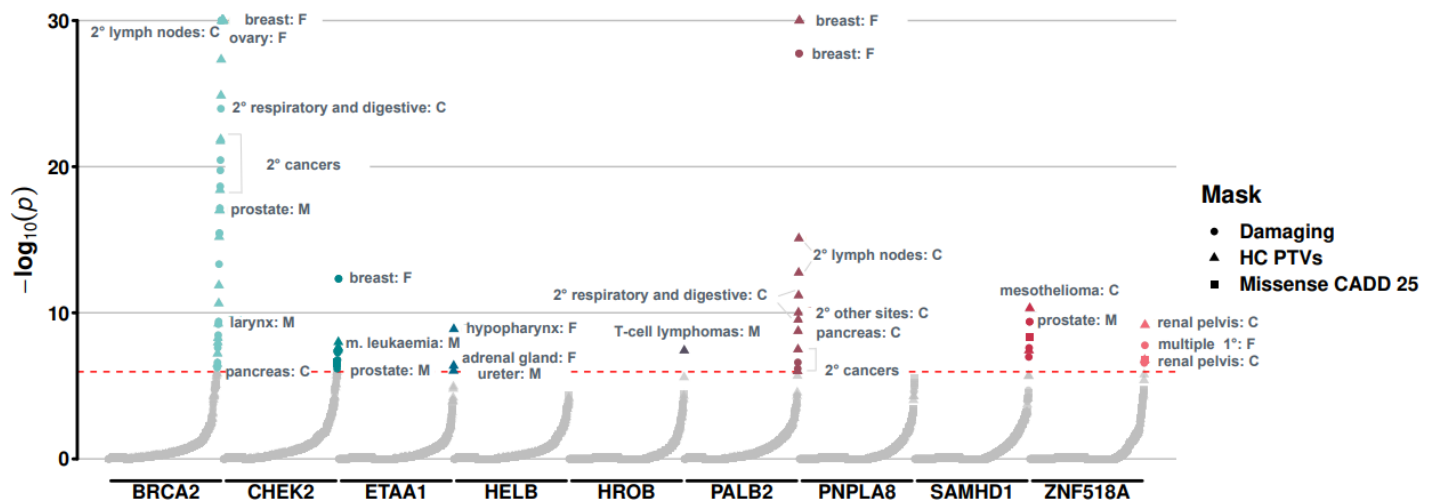


**Figure 3: Forest plot for ANM WES genes with significant gene burden associations for cancer phenotypes.** Exome-wide significant ( $P < 1.08 \times 10^{-6}$ ) genes are displayed, showing sex-stratified and combined results. Hormone sensitive cancers were only tested in males and females separately (Methods). The presented masks were selected based on the most significant association per gene and cancer type. Points and bars indicate OR and 95% CI for specific genes and their variant categories in cancer. Filled symbols indicate a result passing a Bonferroni-corrected significance threshold of  $P < 1.08 \times 10^{-6}$ .

*SAMHD1* associations with cancer appear to be driven by increased risk for multiple site-specific cancers, notably prostate cancer in males, mesothelioma in both males and females, and suggestive evidence for higher breast cancer susceptibility in females (Figure 4, Supplementary Table 11). Although the numbers of mutation carriers diagnosed with each site-specific cancer was small, the majority of these findings persisted using logistic regression

with penalised likelihood estimation, which is more robust to extreme case/control imbalance<sup>34</sup> (**Supplementary Table 11**).

Cancer risk-increasing alleles in *SAMHD1* were associated with later ANM, which is similar to the pattern demonstrated previously for *CHEK2*. This finding is consistent with a mechanism of disrupted DNA damage sensing and apoptosis, resulting in slowed depletion of the ovarian reserve<sup>1</sup>. In addition, we provide robust evidence for a previously described rare variant association for *SAMHD1* with telomere length<sup>35</sup>, highlighting that rare damaging variants cause longer telomere length ( $P=1.4 \times 10^{-59}$ ) (**Supplementary Table 10, Supplementary Figure 7**).



**Figure 4: Genetic susceptibility to premature ovarian ageing and increased risk for diverse cancer types.** Plot showing the association between loss of ANM genes identified in this study and risk of 90 site specific cancers among UK Biobank participants. Summary statistics for cancer associations were obtained using a logistic regression with penalised likelihood estimation that controls for case/control imbalance (Methods)<sup>34</sup>. Associations highlighted in text passed exome-wide significance ( $P < 1.08 \times 10^{-6}$ ). The y-axis is capped at  $-\log_{10}(P) = 30$  for visualisation purposes; un-capped summary statistics can be found in **Supplementary Table 11**. F: females, M: males, C: sex-combined. 1°: primary cancer, 2°: secondary cancer.

## Genetic susceptibility to ANM in mothers influences *de novo* mutation rate in offspring

Our previous common variant analyses demonstrated that many ANM associated variants implicate DNA damage repair (DDR) genes, an observation mirrored here in our rare variant associations. Therefore, we sought to test the hypothesis that inter-individual variation in these DDR processes would influence the mutation rate in germ cells and hence in the offspring. More specifically, we hypothesised that genetic susceptibility to earlier ovarian ageing would be associated with a higher *de novo* mutation (DNM) rate in the offspring. To test this, we analysed

8,089 whole-genome sequenced parent-offspring trios recruited in the rare disease programme of the 100,000 Genome Project (100kGP, **Supplementary Figure 8**). We calculated a polygenic score (PGS) for ANM in the parents based on our previously identified 290 common variants<sup>1</sup> and tested this against the phased DNMR rate in the offspring, adjusted for age. We found that maternal genetic susceptibility to earlier ANM was associated with an increased rate of maternally-derived DNMRs in the offspring (rate ratio = 1.02 per SD of PGS,  $P=6.8 \times 10^{-4}$ ,  $N=8,089$  duos with European ancestry; **Supplementary Table 12**). We confirmed this finding in sensitivity analyses using the same data, in a two-sample Mendelian Randomization (MR) framework that can better model the dose-response relationship of these variants (**Supplementary Table 13**). These results were highly concordant, with all models showing a significant result and no heterogeneity ( $P_{\min}=6.3 \times 10^{-5}$ ). In contrast, the paternal PGS was not associated with paternally-derived DNMRs ( $P=0.51$ ,  $N=8,029$ ) nor was the maternal PGS associated with paternally-derived DNMRs ( $P=0.55$ ).

## Discussion

Our study extends the number of genes implicated in ovarian ageing through the identification of rare, protein-coding variants. Effect sizes ranged from 5.61 years earlier ANM for HC-PTV carriers in *ZNF518A*, to 1.35 years later ANM for women carrying damaging variants in *SAMHD1* compared to a maximum effect size of 1.06 years (median 0.12 years) reported for common variants ( $MAF > 1\%$ )<sup>1</sup>. Several of these effect estimates were comparable to those conferred by *FMR1* premutations, which are currently used as part of the only routinely applied clinical genetics test for premature ovarian insufficiency (POI)<sup>36</sup>. Deleterious variants in three genes (*CHEK2*, *HELB* and *SAMHD1*) were associated with an increase in ANM and therefore represent potential therapeutic targets for enhancing ovarian stimulation in women undergoing *in vitro* fertilisation (IVF) treatment through short-term apoptotic inhibition. Seven of the nine ANM genes identified have known roles in DNA damage repair, and three of these are linked to ANM for the first time (*PALB2*, *ETAA1* and *HROB*): *PALB2* is involved in *BRCA2* localization and stability and compound heterozygous mutations result in Fanconi anaemia and predispose to childhood malignancies<sup>37</sup>. *ETAA1* accumulates at DNA damage sites in response to replication stress<sup>38–41</sup> and *HROB* is involved in homologous recombination by recruiting the *MCM8-MCM9* helicase to sites of DNA damage to promote DNA synthesis<sup>42,43</sup>. Homozygous loss-of-function of *HROB* is associated with POI<sup>44</sup> and infertility in both sexes in mouse models<sup>42</sup>.

Novel biological mechanisms of ovarian ageing were revealed by finding associations with two non-DDR genes (*PNPLA8* and *ZNF518A*): *PNPLA8* is a calcium-independent phospholipase<sup>45–47</sup> and a recessive cause of neurodegenerative mitochondrial disease and mitochondrial myopathy<sup>48–52</sup>; an association with reproductive phenotypes has not been described previously. *ZNF518A* belongs to the zinc finger protein family and is likely a transcriptional regulator for a large number of genes<sup>23</sup>. We found that female carriers of rare protein truncating variants in

*ZNF518A* have shorter reproductive lifespan due to delayed puberty timing and earlier menopause. Enrichment of GWAS signals at *ZNF518A* binding sites suggests that *ZNF518A* regulates the genes involved in reproductive longevity by repression of elements distal to transcription start sites.

While mutation in *SAMHD1* is a common somatic event in a variety of cancers<sup>53–63</sup>, it has not been described as a germline risk factor previously. Recessive inheritance of *SAMHD1* missense and PTV variants have been associated with Aicardi–Goutieres syndrome, a congenital autoimmune disease<sup>64</sup>. Our identified damaging variants in *SAMHD1* increased risk of ‘All cancer’ in males and females, as well as in sex-specific cancers, highlighting *SAMHD1* as a novel risk factor for prostate cancer in males and hormone-sensitive cancers in females. *SAMHD1* has a role in preventing the accumulation of excess deoxynucleotide triphosphates (dNTPs), particularly in non-dividing cells<sup>65</sup>. A regulated dNTP pool is important for the fidelity of DNA repair, thus highlighting additional roles of this gene in facilitation of DNA end resection during DNA replication and repair<sup>65–70</sup>. *SAMHD1* deficiency leads to resistance to apoptosis<sup>71,72</sup>, suggesting that delayed ANM might originate from slowed depletion of ovarian reserve due to disrupted apoptosis, analogous to the mechanism for *CHEK2* that has been reported previously.

Previous studies have demonstrated that parental age is strongly associated with the number of *de novo* mutations in offspring<sup>73</sup>, with the majority of these mutations arising from the high rate of spermatogonial stem cell divisions that underlie spermatogenesis throughout adult life of males<sup>74</sup>. Our current study provides the first direct evidence that maternal mutation rate is heritable, with women at higher genetic risk of earlier menopause transmitting an increased rate of *de novo* mutations to offspring. This could have direct implications for the health of future generations given the widely reported link between *de novo* mutations and increased risk of psychiatric disease and developmental disorders<sup>75–78</sup>. We speculate that if genetic susceptibility to earlier menopause influences *de novo* mutation rate, it is possible that non-genetic risk factors for earlier ANM, such as smoking and alcohol intake, would likely have the same effect<sup>79</sup>. Our observations makes conceptual sense given that menopause timing appears to be primarily driven by the genetic integrity of oocytes and their ability to sustain, detect, repair and respond to acquired DNA damage<sup>1</sup>. These observations also build on earlier work in mice and humans that *BRCA1/2* deficiency increases the rate of double strand breaks in oocytes and reduces ovarian reserve<sup>80–82</sup>.

An important limitation of our study, shared by many other similar large-scale exome sequencing studies, is that we were unable to replicate our findings in an independent cohort. Instead, we aimed to accumulate additional evidence where possible to support our observations and evaluate the biological plausibility of our findings. For example, the identified rare loss of function alleles in *ZNF518A* have the largest effect on ovarian ageing reported to date, which is supported by high expression in fetal germ cells, genome-wide significant common variants at the same locus, and the observation that *ZNF518A* binding sites genome-wide are significantly enriched for common variant ANM association(s). For all identified genes

further experimental studies will ultimately be required to fully understand the biological mechanisms governing the observed effects on ovarian ageing.

## Methodology

### UK Biobank Data Processing and Quality Control

To conduct rare variant burden analyses described in this study, we obtained Whole Exome Sequencing data (WES) for 454,787 individuals from the UK Biobank study<sup>83</sup>. Participants were excluded based on excess heterozygosity, autosomal variant missingness on genotyping arrays  $\geq 5\%$ , or inclusion in the subset of phased samples as defined in Bycroft *et al*<sup>64</sup>. Analysis was restricted to participants with European genetic ancestry, leaving a total of 421,065 individuals. Variant quality control (QC) and annotation were performed using the UK Biobank Research Analysis Platform (RAP; <https://ukbiobank.dnanexus.com/>), a cloud-based central data repository for UK Biobank WES and phenotypic data. Besides the QC described by Backman *et al*.<sup>83</sup>, we performed additional steps using custom applets designed for the RAP. Firstly, we processed provided population-level Variant Call Format (VCF) files by splitting and left-correcting multi-allelic variants into separate alleles using 'bcftools norm'<sup>85</sup>. Secondly, we performed genotype-level filtering applying 'bcftools filter' separately for Single Nucleotide Variants (SNVs) and Insertions/Deletions (InDels) using a missingness-based approach. Using this approach, we set to missing (i.e. ./.) all SNV genotypes with depth  $< 7$  and genotype quality  $< 20$  or InDel genotypes with a depth  $< 10$  and genotype quality  $< 20$ . Next, we applied a binomial test to assess an expected alternate allele contribution of 50% for heterozygous SNVs; we set to missing all SNV genotypes with a binomial test p. value  $\leq 1 \times 10^{-3}$ . Following genotype-level filtering we recalculated the proportion of individuals with a missing genotype for each variant and filtered all variants with a missingness value  $> 50\%$ . The variant annotation was performed using the ENSEMBL Variant Effect Predictor (VEP) v104<sup>86</sup> with the '--everything' flag and plugins for CADD<sup>87</sup> and LOFTEE<sup>88</sup> enabled. For each variant we prioritised the highest impact individual consequence as defined by VEP and one ENSEMBL transcript as determined by whether or not the annotated transcript was protein-coding, MANE select v0.97, or the VEP Canonical transcript. Following annotation, variants were categorised based on their predicted impact on the annotated transcript. Protein Truncating Variants (PTVs) were defined as all variants annotated as stop gained, frameshift, splice acceptor, and splice donor. Missense variant consequences are identical to those defined by VEP. Only autosomal or chrX variants within ENSEMBL protein-coding transcripts and within transcripts included on the UKBB ES assay<sup>83</sup> were retained for subsequent burden testing.

### Exome-wide association analyses in the UK Biobank

In order to perform rare variant burden tests, we used a custom implementation of BOLT-LMM v2.3.6<sup>89</sup> for the RAP. Two primary inputs are required by BOLT-LMM: i) a set of genotypes with minor allele count  $> 100$  derived from genotyping arrays to construct a null linear mixed effects model and ii) a larger set of variants collapsed on ENSEMBL transcript to perform association

tests. For the former, we queried genotyping data available on the RAP and restricted to an identical set of individuals included for rare variant association tests. For the latter, and as BOLT-LMM expects imputed genotyping data as input rather than per-gene carrier status, we created dummy genotype files where each variant represents one gene and individuals with a qualifying variant within that gene are coded as heterozygous, regardless of the number of variants that individual has in that gene.

To test a range of variant annotation categories for  $MAF < 0.1\%$ , we created dummy genotype files for high confidence PTVs as defined by LOFTEE, missense variants with  $CADD \geq 25$ , and damaging variants that included both high confidence PTVs and missense variants with  $CADD \geq 25$ . For each phenotype tested, BOLT-LMM was then run with default parameters other than the inclusion of the 'lmmInfOnly' flag. To derive association statistics for individual markers, we also provided all 26,657,229 individual markers regardless of filtering status as input to BOLT-LMM. All tested phenotypes were run as continuous traits corrected by age,  $age^2$ , sex, the first ten genetic principal components as calculated in Bycroft *et al*<sup>64</sup> and study participant ES batch as a categorical covariate (either 50k, 200k, or 450k).

For discovery analysis in the primary trait of interest, age at natural menopause, we analysed 17,475 protein-coding genes with the minimum of 10 rare allele carriers in at least one of the masks tested using BOLT-LMM (**Supplementary Table 1**). The significant gene-level associations for ANM were identified applying Bonferroni correction for the number of masks with  $MAC \geq 10$  ( $N=46,251$  masks) in 17,475 protein-coding genes ( $P: 0.05/46,251 = 1.08 \times 10^{-6}$ ) (**Supplementary Table 2**). The age at natural menopause results obtained via BOLT-LMM are available in **Supplementary Table 1**. Furthermore, in order to compare and explain potential differences between our WES results and the previously published one<sup>9</sup>, we ran the above described approach using  $MAF < 1\%$ , a cutoff applied by Ward *et al*. (**Supplementary Table 3, Supplementary Note**).

To generate accurate odds ratio and standard error estimates for binary traits, we also implemented a generalised linear model using the statsmodels package<sup>90</sup> for python in a three step process. First, a null model was run with the phenotype as a continuous trait, corrected for control covariates as described above. Second, we regressed carrier status for individual genes on the residuals of the null model to obtain a preliminary  $P$  value. Thirdly, all genes were again tested using a full model to obtain odds ratios and standard errors with the family set to 'binomial'. Generalised linear models utilised identical input to BOLT-LMM converted to a sparse matrix.

## Phenotype derivation

Age at natural menopause was derived for individuals within the UK Biobank, who were deemed to have undergone natural menopause, i.e. not affected by surgical or pharmaceutical interventions, as follows:

Firstly, European female participants ( $n=245,820$ ) who indicated during any of the attended visits having had a hysterectomy were collated (fields 3591 and 2724) and their reported hysterectomy ages were extracted (field 2824) and the median age was kept ( $n=47,218$  and

46,260 with reported ages). The same procedure was followed for participants indicating having undergone a bilateral oophorectomy (surgery field 2834 and age field 3882, n=20,495 and 20,001 with reported ages).

For individuals having indicated the use of hormone replacement therapy (HRT; field 2814), HRT start and end ages were collated (fields 3536 and 3546, accordingly) across the different attended visits (n=98,104). In cases where the reported chronological HRT age at later attended visits was greater than that at previous visits, the later instances were prioritised, i.e. as they would potentially indicate an updated use of HRT. In cases where different HRT ages were reported, but not in chronologically increasing order, the median age was kept.

Menopausal status was determined using data across instances (field 2724) and prioritising the latest reported data, to account for changes in menopause status. For participants indicating having undergone menopause, their reported ages at menopause were collated (field 3581) using the same procedure as for HRT ages (n=158,264).

Exclusions were then applied to this age at menopause, as follows:

- Participants reporting undergoing a hysterectomy and/or oophorectomy, but not the age at which this happened (n=958 and 494, accordingly)
- Participants reporting multiple hysterectomy and/or oophorectomy ages, which were more than 10 years apart (n=38 and 23, accordingly)
- Participants reporting multiple HRT start and/or end ages, which were not in chronologically ascending order and were more than 10 years apart (n=124 and 137, accordingly)
- Participants reporting multiple ages at menopause, which were not in chronologically ascending order and were more than 10 years apart (n=73) and participants who reported both having and not having been through menopause and no other interventions (n=98)
- Participants having undergone a hysterectomy/oophorectomy before or during the year they report undergoing menopause
- Participants starting HRT prior to undergoing menopause and participants reporting HRT use, with no accompanying dates

The resulting trait was representative of an age at natural menopause (ANM, n=115,051) and was used in downstream analyses. Two additional ANM traits were also calculated, windsorized one by coding everyone reporting an ANM younger than 34, as 34 used in the discovery analysis as the primary phenotype (n=115,051 total, reduced to 106,973 after covariate-resulting exclusions), and one by only including participants reporting ANM between 40 and 60, inclusive (n=104,506), treated as a sensitivity analysis.

All manipulations were conducted in R (v4.1.2) on the UKB Research Analysis Platform (RAP; <https://ukbiobank.dnanexus.com/>).

## Phenome-wide association analysis

In order to test the association of ANM identified genes in other phenotypes, we processed additional reproductive ageing-related phenotypes, including age at menarche, cancer, telomere length (TL) and sex hormones (SH). All tested phenotypes were run as either continuous (age at menarche, TL and SH) or binary traits (cancer) corrected by age, age<sup>2</sup>, sex, the first ten genetic principal components as calculated in Bycroft *et al.*<sup>64</sup>, and study participant ES batch as a categorical covariate (either 50k, 200k, or 450k). Phenotype definitions and processing used in this study are described in **Supplementary Tables 8 and 9**. Only the first instance (initial visit) was used for generating all phenotype definitions unless specifically noted in **Supplementary Table 8**. In case of cancer-specific analysis data from cancer registries, death records, hospital admissions and self-reported were harmonised to ICD10 coding. If a participant had a code for any of the cancers recorded in ICD10 (C00-C97) then they were counted as a case for this phenotype. Minimal filtering was performed on the data, with only those cases where a diagnosis of sex-specific cancer was given in contrast to the sex data contained in UK Biobank record 31, was a diagnosis not used. For more information on cancer-specific analysis refer to **Supplementary Tables 9 and 11**.

## Cancer PheWAS Associations

To test for an association between genes we identified as associated with menopause timing (**Supplementary Table 2, Figure 1**) and 90 individual cancers as included in cancer registries, death records, hospital admissions and self-reported data provided by UK Biobank (e.g. breast, prostate, etc.) we utilised a logistic model with identical covariates as used during gene burden testing (N = 2430 tests) (**Supplementary Tables 9 and 11**). As standard logistic regression can lead to inflated *P* value estimates in cases of severe case/control imbalance<sup>91</sup>, we also performed a logistic regression with penalised likelihood estimation as described by Firth<sup>34</sup> (**Supplementary Table 12**). Models were run as discussed in Kosmidis *et al.*<sup>92</sup> using the 'brglm2' package implemented in R. brglm2 was run via the 'glm' function with default parameters other than "family" set to "binomial", "method" set to "brglmFit", and "type" set to "AS\_mean".

## WES sensitivity analysis using REGENIE

To replicate the primary findings and account for potential bias that could be introduced by exclusively using one discovery approach, a second analyst independently derived the age at menopause phenotype using a previously published method<sup>93</sup> and conducted additional burden association analysis using the REGENIE regression algorithm (REGENIEv2.2.4; <https://github.com/rgcgithub/regenie>). REGENIE implements a generalised mixed-model region-based association test that can account for population stratification and sample relatedness in large-scale analyses. REGENIE runs in 2 steps<sup>94</sup>, which we implemented on the UKBiobank RAP: In the first step, genetic variants are aggregated into gene specific units for each class of variant called masks. We selected variants in CCDS transcripts deemed to be high confidence by LOFTEE<sup>88</sup> with MAF<0.1% and annotated using VEP<sup>86</sup>. We created three masks,



independently of primary analysis group: (1) loss-of-function (LOF) variants (stop-gain, frameshift, or abolishing a canonical splice site (-2 or +2 bp from exon, excluding the ones in the last exon)) or missense variants with CADD score >30, (2) LOF or missense variants with CADD score >25, (3) all missense variants. In the second step, the three masks were tested for association with ANM. We applied an inverse normal rank transformation to ANM and included recruitment centre, sequence batch and 40 principal components as covariates. For each gene, we present results for the transcript with the smallest burden *P* value. The results for the sensitivity analysis performed via REGENIE are available in **Supplementary Table 1**.

## Common variant GWAS lookups

Genes within 500kb upstream and downstream of the 290 lead SNPs from the latest GWAS of ANM<sup>1</sup> were extracted from the exome-wide analysis. There were a total of 2149 genes within the GWAS regions. Burden tests in these genes with a Bonferroni corrected *P* value of  $<2.3 \times 10^{-5}$  (0.05/2149) were highlighted. The results are available in **Supplementary Table 4**.

## Analysis of GWAS and WES genes expression profiles in human female germ cells at various stages of development

We studied the mRNA abundance of WES genes during various stages of human female germ cell development using single-cell RNA sequencing data (**Supplementary Tables 6 and 7**). We used the processed single cell RNA resequencing datasets from two published studies. This included single-cell RNA sequencing data from foetal primordial germ cells of human female embryos (Accession code: GSE86146<sup>95</sup>), and from oocyte and granulosa cell fractions during various stages of follicle development (Accession code: GSE107746<sup>96</sup>). A pseudo score of 1 was added to all values before log transformation of the dataset. The samples from fetal germ cells (FGCs) were categorised into sub-clusters as defined in the original study. The study by Li *et al*<sup>95</sup> had identified 17 clusters by performing a t-distributed stochastic neighbour embedding (t-SNE) analysis and using expression profiles of known marker genes for various stages of fetal germ cell development. In our analysis we have included four clusters of female FGCs (Mitotic, Retinoic Acid (RA) responsive, Meiotic, Oogenesis) and four clusters containing somatic cells in the fetal gonads (Endothelial, Early\_Granulosa, Mural\_Granulosa, Late\_Granulosa). Software packages for R - tidyverse (<https://www.tidyverse.org/>), pheatmap, (<https://CRAN.R-project.org/package=pheatmap>), reshape2 (<https://github.com/hadley/reshape>), were used in processing and visualising the data.

## Functional enrichment tests for ZNF518A transcription factor binding sites using fGWAS and SLDP

fGWAS (v.0.3.6), a hierarchical model for joint analysis of GWAS and genomic annotations, was implemented to test the functional enrichment of ANM GWAS hits in ZNF518A transcription factor binding sites<sup>26</sup>. The fGWAS input file contained the ANM GWAS summary stats derived from the Reprogen study<sup>1</sup> annotated for ZNF518A binding sites. The ZNF518A annotation file

was derived from the ENCODE ChIP-seq data from human HEK293 cell line<sup>97</sup> the optimal independent discovery rate peak calling against hg19 [ENCF415VBF] was used. The ANM GWAS hits were annotated for the presence/absence of the *ZNF518A* transcription factor binding sites in a binary way (0, 1), with '1' if the SNP falls within the transcription factor binding site and '0' otherwise. The fGWAS tool available from <https://github.com/joepickrell/fgwas> and was run in annotation mode "-w" for the describe *ZNF518A* annotation. Detailed description of fGWAS methodology is available in Pickrell *et al*, 2014<sup>26</sup>. In short, the genome is split into independent blocks, which are allowed to contain either a single polymorphism that causally influences the trait or none. fGWAS then models the prior probability that any given block contains an association and the conditional prior probability that any given SNP in the block is the causal one, with probabilities allowed to vary according to functional annotations. The priors are then estimated using an empirical Bayes approach. The fGWAS output contained the maximum likelihood parameter estimates for each parameter in the model, in this case *ZNF518A*, with the lower and upper bound of the 95% confidence interval (CI) on the parameter. The P value was calculated from lower and upper CI in 3 following steps: (1) Standard error (SE) calculation:  $SE = (Upper\ CI - Lower\ CI)/(2*1.96)$ ; (2) Test statistics calculation:  $Z = Estimate / SE$ ; and (3) P value calculation:  $P = exp(-0.717*Z - 0.416*Z^2)$ .

Signed LD profile (SLDP) regression was applied to explore the directional effect of a signed functional annotation, *ZNF518A*, on a heritable trait like ANM using GWAS summary statistics. More specifically, we tested whether alleles that are predicted to increase the binding of the transcription factor *ZNF518A* have a genome-wide tendency to increase or decrease timing of menopause in women. The SLDP tool was installed from <https://github.com/yakirr/sldp>, with the comprehensive methodological steps described in Reshef *et al*, 2018<sup>27</sup>. For the analysis to be conducted, SLDP required GWAS summary statistics for ANM, signed LD profiles for *ZNF518A* binding, signed background model and reference panel in a SLDP compatible format. For the reference we used a 1000 Genomes Phase 3 European reference panel in *plink* format, which contained approximately 10M SNPs and 500 people and was available for download at the '[refpanel](#)' page. The ANM GWAS summary statistics, available from our latest Reprogen study<sup>1</sup>, was pre-processed using the '*preprocesspheno*' tool from the SLDP package. To conduct this step, we also obtained the list of regression SNPs along with the LD scores for the reference panel from the '[refpanel](#)' page. The pre-processing step included filtering down to SNPs that are also present in the reference panel, harmonising alleles to the reference, and multiplying the summary statistics by the SLDP regression weights. In addition, we applied the '*preprocessrefpanel*' tool to compute a truncated singular value decomposition (SVD) for each LD block in the reference panel. These SVDs were later used to weight the SLDP regression. The *ZNF518A* annotation file was obtained from the ENCODE CHIP-seq analysis, as described above, and preprocessed using the '*preprocessanno*' tool that turns signed functional annotations into signed LD profiles. Prior to running SLDP, we also obtained the signed background LD profiles that enabled us to control for systematic signed effects of minor alleles, which could arise from either population stratification or negative selection. SLDP was then run on our data using '*sldp*' function. To explore the relevance of *ZNF518A* for menopause timing in comparison to other transcription regulators, we tested whether genome-wide sequence changes introduced by SNP alleles identified in ANM GWAS increase or decrease binding of

additional 382 transcription factors (TFs). The preprocessed annotation files for 382 TFs derived from ENCODE CHIP-seq experiments, were available for download at the [annotation data page](#). The results are available in **Supplementary Table 5**.

## Functional analysis of ZNF518a binding sites

*ZNF518A* peaks were derived from unique genomic regions in ENCODE accession ENCFF415VBF described above. Quantification of ChIP-seq signal by aligning paired-end replicates (ENCFF174HBR, ENCFF574GQY, ENCFF808AJP, ENCFF453FDD) to the hg19 genome with Bowtie2 v2.3.5.1<sup>98</sup> with options “-I 0 -X 1000 –no-discordant –no-mixed”, reads were filtered for those with MAPQ > 30 with samtools v1.10. Assessment of H3K27ac<sup>29</sup> and chromatin accessibility by ATAC-seq<sup>30</sup> in day 4 human primordial germ cell like cells (hPGCLCs) at *ZNF518A* peaks was performed. For H3K27ac single end reads from accessions GSM4257216, GSM4257217, GSM4257218 were obtained and aligned with Bowtie2 v2.3.5.1 with default settings and MAPQ > 30 reads retained as above. For ATAC-seq paired-end reads were obtained from accessions GSM3406938, GSM3406939 and mapped and filtered as *ZNF518A* reads above.

Quantification of ChIP-seq and ATAC-seq signals for peak heights, heatmaps was performed with <https://github.com/owensnick/GenomeFragments.jl>. Peak to TSS distances were calculated against Gencode v36 release liftover to hg19 using GenomicFeatures.jl and <https://github.com/owensnick/ProximityEnrichment.jl>. We consider four categories of peaks: TSS intersecting, TSS proximal (TSS < 2000kb, outside gene body), Gene body intersecting, Intergenic and Distal (TSS > 5kb).

To perform *de novo* motif discovery we used Homer v4.11.1<sup>99</sup> using findMotifsGenome.pl with options “hg19 -size 200”. We ran this on all *ZNF518A* peaks, distal peaks and those intersecting TSS, we recovered a motif matching JASPAR<sup>28</sup> unvalidated motif UN0199.1 in all peak sets apart from those intersecting TSS. We then used <https://github.com/exeter-tfs/MotifScanner.jl> to quantify the occurrence of all instances of motif UN0199.1 in *ZNF518A* peaks. We downloaded the 18-state ChromHMM<sup>100</sup> models for all 833 biosamples in Epimap<sup>31</sup> from <http://compbio.mit.edu/epimap/>. We calculated the intersection between each state in each biosample and either all *ZNF518A* peaks or distal *ZNF518A* peaks using GenomicFeatures.jl. We calculated odds ratios from contingency tables using the approximation of bedtools<sup>101</sup> and Giggle<sup>102</sup>, by estimating total genomic intervals as hg19 genome size divided by the sum of the mean *ZNF518A* peak size and the chromatin state interval size.

## *De novo* mutation rate analyses

We calculated polygenic scores (PGSs) in participants from the rare disease programme of the 100,000 Genome Project (100kGP) v14. There are 77,901 individuals in the Aggregated Variant Calls (aggV2) after excluding participants whose genetically inferred sex is not consistent with

their phenotypic sex. We restricted the PGS analysis to individuals of European ancestry, which was predicted by the Genomics England Bioinformatics team using a random forest model based on genetic principal components (PCs) generated by projecting aggV2 data onto the 1000 Genomes phase 3 PC loadings. We removed one sample in each pair of related probands with kinship coefficient  $> 1/(2^{4.5})$ , i.e. up to and including third degree relationships. Probands with the highest number of relatives were removed first. Similarly, we retained unrelated mothers and fathers of these unrelated probands. It left us with 8,089 mother-offspring duos and 8,029 father-offspring duos.

We used the lead variants (or proxies, as described below) for genome-wide significant loci previously reported for ANM<sup>1</sup> to calculate PGS in the parents. In 100kGP, we removed variants with minor allele frequency (MAF)  $< 0.5\%$  or missing rate  $> 5\%$  from the aggV2 variants prepared by the Genomics England bioinformatics team. For lead variants that did not exist in 100kGP, we used the most significant proxy variants with linkage disequilibrium (LD)  $r^2 > 0.5$  if available in 100kGP. This resulted in a PGS constructed from 287 of the 290 previously reported loci. We regressed out 20 genetic PCs that were calculated within the European subset from the PGS and scaled the residuals to have mean = 0 and standard deviation = 1. Higher PGS indicates later age at menopause.

De novo mutations (DNMs) were called in 10,478 parent offspring trios by the Genomics England Bioinformatics team. The detailed analysis pipeline is documented at: <https://research-help.genomicsengland.co.uk/display/GERE/De+novo+variant+research+dataset>. Extensive quality control (QC) and filtering were applied by Kaplanis *et al.* as described previously<sup>103</sup>. De novo single nucleotide variants (dnSNVs) were phased using a read-based approach based on heterozygous variants near the DNM that were able to be phased to a parent. About one third of the dnSNVs were phased, of which three quarters were paternally phased (**Supplementary Figure 8, Supplementary Table 12**).

In association models, we accounted for parental age, the primary determinant of the number of DNMs, and various data quality metrics as described in<sup>103</sup>:

- Mean coverage for the child, mother and father (child\_mean\_RD, mother\_mean\_RD, father\_mean\_RD)
- Proportion of aligned reads for the child, mother and father (child\_prop\_aligned, mother\_prop\_aligned, father\_prop\_aligned)
- Number of SNVs called for child, mother and father (child\_SNVs, mother\_SNVs, father\_SNVs)
- Median variant allele fraction of DNMs called in child (median\_VAF)
- Median 'Bayes Factor' as outputted by Platypus for DNMs called in the child. This is a metric of DNM quality (median\_BF).

We first tested the association between parental PGSs and total dnSNV count in the offspring in a Poisson regression:

$$dnSNVs_{total} = \beta_0 + \beta_1 paternal\_PGS + \beta_2 maternal\_PGS +$$

$$\beta_3 \text{paternal\_age} + \beta_4 \text{maternal\_age} + \\ \beta_5 \text{child\_mean\_RD} + \beta_6 \text{mother\_mean\_RD} + \beta_7 \text{father\_mean\_RD} + \\ \beta_8 \text{child\_prop\_aligned} + \beta_9 \text{mother\_prop\_aligned} + \beta_{10} \text{father\_prop\_aligned} + \\ \beta_{11} \text{child\_snvs} + \beta_{12} \text{mother\_snvs} + \beta_{13} \text{father\_snvs} + \\ \beta_{14} \text{median\_VAF} + \beta_{15} \text{median\_BF}$$

We also fitted Poisson regression models to test the association between the PGS of one of the parents and the dnSNVs in the offspring that were phased to the relevant parent.

The paternal model included paternal PGS, age, and data quality metrics that are related to the proband and the father:

$$\text{dnSNVs\_paternal} = \beta_0 + \beta_1 \text{paternal\_PGS} + \beta_2 \text{paternal\_age} + \\ \beta_3 \text{child\_mean\_RD} + \beta_4 \text{father\_mean\_RD} + \\ \beta_5 \text{child\_prop\_aligned} + \beta_6 \text{father\_prop\_aligned} + \\ \beta_7 \text{child\_snvs} + \beta_8 \text{father\_snvs} + \\ \beta_9 \text{median\_VAF} + \beta_{10} \text{median\_BF}$$

Similarly, the maternal model was as follows:

$$\text{dnSNVs\_maternal} = \beta_0 + \beta_1 \text{maternal\_PGS} + \beta_2 \text{maternal\_age} + \\ \beta_3 \text{child\_mean\_RD} + \beta_4 \text{mother\_mean\_RD} + \\ \beta_5 \text{child\_prop\_aligned} + \beta_6 \text{mother\_prop\_aligned} + \\ \beta_7 \text{child\_snvs} + \beta_8 \text{mother\_snvs} + \\ \beta_9 \text{median\_VAF} + \beta_{10} \text{median\_BF}$$

Finally, as a sanity check, we assessed the association between the maternal PGS and paternally phased dnSNVs, and vice versa:

$$\text{dnSNVs\_paternal} = \beta_0 + \beta_1 \text{maternal\_PGS} + \beta_2 \text{paternal\_age} + \\ \beta_3 \text{child\_mean\_RD} + \beta_4 \text{father\_mean\_RD} + \\ \beta_5 \text{child\_prop\_aligned} + \beta_6 \text{father\_prop\_aligned} + \\ \beta_7 \text{child\_snvs} + \beta_8 \text{father\_snvs} + \\ \beta_9 \text{median\_VAF} + \beta_{10} \text{median\_BF}$$

$$\text{dnSNVs\_maternal} = \beta_0 + \beta_1 \text{paternal\_PGS} + \beta_2 \text{maternal\_age} + \\ \beta_3 \text{child\_mean\_RD} + \beta_4 \text{mother\_mean\_RD} + \\ \beta_5 \text{child\_prop\_aligned} + \beta_6 \text{mother\_prop\_aligned} + \\ \beta_7 \text{child\_snvs} + \beta_8 \text{mother\_snvs} + \\ \beta_9 \text{median\_VAF} + \beta_{10} \text{median\_BF}$$

## Mendelian Randomization

### Instrumental variable selection

MR analysis was applied to examine the likelihood of a causal effect of polygenic score (PGS) of age at natural menopause on the risk of de novo mutation rates in the offspring (**Supplementary Table 13**). In this approach, genetic variants that are significantly associated with an exposure of interest are used as instrumental variables (IVs) to test the causality of that exposure on the outcome of interest<sup>104–106</sup>. For a genetic variant to be a reliable instrument, the following assumptions should be met: (1) the genetic instrument is associated with the exposure of interest, (2) the genetic instrument should not be associated with any other competing risk factor that is a confounder, and (3) the genetic instrument should not be associated with the outcome, except via the causal pathway that includes the exposure of interest<sup>104,107</sup>. Genotypes at all variants were aligned to designate the ANM PGS-increasing alleles as the effect alleles as described above and this was used as a genetic instrument of interest. The effect sizes of genetic instruments (genotypes in the mother) on maternally phased de novo SNVs in the offspring estimated in 8,089 duos were obtained from Genomics England.

### MR Frameworks

The MR analysis was conducted using the inverse-variance weighted (IVW) model as the primary model due to the highest statistical power<sup>108</sup>. However, as it does not correct for heterogeneity in outcome risk estimates between individual variants<sup>109</sup>, we applied a number of sensitivity MR methods that better account for heterogeneity<sup>110</sup>. These include MR Egger to identify and correct for unbalanced heterogeneity ('horizontal pleiotropy'), indicated by a significant Egger intercept ( $P < 0.05$ )<sup>111</sup>, and weighted median (WM) and penalised weighted median (PWM) models to correct for balanced heterogeneity<sup>112</sup>. In addition, we introduced the MR Radial method to exclude variants from each model in cases where they are recognized as outliers<sup>113</sup>. The results were considered as significant based on the  $P$  value significance consistency across different primary and sensitivity models applied. The results are available in **Supplementary Table 13**. Finally, in order to calculate the effect of ANM on offspring de novo mutation rate when comparing women with ANM at two extremes of the ANM distribution curve, we multiplied the effect obtained by MR IVW, i.e. a de novo count beta per 1 year change in ANM, by 20, an arbitrary number that compares women with ANM 20 years apart.

## Acknowledgements

This work was funded by the Medical Research Council (Unit programs: MC\_UU\_12015/2, MC\_UU\_00006/2, MC\_UU\_12015/1, and MC\_UU\_00006/1). For the purpose of open access, the author has applied a Creative Commons Attribution (CC BY) licence to any Author Accepted Manuscript version arising. This research was conducted using the UK Biobank Resource under application 9905 (University of Cambridge) and 9072 and 871 (University of Exeter).

Ajuna Azad and Eva Hoffmann were supported by the ERC (724718-ReCAP), Novo Nordisk Foundation (NNF15COC0016662), the Independent Research Foundation Denmark (0134-00299B), and a grant from the Danish National Research Foundation Centre (6110-00344B).

Saleh Shekari was supported by the QUEX Institute (University of Exeter, UK and the University of Queensland, Australia). Anna Murray, Caroline Wright and Michael Weedon are supported by the Medical Research Council (MR/T00200X/1). The authors acknowledge the use of the University of Exeter High-Performance Computing facility in carrying out this work, funded by a MRC Clinical Research Infrastructure award (MRC Grant: MR/M008924/1).

This research was made possible through access to the data and findings generated by the 100,000 Genomes Project. The 100,000 Genomes Project is managed by Genomics England Limited (a wholly owned company of the Department of Health and Social Care). The 100,000 Genomes Project is funded by the National Institute for Health Research and NHS England. The Wellcome Trust, Cancer Research UK and the Medical Research Council have also funded research infrastructure. The 100,000 Genomes Project uses data provided by patients and collected by the National Health Service as part of their care and support.”

## Disclosures

John Perry is an employee and shareholder of Adrestia Therapeutics

## Data availability

All data produced using the UK Biobank resource will be returned to the UK Biobank returns catalogue upon publication.

## Extended authorship list

Ambrose, J. C.<sup>8</sup>; Arumugam, P.<sup>8</sup>; Bevers, R.<sup>8</sup>; Bleda, M.<sup>8</sup>; Boardman-Pretty, F.<sup>8,9</sup>; Boustred, C. R.<sup>8</sup>; Brittain, H.<sup>8</sup>; Brown, M.A.; Caulfield, M. J.<sup>8,9</sup>; Chan, G. C.<sup>8</sup>; Giess A.<sup>8</sup>; Griffin, J. N.; Hamblin, A.<sup>8</sup>; Henderson, S.<sup>8,9</sup>; Hubbard, T. J. P.<sup>8</sup>; Jackson, R.<sup>8</sup>; Jones, L. J.<sup>8,9</sup>; Kasperaviciute, D.<sup>8,9</sup>; Kayikci, M.<sup>8</sup>; Kousathanas, A.<sup>7?</sup>; Lahnstein, L.<sup>8</sup>; Lakey, A.; Leigh, S. E. A.<sup>8</sup>; Leong, I. U. S.<sup>8</sup>; Lopez, F. J.<sup>8</sup>; Maleady-Crowe, F.<sup>8</sup>; McEntagart, M.<sup>8</sup>; Minneci F.<sup>8</sup>; Mitchell, J.<sup>8</sup>; Moutsianas, L.<sup>8,9</sup>; Mueller, M.<sup>8,9</sup>; Murugaesu, N.<sup>8</sup>; Need, A. C.<sup>8,9</sup>; O'Donovan P.<sup>8</sup>; Odhams, C. A.<sup>8</sup>; Patch, C.<sup>8,9</sup>; Perez-Gil, D.<sup>8</sup>; Pereira, M. B.<sup>8</sup>; Pullinger, J.<sup>8</sup>; Rahim, T.<sup>8</sup>; Rendon, A.<sup>8</sup>; Rogers, T.<sup>8</sup>; Savage, K.<sup>8</sup>; Sawant, K.<sup>8</sup>; Scott, R. H.<sup>8</sup>; Siddiq, A.<sup>8</sup>; Sieghart, A.<sup>8</sup>; Smith, S. C.<sup>8</sup>; Sosinsky, A.<sup>8,9</sup>; Stuckey, A.<sup>8</sup>; Tanguy M.<sup>8</sup>; Taylor Tavares, A. L.<sup>8</sup>; Thomas, E. R. A.<sup>8,9</sup>; Thompson, S. R.<sup>8</sup>; Tucci, A.<sup>8,9</sup>; Welland, M. J.<sup>8</sup>; Williams, E.<sup>8</sup>; Witkowska, K.<sup>8,9</sup>; Wood, S. M.<sup>8,9</sup>; Zarowiecki, M.<sup>8</sup>.

## References

1. Ruth, K. S. *et al.* Genetic insights into biological mechanisms governing human ovarian ageing. *Nature* **596**, 393–397 (2021).
2. Perry, J. R. B., Murray, A., Day, F. R. & Ong, K. K. Molecular insights into the aetiology of female reproductive ageing. *Nat Rev Endocrinol* **11**, 725–734 (2015).
3. te Velde, E. R. & Pearson, P. L. The variability of female reproductive ageing. *Hum Reprod Update* **8**, 141–154 (2002).
4. te Velde, E. R., Dorland, M. & Broekmans, F. J. Age at menopause as a marker of reproductive ageing. *Maturitas* **30**, 119–125 (1998).
5. Oktem, O. & Oktay, K. The ovary: anatomy and function throughout human life. *Ann N Y Acad Sci* **1127**, 1–9 (2008).
6. Wallace, W. H. B. & Kelsey, T. W. Human ovarian reserve from conception to the menopause. *PLoS One* **5**, (2010).
7. Lambalk, C. B., van Disseldorp, J., de Koning, C. H. & Broekmans, F. J. Testing ovarian reserve to predict age at menopause. *Maturitas* **63**, 280–291 (2009).
8. Podfigurna-Stopa, A. *et al.* Premature ovarian insufficiency: the context of long-term effects. *J Endocrinol Invest* **39**, 983–990 (2016).
9. Ward, L. D. *et al.* Rare coding variants in DNA damage repair genes associated with timing of natural menopause. *HGG Adv* **3**, (2021).
10. Tam, V. *et al.* Benefits and limitations of genome-wide association studies. *Nat Rev Genet* **20**, 467–484 (2019).
11. He, C. *et al.* Genome-wide association studies identify loci associated with age at menarche and age at natural menopause. *Nat Genet* **41**, 724–728 (2009).
12. Lunetta, K. L. *et al.* Genetic correlates of longevity and selected age-related phenotypes: a genome-wide association study in the Framingham Study. *BMC Med Genet* **8 Suppl 1**, (2007).
13. Horikoshi, M. *et al.* Elucidating the genetic architecture of reproductive ageing in the Japanese population. *Nat Commun* **9**, (2018).
14. Stolk, L. *et al.* Loci at chromosomes 13, 19 and 20 influence age at natural menopause. *Nat Genet* **41**, 645–647 (2009).
15. Hartge, P. Genetics of reproductive lifespan. *Nat Genet* **41**, 637–638 (2009).
16. Perry, J. R. B. *et al.* A genome-wide association study of early menopause and the combined impact of identified variants. *Hum Mol Genet* **22**, 1465–1472 (2013).
17. Day, F. R. *et al.* Large-scale genomic analyses link reproductive aging to hypothalamic signaling, breast cancer susceptibility and BRCA1-mediated DNA repair. *Nat Genet* **47**, 1294–1303 (2015).
18. Stolk, L. *et al.* Meta-analyses identify 13 loci associated with age at menopause and highlight DNA repair and immune pathways. *Nat Genet* **44**, 260–268 (2012).
19. Perry, J. R. B. *et al.* DNA mismatch repair gene MSH6 implicated in determining age at natural menopause. *Hum Mol Genet* **23**, 2490–2497 (2014).



20. Murray, A. *et al.* Common genetic variants are significant risk factors for early menopause: results from the Breakthrough Generations Study. *Hum Mol Genet* **20**, 186–192 (2011).
21. Szustakowski, J. D. *et al.* Advancing human genetics research and drug discovery through exome sequencing of the UK Biobank. *Nat Genet* **53**, 942–948 (2021).
22. Loh, P. R. *et al.* Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat Genet* **47**, 284–290 (2015).
23. Maier, V. K. *et al.* Functional Proteomic Analysis of Repressive Histone Methyltransferase Complexes Reveals ZNF518B as a G9A Regulator. *Mol Cell Proteomics* **14**, 1435–1446 (2015).
24. Dunham, I. *et al.* An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
25. Davis, C. A. *et al.* The Encyclopedia of DNA elements (ENCODE): data portal update. *Nucleic Acids Res* **46**, D794–D801 (2018).
26. Pickrell, J. K. Joint analysis of functional genomic data and genome-wide association studies of 18 human traits. *Am J Hum Genet* **94**, 559–573 (2014).
27. Reshef, Y. A. *et al.* Detecting genome-wide directional effects of transcription factor binding on polygenic disease risk. *Nat Genet* **50**, 1483–1493 (2018).
28. Castro-Mondragon, J. A. *et al.* JASPAR 2022: the 9th release of the open-access database of transcription factor binding profiles. *Nucleic Acids Res* **50**, D165–D173 (2022).
29. Chen, D. *et al.* Human Primordial Germ Cells Are Specified from Lineage-Primed Progenitors. *Cell Rep* **29**, 4568-4582.e5 (2019).
30. Chen, D. *et al.* The TFAP2C-Regulated OCT4 Naive Enhancer Is Involved in Human Germline Formation. *Cell Rep* **25**, 3591-3602.e5 (2018).
31. Boix, C. A., James, B. T., Park, Y. P., Meuleman, W. & Kellis, M. Regulatory genomic circuitry of human disease loci by integrative epigenomics. *Nature* **590**, 300–307 (2021).
32. Day, F. R. *et al.* Genomic analyses identify hundreds of variants associated with age at menarche and support a role for puberty timing in cancer risk. *Nat Genet* **49**, 834–841 (2017).
33. Day, F. R. *et al.* Shared genetic aetiology of puberty timing between sexes and with health-related outcomes. *Nat Commun* **6**, (2015).
34. Firth, D. Bias reduction of maximum likelihood estimates. *Biometrika* **80**, 27–38 (1993).
35. Codd, V. *et al.* Polygenic basis and biomedical consequences of telomere length variation. *Nat Genet* **53**, 1425–1433 (2021).
36. Sherman, S. L. Premature ovarian failure in the fragile X syndrome. *Am J Med Genet* **97**, 189–194 (2000).
37. Reid, S. *et al.* Biallelic mutations in PALB2 cause Fanconi anemia subtype FA-N and predispose to childhood cancer. *Nat Genet* **39**, 162–164 (2007).
38. Bass, T. E. *et al.* ETAA1 acts at stalled replication forks to maintain genome integrity. *Nat Cell Biol* **18**, 1185–1195 (2016).
39. Feng, S. *et al.* Ewing Tumor-associated Antigen 1 Interacts with Replication Protein A to Promote Restart of Stalled Replication Forks. *J Biol Chem* **291**, 21956–21962 (2016).

40. Haahr, P. *et al.* Activation of the ATR kinase by the RPA-binding protein ETAA1. *Nat Cell Biol* **18**, 1196–1207 (2016).
41. Saldivar, J. C. *et al.* An intrinsic S/G 2 checkpoint enforced by ATR. *Science* **361**, 806–810 (2018).
42. Hustedt, N. *et al.* Control of homologous recombination by the HROB-MCM8-MCM9 pathway. *Genes and Development* **33**, 1397–1415 (2019).
43. Huang, J. W. *et al.* MCM8IP activates the MCM8-9 helicase to promote DNA synthesis and homologous recombination upon DNA damage. *Nat Commun* **11**, (2020).
44. Tucker, E. J. *et al.* Meiotic genes in premature ovarian insufficiency: variants in HROB and REC8 as likely genetic causes. *Eur J Hum Genet* **30**, 219–228 (2022).
45. Hara, S., Yoda, E., Sasaki, Y., Nakatani, Y. & Kuwata, H. Calcium-independent phospholipase A 2  $\gamma$  (iPLA 2  $\gamma$ ) and its roles in cellular functions and diseases. *Biochim Biophys Acta Mol Cell Biol Lipids* **1864**, 861–868 (2019).
46. Liu, G. Y. *et al.* The phospholipase iPLA 2  $\gamma$  is a major mediator releasing oxidized aliphatic chains from cardiolipin, integrating mitochondrial bioenergetics and signaling. *J Biol Chem* **292**, 10672–10684 (2017).
47. Mancuso, D. J. *et al.* Genetic ablation of calcium-independent phospholipase A2 $\gamma$  leads to alterations in mitochondrial lipid metabolism and function resulting in a deficient mitochondrial bioenergetic phenotype. *J Biol Chem* **282**, 34611–34622 (2007).
48. Shukla, A., Saneto, R. P., Hebbar, M., Mirzaa, G. & Girisha, K. M. A neurodegenerative mitochondrial disease phenotype due to biallelic loss-of-function variants in PNPLA8 encoding calcium-independent phospholipase A2 $\gamma$ . *Am J Med Genet A* **176**, 1232–1237 (2018).
49. Mancuso, D. J. *et al.* Genetic ablation of calcium-independent phospholipase A2 $\gamma$  leads to alterations in hippocampal cardiolipin content and molecular species distribution, mitochondrial degeneration, autophagy, and cognitive dysfunction. *J Biol Chem* **284**, 35632–35644 (2009).
50. Mancuso, D. J. *et al.* Genetic ablation of calcium-independent phospholipase A2 $\gamma$  prevents obesity and insulin resistance during high fat feeding by mitochondrial uncoupling and increased adipocyte fatty acid oxidation. *J Biol Chem* **285**, 36495–36510 (2010).
51. Saunders, C. J. *et al.* Loss of function variants in human PNPLA8 encoding calcium-independent phospholipase A2  $\gamma$  recapitulate the mitochondriopathy of the homologous null mouse. *Hum Mutat* **36**, 301–306 (2015).
52. Masih, S., Moirangthem, A. & Phadke, S. R. Homozygous Missense Variation in PNPLA8 Causes Prenatal-Onset Severe Neurodegeneration. *Mol Syndromol* **12**, 174–178 (2021).
53. Clifford, R. *et al.* SAMHD1 is mutated recurrently in chronic lymphocytic leukemia and is involved in response to DNA damage. *Blood* **123**, 1021–1031 (2014).
54. Schott, K. *et al.* SAMHD1 in cancer: curse or cure? *J Mol Med (Berl)* **100**, 351–372 (2022).
55. Sjöblom, T. *et al.* The consensus coding sequences of human breast and colorectal cancers. *Science* **314**, 268–274 (2006).
56. Parsons, D. W. *et al.* The genetic landscape of the childhood cancer medulloblastoma. *Science* **331**, 435–439 (2011).

57. Jones, S. *et al.* Core signaling pathways in human pancreatic cancers revealed by global genomic analyses. *Science* **321**, 1801–1806 (2008).
58. Imielinski, M. *et al.* Mapping the hallmarks of lung adenocarcinoma with massively parallel sequencing. *Cell* **150**, 1107–1120 (2012).
59. Herold, N. *et al.* With me or against me: Tumor suppressor and drug resistance activities of SAMHD1. *Exp Hematol* **52**, 32–39 (2017).
60. Coggins, S. A., Mahboubi, B., Schinazi, R. F. & Kim, B. SAMHD1 Functions and Human Diseases. *Viruses* **12**, (2020).
61. Chen, Z., Hu, J., Ying, S. & Xu, A. Dual roles of SAMHD1 in tumor development and chemoresistance to anticancer drugs. *Oncol Lett* **21**, (2021).
62. Xagoraris, I. *et al.* Expression of the novel tumour suppressor sterile alpha motif and HD domain-containing protein 1 is an independent adverse prognostic factor in classical Hodgkin lymphoma. *Br J Haematol* **193**, 488–496 (2021).
63. Merati, M. *et al.* Aggressive CD8(+) epidermotropic cutaneous T-cell lymphoma associated with homozygous mutation in SAMHD1. *JAAD Case Rep* **1**, 227–229 (2015).
64. Rice, G. I. *et al.* Mutations involved in Aicardi-Goutières syndrome implicate SAMHD1 as regulator of the innate immune response. *Nat Genet* **41**, 829–832 (2009).
65. Franzolin, E. *et al.* The deoxynucleotide triphosphohydrolase SAMHD1 is a major regulator of DNA precursor pools in mammalian cells. *Proc Natl Acad Sci U S A* **110**, 14272–14277 (2013).
66. Kumar, D. *et al.* Mechanisms of mutagenesis in vivo due to imbalanced dNTP pools. *Nucleic Acids Res* **39**, 1360–1371 (2011).
67. Coquel, F. *et al.* SAMHD1 acts at stalled replication forks to prevent interferon induction. *Nature* **557**, 57–61 (2018).
68. Daddacha, W. *et al.* SAMHD1 Promotes DNA End Resection to Facilitate DNA Repair by Homologous Recombination. *Cell Rep* **20**, 1921–1935 (2017).
69. Mathews, C. K. Deoxyribonucleotide metabolism, mutagenesis and cancer. *Nat Rev Cancer* **15**, 528–539 (2015).
70. Mao, Z., Bozzella, M., Seluanov, A. & Gorbunova, V. DNA repair by nonhomologous end joining and homologous recombination during cell cycle in human cells. *Cell Cycle* **7**, 2902–2906 (2008).
71. Bonifati, S. *et al.* SAMHD1 controls cell cycle status, apoptosis and HIV-1 infection in monocytic THP-1 cells. *Virology* **495**, 92–100 (2016).
72. Kodigepalli, K. M., Li, M., Liu, S. L. & Wu, L. Exogenous expression of SAMHD1 inhibits proliferation and induces apoptosis in cutaneous T-cell lymphoma-derived HuT78 cells. *Cell Cycle* **16**, 179–188 (2017).
73. Rahbari, R. *et al.* Timing, rates and spectra of human germline mutation. *Nat Genet* **48**, 126–133 (2016).
74. Crow, J. F. The origins, patterns and implications of human spontaneous mutation. *Nat Rev Genet* **1**, 40–47 (2000).
75. Iossifov, I. *et al.* The contribution of de novo coding mutations to autism spectrum disorder. *Nature* **515**, 216–221 (2014).
76. Fromer, M. *et al.* De novo mutations in schizophrenia implicate synaptic networks. *Nature* **506**, 179–184 (2014).

77. Wang, W., Corominas, R. & Lin, G. N. De novo Mutations From Whole Exome Sequencing in Neurodevelopmental and Psychiatric Disorders: From Discovery to Application. *Front Genet* **10**, (2019).
78. Coe, B. P. *et al.* Neurodevelopmental disease genes implicated by de novo mutation and copy number variation morbidity. *Nat Genet* **51**, 106–116 (2019).
79. Linschooten, J. O. *et al.* Paternal lifestyle as a potential source of germline mutations transmitted to offspring. *FASEB J* **27**, 2873–2879 (2013).
80. Titus, S. *et al.* Impairment of BRCA1-related DNA double-strand break repair leads to ovarian aging in mice and humans. *Sci Transl Med* **5**, (2013).
81. Miao, Y. *et al.* BRCA2 deficiency is a potential driver for human primary ovarian insufficiency. *Cell Death Dis* **10**, (2019).
82. Lin, W., Titus, S., Moy, F., Ginsburg, E. S. & Oktay, K. Ovarian Aging in Women With BRCA Germline Mutations. *J Clin Endocrinol Metab* **102**, 3839–3847 (2017).
83. Backman, J. D. *et al.* Exome sequencing and analysis of 454,787 UK Biobank participants. *Nature* **599**, 628–634 (2021).
84. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
85. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *Gigascience* **10**, (2021).
86. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol* **17**, (2016).
87. Rentzsch, P., Schubach, M., Shendure, J. & Kircher, M. CADD-Splice-improving genome-wide variant effect prediction using deep learning-derived splice scores. *Genome Med* **13**, (2021).
88. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
89. Loh, P. R. *et al.* Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat Genet* **47**, 284–290 (2015).
90. Seabold, S. P. J. Statsmodels: Econometric and Statistical Modeling with Python. *SCIPY* (2010).
91. Ma, C., Blackwell, T., Boehnke, M. & Scott, L. J. Recommended joint and meta-analysis strategies for case-control association testing of single low-count variants. *Genet Epidemiol* **37**, 539–550 (2013).
92. Kosmidis, I. K. P. E. C. S. N. Mean and median bias reduction in generalized linear models. *Statistics and Computing* **30**, 43–59 (2019).
93. Ruth, K. S. *et al.* Events in Early Life are Associated with Female Reproductive Ageing: A UK Biobank Study. *Sci Rep* **6**, (2016).
94. Mbatchou, J. *et al.* Computationally efficient whole-genome regression for quantitative and binary traits. *Nat Genet* **53**, 1097–1103 (2021).
95. Li, L. *et al.* Single-Cell RNA-Seq Analysis Maps Development of Human Germline Cells and Gonadal Niche Interactions. *Cell Stem Cell* **20**, 858-873.e4 (2017).
96. Zhang, Y. *et al.* Transcriptome Landscape of Human Folliculogenesis Reveals Oocyte and Granulosa Cell Interactions. *Mol Cell* **72**, 1021-1034.e4 (2018).
97. Lou, S. *et al.* TopicNet: a framework for measuring transcriptional regulatory network change. *Bioinformatics* **36**, I474–I481 (2020).

98. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**, 357–359 (2012).
99. Heinz, S. *et al.* Simple combinations of lineage-determining transcription factors prime cis-regulatory elements required for macrophage and B cell identities. *Mol Cell* **38**, 576–589 (2010).
100. Ernst, J. & Kellis, M. ChromHMM: automating chromatin-state discovery and characterization. *Nat Methods* **9**, 215–216 (2012).
101. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).
102. Layer, R. M. *et al.* GIGGLE: a search engine for large-scale integrated genome analysis. *Nat Methods* **15**, 123–126 (2018).
103. Kaplanis, J. *et al.* Genetic and chemotherapeutic causes of germline hypermutation. *bioRxiv* (2021).
104. Smith, G. D. & Ebrahim, S. “Mendelian randomization”: can genetic epidemiology contribute to understanding environmental determinants of disease? *Int J Epidemiol* **32**, 1–22 (2003).
105. Burgess, S., Foley, C. N., Allara, E., Staley, J. R. & Howson, J. M. M. A robust and efficient method for Mendelian randomization with hundreds of genetic variants. *Nat Commun* **11**, (2020).
106. Lawlor, D. A., Harbord, R. M., Sterne, J. A. C., Timpson, N. & Smith, G. D. Mendelian randomization: using genes as instruments for making causal inferences in epidemiology. *Stat Med* **27**, 1133–1163 (2008).
107. Burgess, S. *et al.* Guidelines for performing Mendelian randomization investigations. *Wellcome Open Res* **4**, (2020).
108. Slob, E. A. W. & Burgess, S. A comparison of robust Mendelian randomization methods using summary data. *Genet Epidemiol* **44**, 313–329 (2020).
109. Bowden, J. *et al.* A framework for the investigation of pleiotropy in two-sample summary data Mendelian randomization. *Stat Med* **36**, 1783–1802 (2017).
110. Burgess, S., Bowden, J., Fall, T., Ingelsson, E. & Thompson, S. G. Sensitivity Analyses for Robust Causal Inference from Mendelian Randomization Analyses with Multiple Genetic Variants. *Epidemiology* **28**, 30–42 (2017).
111. Bowden, J., Smith, G. D. & Burgess, S. Mendelian randomization with invalid instruments: effect estimation and bias detection through Egger regression. *Int J Epidemiol* **44**, 512–525 (2015).
112. Bowden, J., Davey Smith, G., Haycock, P. C. & Burgess, S. Consistent Estimation in Mendelian Randomization with Some Invalid Instruments Using a Weighted Median Estimator. *Genet Epidemiol* **40**, 304–314 (2016).
113. Bowden, J. *et al.* Improving the visualization, interpretation and analysis of two-sample summary data Mendelian randomization via the Radial plot and Radial regression. *Int J Epidemiol* **47**, 1264–1278 (2018).