

1 **Gene expression profiling of peripheral blood and endometrial cancer risk**
2 **factors: systems epidemiology approach in the NOWAC Postgenome**
3 **Cohort Study.**

4 Oxana Gavriilyuk¹, Igor Snapkov^{1,2}, Jean-Christophe Thalabard³, Lars Holden⁴, Marit
5 Holden⁴, Hege M. Bøvelstad¹, Vanessa Dumeaux⁵, Eiliv Lund^{1,6}

6 **Corresponding author:** Oxana A. Gavriilyuk; Institutt for samfunnsmedisin, Det
7 helsevitenskapelige fakultet, UiT Norges arktiske universitet, 9037 Tromsø, Norway; phone:
8 +47 77645126; e-mail: oxana.gavriilyuk@uit.no

9 **Running title: Gene expression and endometrial cancer risk**

10 **Keywords:** Endometrial cancer; gene expression; systems epidemiology; risk factors;
11 prospective study;

12 **Abbreviations:**

13 LNYM - lifetime number of years of menstruation

14 BMI - body mass index

15 CI - confidence interval

16 EC - endometrial cancer

17 MHT - menopausal hormone therapy

18 NOWAC - The Norwegian Women and Cancer Study

19 OC - oral contraceptives

20 FDR - false discovery rate

21 GSEA - gene set enrichment analysis

22 **Affiliations:**

23 1. Department of Community Medicine, Faculty of Health Sciences, UiT, The Arctic University of
24 Norway, Tromsø, Norway

- 1 2. Department of Immunology, University of Oslo and Oslo University Hospital-Rikshospitalet, Oslo,
- 2 Norway
- 3 3. MAP5, MRC CNRS 8145, Université Paris Descartes, Sorbonne Paris Cite, Paris, France MAP5, MRC
- 4 CNRS 8145, Université Paris Descartes, Sorbonne Paris Cite, Paris, France
- 5 4. Norwegian Computing Center, Oslo, Norway
- 6 5. PERFORM Center, Concordia University, Montreal, Canada
- 7 6. The Cancer Registry of Norway, Oslo, Norway

8

9 **Funding:** This study was supported by funding from the Northern Norway Regional Health
10 Authority (Helse-Nord RHF) and by a grant from the European Research Council (ERC-AdG
11 232997 TICE).

12 **Conflict of interest:** The authors declare that they have no conflict of interest.

13 **Author's contributions:**

14 OG, EL, VD, and IS designed the study and interpreted the results. EL is a PI of the NOWAC
15 Study. OG and IS constructed the tables and drafted the manuscript. MH and LH carried out
16 the statistical analyses and contributed to drafting the manuscript. HMB participated in initial
17 statistical analysis and critically revised the manuscript. VD and JCT contributed to the
18 interpretation of the data, and the drafting and revision of the manuscript. All authors read and
19 approved the final manuscript.

20

21 **Abstract**

22 **Background:**

23 Increasing incidence of endometrial cancer (EC), the most common gynecologic cancer in the
24 world, requires extensive search for novel preventive tools and early intervention approaches.
25 Several factors, including parity status, breastfeeding duration, use of oral contraceptives

1 (OC), coffee consumption, BMI, use of hormone replacement therapy (HRT), and lifetime
2 number of years of menstruation have previously been reported to modify EC risk. However,
3 establishment of reliable predictive models is impossible without knowledge on genetic
4 changes prior to diagnosis. In this work, we aimed to establish if known EC risk factors
5 influence peripheral blood gene expression in a prospective design.

6 **Methods:**

7 First, we selected variables that were shown to have an impact on EC risk in the whole
8 Norwegian Women and Cancer (NOWAC) cohort (165 000 women). Then, we tested the
9 association between these variables and changes in gene expression profiles in blood in a
10 nested case-control study (79 case-control pairs) of women from the NOWAC postgenome
11 cohort. Lastly, we undertook a gene set enrichment analysis (GSEA).

12 **Results:**

13 When we looked at overall gene expression, we found no difference between EC cases and
14 controls. Introduction of parity status into the statistical model, revealed changes in
15 expression of 1379 genes (false discovery rate (FDR) 20%) in controls, while we did not
16 observe any expression changes in cases. 27 genes (FDR 20%) were associated with BMI
17 increase in controls, whereas there was no association between changes in BMI and gene
18 expression in women with EC. In GSEA, the major part of significantly enriched gene sets
19 (2407, FDR 20%) were attributed to parity increase among cancer-free women.

20 **Conclusions:**

21 We found that increased number of parities and elevated BMI do not change peripheral blood
22 gene expression in women diagnosed with EC later in life. The descriptive study design does
23 not allow us to provide accurate explanation of our findings in biologic terms but this work
24 brings solid background for further research on the development of predictive EC risk models.

25

1 **Introduction**

2 Endometrial cancer (EC) is s the most common and second most lethal gynecological cancer
3 among women worlwide with 774 new cases registered in 2016 in Norway [1, 2]. While the
4 incidence and mortality rates of several other cancers have plateaued or decreased in the last
5 decade, the incidence of EC has been rising globally, with the highest increase found in
6 countries that have undergone a drastic change in living standards and lifestyles [3, 4]. In
7 particular, changes in reproductive factors (eg declines in parity) combined with increase in
8 obesity prevalence might explain the rise in EC incidence associated with socioeconomic
9 transition [5].

10 Both descriptive and analytic epidemiological approaches can help health officials
11 appropriately target prevention and control activities, but provide only limited information on
12 the biological processes underlying the cause-effect relationship between a risk factor and the
13 disease. Technological advances in genomic profiling provide epidemiologists with the
14 opportunity to integrate molecular data into etiologic studies in order to decypher the
15 carcinogenic process.

16 So far, most studies have been focused on the development of tools to reclassify
17 endometrial tumors according to their molecular features and stratify patients according to
18 risk of metastases and recurrence [6-8]. We recently showed that an in-depth analysis of
19 multi-tissue genomic profiles could be a crucial addition to the personalized assessment of an
20 individual's biology and health [9].

21 As a major defense and transport system, blood cells can adjust expression of their
22 genes in response to various environmental factors and pathological conditions. We
23 previously highlighted several specific behavioral programs, such as metabolism or signaling,
24 deregulated in the individual's blood cells that are associated with biological and/or

1 pathological responses to a given condition in the general population (eg smoking) [10],
2 cancer-related risk factors (article in press), and health status (eg breast cancer diagnosis) [12,
3 13].

4 In this study, we aim to investigate the associations between known EC risk factors
5 and gene expression changes in blood cells differential between women who will develop EC
6 and controls. This will provide insight into biological processes underlying lifestyle/exposure
7 EC risk factors that might explain incidence of EC.

8

9 **Material and methods**

10 **The NOWAC Study**

11 The NOWAC Study is a national population-based cohort study which include Norwegian
12 women aged between 30-70 randomly drawn from the Norwegian Central Population Register
13 [14]. Starting from 1991 and within 4-6 years intervals, these women filled in questionnaires
14 with focus on lifestyle and health. Of this original cohort of about 172 000 women, more than
15 45 000 women born between 1943 and 1957 provided blood samples between 2003 and 2006
16 and filled an additional 2-page questionnaire to constitute the NOWAC Postgenome Cohort
17 [14]. PAXgene tubes (PreAnalytiX GmbH, Hembrechtikon, Switzerland) were used to prevent
18 RNA degradation after blood sampling and allowed genome-wide analyses of blood gene
19 expression profiles. Through the linkage to the Cancer Registry of Norway and Register of
20 death certificates in Statistics Norway, we have identified 88 women from the NOWAC
21 Postgenome Cohort who developed EC between the time of blood sampling and December
22 31, 2008 (end of follow-up). Of these 88 individuals, four were excluded, as blood samples
23 were not received and stored at -80C within 4 days after blood collection . To ensure the same
24 storage time and age between cases and controls, the controls, who did not receive any cancer

1 diagnosis, were drawn at random from the NOWAC Postgenome Cohort but matched by time
2 of blood collection and birth year. Matched case-control pairs of blood samples (n=84 pairs)
3 were sent to the Genomics Core Facility at the Norwegian University of Science and
4 Technology (NTNU) for microarray gene expression profiling in January 2011.

5 **Assessment of covariates and calculation of lifetime number of years of menstruation**

6 Information on the covariates age at menarche, age at menopause, number of full-term
7 pregnancies, duration of breastfeeding, height, weight, oral contraceptive use, and smoking
8 status was taken from NOWAC questionnaires (series from years 2002-2005). Self-reported
9 height and weight were used to calculate BMI in kg/m². Parity and breastfeeding variables are
10 generally reported to have a good validity in NOWAC Study [15]. The lifetime number of
11 years of menstruation (LNYM) count the number of years between age at menarche and age
12 at menopause, minus the cumulative duration of full-term pregnancies (calculated as the
13 number of full-term pregnancies, including live and stillbirths, times 0.75 years), duration of
14 breastfeeding (calculated as the cumulative number of months of breastfeeding in all
15 pregnancies), and duration of OC use [16].

16 **Laboratory procedures**

17 In order to minimize the technical variability, each control sample was processed together
18 with the matching case sample throughout RNA extraction, amplification and hybridization.
19 Total RNA was isolated using the PAXgene Blood RNA Isolation Kit (Preanalytix, Qiagen,
20 Hilden, Germany) following the manufacturer's instructions. RNA quantity and purity were
21 assessed by the NanoDrop ND1000 spectrophotometer (Thermo Scientific, Wilmington,
22 Delaware, USA) and Agilent bioanalyzer (Agilent Technologies, Palo Alto, CA, USA) RNA
23 amplification was performed in 96-well plates using 300 ng of total RNA and the Illumina®
24 TotalPrep™-96 RNA Amplification Kit (Ambion Inc., Austin, TX, USA). Genome-wide
25 RNA profiles were obtained using IlluminaHumanHT-12 chips version 4.

1 **Outlier removal**

2 Initial quality control made at the NTNU laboratory revealed two technical outliers (one 5S
3 type degradation and one failed cRNA synthesis). These were removed along with their
4 matching samples. Further, one case was later found to have two cancer diagnoses and was
5 excluded along with its matching control. The data was carefully investigated to identify and
6 remove technical outliers using the standard operating procedure for outlier removal described
7 in [17]. In particular, probes related to genes in the human leukocyte antigen (HLA) systems
8 were excluded (38 probes). These genes are known to be expressed strongly and have a high
9 variance, which could affect multivariate analyses used in the outlier search. For individuals
10 that were borderline outlier candidates, we excluded the ones with quality measures outside
11 the range of the following thresholds: RIN value <7 , 260/280 ratio <2 , 260/230 ratio <1.7 , and
12 $50 < \text{RNA} < 500$. Based on this outlier search we identified two additional individuals as
13 technical outliers and excluded them along with their matching individual. In total, we
14 investigated blood profiles from 79 women who developed EC and 79 age-matched controls.

15 **Microarray data preprocessing and normalization**

16 Microarray data preprocessing and analysis were performed using R v3.3.1 ([http://cran.r-](http://cran.r-project.org)
17 [project.org](http://cran.r-project.org)), and tools from the Bioconductor project (<http://www.bioconductor.org>), adapted
18 to our needs.

19 Expression profiles including 47 248 probes were adjusted for background noise using
20 the negative control probes [18]. The data was further log₂-transformed using a variance
21 stabilizing technique [19] and finally normalized using quantile normalization. We retained
22 probes present in at least 70% of the samples. If a gene was represented by more than one
23 probe, the average expression of the probes was used as expression value for the gene. The
24 probes were translated to genes using the lumiHumanIDMapping database [20]. In total, the

1 final data set included the expression for 7104 unique genes. Finally, we computed the
2 differences in log₂ gene expression levels for each case-control pair.

3 **Statistical methods**

4 We used linear models (Bioconductor R-package Limma [21]) to evaluate the significance of
5 gene-wise expression differences between cases and controls using the following intercept-
6 only model:

$$7 \quad Y_{g,c} = \alpha + \varepsilon_{g,c}$$

8 where $Y_{g,c} = Y_{g,c}^{case} - Y_{g,c}^{ctrl}$ is the difference in log₂ gene expression level for case-control
9 pair c and gene g and $\varepsilon_{g,c}$ is normally distributed with zero expectation.

10 Using the same approach, we tested whether gene expression changes between cases and
11 controls were associated with EC risk factors (parity, LNYM, coffee, BMI, age of menopause,
12 or OC use) using the following model:

$$13 \quad Y_{gc} = \alpha + \beta_V^{ctrl} V_c^{ctrl} + \beta_V^{case} V_c^{case} + \varepsilon_{g,c}$$

14 where V_c^{ctrl} and V_c^{case} are the parity, coffee, LNYM, BMI, age of menopause, or OC use
15 variable for case-control pair c in control and cases, respectively.

16 Similarly, we applied a gene set approach to determine changes in gene signatures
17 between EC cases and controls associated with known risk factors. In this approach, the
18 dependent variable is the difference in enrichment scores for case-control pair c and a defined
19 gene set. The enrichment scores for eight collections (C1-C7 and H) of gene sets from the
20 Molecular Signatures Database v6.1 (<http://software.broadinstitute.org/gsea/index.jsp>) [22]
21 are obtained using the GSEA Bioconductor/R package [23]. P-values were adjusted for
22 multiple testing using the Benjamini-Hochberg procedure for controlling FDR [26].

23 Distributions of known EC risk factors were compared between cases and matched
24 controls using independent two-sided sample t-tests, Mann-Whitney U tests, and Chi square
25 tests (R statistical package).

1 **Ethics Statement**

2 The NOWAC study was approved by the Norwegian Data Inspectorate and the Regional
3 Ethical Committee of North Norway (REK). The study was conducted in compliance with the
4 Declaration of Helsinki and all participants gave written informed consent. The linkages of
5 the NOWAC database to national registries such as the Cancer Registry of Norway and
6 registries on death and emigration was approved by the Directorate of Health. The women
7 were informed about these linkages. Furthermore, the collection and storing of human
8 biological material was approved by the REK in accordance with the Norwegian Biobank
9 Act. Women were informed in the letter of introduction that the blood samples would be used
10 for gene expression analyses.

11

12 **Results**

13 **Study population**

14 In total, blood gene expression profiles were analyzed from 79 women diagnosed with EC
15 after blood collection and 79 women, matched by year of birth and time of blood sampling,
16 who did not receive any cancer diagnosis within the same interval after blood collection.
17 There was no significant difference in age, menarche onset, and number of children between
18 the two groups (Table 1). In both group, about half of the women were premenopausal. In
19 general, cases had later occurrence of menopause compared to controls. Additionally, we
20 observed a difference in LNYM with trend towards increase in cases. In our study, cumulative
21 breastfeeding duration was the lowest in controls. This can be explained by a large spread in
22 reported time of breastfeeding by women with EC. There were more OC users among cases.
23 Our study population was overweight, particularly, 64.6% of women with EC diagnosis had
24 BMI >25. Controls drank slightly more coffee.

1 **Differential blood gene expression profiles associated with EC diagnosis and risk factors** 2 **in cases and controls**

3 After preprocessing, the study dataset included expression values for 7104 genes. In overall
4 analysis we were unable to identify any significant differences in gene expression profiles
5 between cases and controls. The same negative result was obtained when we performed
6 separate analysis for each year before diagnosis.

7 Then, we tested the hypothesis that the expression of some genes in the blood of either
8 cases or controls might be influenced by a set of variables modulating EC risk. Among the
9 cases, there was no relationship between any of the variables used and log gene expression. In
10 controls, we observed no differentially expressed genes when using a model that included
11 either of the following: coffee consumption, age of menopause, use of OC. Variations in BMI
12 and number of pregnancies had the strongest impact on the gene expression in controls.
13 Increasing parity was related to expression differences of 1379 genes (FDR 20%). Higher
14 BMI altered the expression of 8, 17, and 27 genes at FDR 10%, 15%, and 20% respectively.
15 Of note, for both BMI and parity, the major number of top 10 genes were downregulated
16 (Table 2).

17 **Gene set enrichment analysis**

18 For GSEA, we used all collections available at MSigDB (Molecular Signatures Database,
19 <http://software.broadinstitute.org/gsea/index.jsp>) [27]. In women diagnosed with EC, totally,
20 we identified 3 significantly enriched gene sets (FDR 20%), of which 1 was parity-associated;
21 1 was enriched along with BMI increase, and 1 was linked to the use of OC (Table S1).
22 Among the controls, GSEA revealed 2415 enriched gene sets (FDR 20%), where 2407 were
23 attributed to parity (Figure 2). Remarkably, the biggest part of these gene sets (786 at FDR
24 15% and 1184 at FDR 20%) were from C7 collection (immunologic gene sets).

1 **Discussion**

2 In a large prospective cohort, we studied the possibility to trace blood gene expression
3 changes prior EC diagnosis. We did not observe any significant differences in expression
4 signatures between cancer-free controls and women with EC when compared directly.
5 Interestingly, BMI and parity, also known to be associated with EC incidence, were
6 associated with significant changes in blood expression profiles in controls but not in cases.
7 To our knowledge, this is the first study demonstrating pre-diagnostic blood gene expression
8 differences between EC cases and matched controls utilizing systems epidemiology approach.

9 Association between high BMI and increased risk of EC development is well
10 established [28]. Moreover, overweight patients with EC have higher risk of death with
11 relative risk up to 6.25 in patients with BMI>40 compared to normal weight women [29]. The
12 main obesity-associated pathways contributing to EC development include augmented
13 estrogen and estrogen metabolites synthesis, presence of chronic inflammation, and insulin
14 resistance [28]. In our single gene analysis, unexpectedly, we did not observed changes in
15 expression associated with BMI increase in EC cases. In turn, we observed a number of genes
16 demonstrating differential expression in controls. This finding might be explained by the fact
17 that the difference in BMI between cases and controls in the study cohort was modest (Table
18 1). In GSEA, we identified 3 significantly enriched gene sets from C2 collection (1 in cases
19 and 2 in controls), and 1 gene set from H collection. Notably,
20 “HALLMARK_TGF_BETA_SIGNALING” gene set from H collection was enriched when
21 the BMI of controls increased. This finding is in line with other studies demonstrating the
22 involvement of disturbed TGF β signaling in EC development and progression [30, 31].

23 Not surprisingly, the major disparity on both single gene and gene sets levels in our
24 study was connected to the number of pregnancies. There is a large body of evidence showing
25 the negative correlation between parity and risk of EC [32-34]. However, recent meta-analysis

1 report nonlinear association between number of children and RR [5]. Indeed, in the entire
2 NOWAC cohort, we found decrease in EC incidence rate in women with 1, 2 or 3 children
3 compared to nulliparous (Figure S1). The elevated incidence rate of EC in women with 4 and
4 more children is attributed to the low number of women with high number of pregnancies in
5 the cohort and, therefore, limited sample size. Similar observations were reported by other
6 studies [33]. Reduced time of estrogen exposure with increased parity is considered a major
7 protective mechanism [35]. Additionally, shedding of the endometrium resulting in
8 mechanical elimination of potentially premalignant cells is well described [5]. In the current
9 study, we observed significant enrichment of a large number of immunologic gene sets (C7
10 collection in MSigDB) in controls with growing number of parities. Based on this finding, it
11 is possible to assume that changes in the immune system associated with pregnancy may be
12 yet another explanation of parity-dependent protection against EC. Moreover, this protective
13 effect grows cumulatively with every new child. Previously, our group published similar
14 observations on parity and BC protection in a larger cohort (article in press). Nevertheless,
15 taking into account the complexity of gene sets data, limited sample size, and per se
16 explorative design of this study, it is impossible to provide clear hypothesis on how the
17 immune system changes in pregnancy contribute to EC protection. Therefore, further studies
18 using both laboratory and epidemiologic design, which address the long-term effects of
19 immune processes on endometrial tumorigenesis, are warranted.

20 In cases, only “HALLMARK_DNA_REPAIR” gene set was significantly associated
21 with high parity (FDR 10%). Alterations in DNA repair machinery are well known to play a
22 major role in EC carcinogenesis [36]. Thus, it is possible to hypothesize that mutations
23 influencing DNA repair mechanisms could abrogate protective effect of parity on EC and lead
24 to cancer development even in women who have given several births.

1 In current work, other factors with published evidences of involvement in EC
2 development had relatively weak impact on pre-diagnostic blood gene expression.

3 It has been demonstrated in NOWAC and by others that increased coffee consumption
4 inversely associated with EC risk [9, 17]. Here we identified one significant gene set
5 “HALLMARK_IL2_STAT5_SIGNALING” (FDR 20%) related to coffee drinking in
6 controls. Recently, Gotthardt and co-authors demonstrated that maintenance of stable STAT5
7 level is necessary for tumor surveillance by NK cells [37]. STAT5 depleted NKs were shown
8 to promote tumor development. Hence, impact of coffee compounds on STAT5 metabolism
9 in immune cells can be an additional biologic substrate of protective functions.

10 In OC users among cases, we revealed significant enrichment of
11 “REACTOME_HYALURONAN_METABOLISM” gene set (FDR 20%). Interestingly,
12 despite the low significance level, second gene set from the top
13 “REACTOME_HYALURONAN_UPTAKE_AND_DEGRADATION” was also related to
14 hyaluronic acid metabolism. Elevated levels of hyaluronic acid in both tumor tissue and
15 serum have been demonstrated to be involved in EC progression [38, 39]. Therefore,
16 monitoring of hyaluronan in the blood of women using OC might be a valuable tool in
17 endometrial cancer screening.

18 To what extent circulating blood cells can reflect processes that occur in tumors is still
19 an open question. Being an easily accessible tissue, blood may serve as an ideal tool for
20 disease prognosis, monitoring and assessment of the treatment. In this work, we attempted to
21 discover gene expression changes in circulating cells that can be identified long before
22 diagnosis of EC. It is important to emphasize that findings reported here need further
23 investigation as most of the information available on role of different genes in tumorigenesis
24 is based on tissue studies and, therefore, cannot be entirely extrapolated to the blood cells.

1 The main strengths of the study include prospective design, population
2 representativeness of the cohort, complete information on cancer status, emigration and
3 mortality obtained from national registries. The systems approach of testing epidemiologic
4 data in the sub cohort using gene expression profiles reduces the probability of false positive
5 findings.

6 Relatively small sample size and the lack of a common algorithm for the gene
7 expression analysis are limitations of this work. In addition, FDR levels we accepted were
8 higher than recommended but this can be justified by relatively small sample size and low p-
9 values of genes and gene sets included.

10 In conclusion, we identified a number of differences in gene set enrichment profiles
11 between cancer-free women and women with EC prior diagnosis in relation to known risk
12 factors for EC. We believe that this integrated analysis may provide a promising background
13 for developing a new multilevel prediction model of EC risk at population level. However,
14 this approach should be further tested in a bigger sample size and in different populations.

15

16 **References:**

- 17 [1] Ferlay J, Soerjomataram I, Dikshit R, Eser S, Mathers C, Rebelo M, Parkin DM, Forman D and Bray F.
18 Cancer incidence and mortality worldwide: sources, methods and major patterns in GLOBOCAN 2012.
19 Int J Cancer 2015; 136: E359-386.
- 20 [2] Cancer in Norway 2016. Cancer Registry of Norway.
21 <https://www.kreftregisteret.no/globalassets/cancer-in-norway/2016/cin-2106-070218.pdf>.
- 22 [3] Torre LA, Siegel RL, Ward EM and Jemal A. Global Cancer Incidence and Mortality Rates and Trends-
23 -An Update. Cancer Epidemiol Biomarkers Prev 2016; 25: 16-27.
- 24 [4] Lortet-Tieulent J, Ferlay J, Bray F and Jemal A. International Patterns and Trends in Endometrial
25 Cancer Incidence, 1978-2013. J Natl Cancer Inst 2017;
- 26 [5] Doepke M. Accounting for fertility decline during the transition to growth. Journal of Economic
27 Growth 2004; 9: 347-383.
- 28 [6] O'Mara TA, Zhao M and Spurdle AB. Meta-analysis of gene expression studies in endometrial cancer
29 identifies gene expression profiles associated with aggressive disease and patient outcome. Sci Rep
30 2016; 6: 36677.
- 31 [7] Saghir FS, Rose IM, Dali AZ, Shamsuddin Z, Jamal AR and Mokhtar NM. Gene expression profiling
32 and cancer-related pathways in type I endometrial carcinoma. Int J Gynecol Cancer 2010; 20: 724-731.
- 33 [8] Sung CO and Sohn I. The expression pattern of 19 genes predicts the histology of endometrial
34 carcinoma. Sci Rep 2014; 4: 5174.

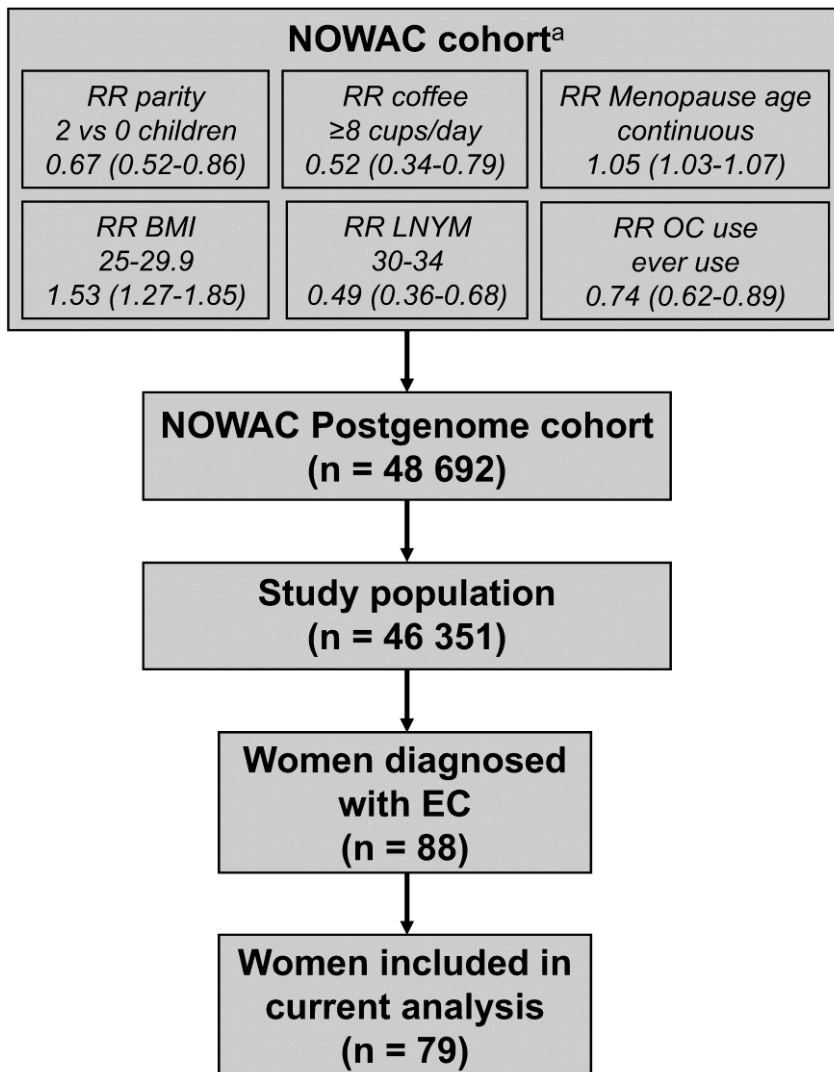
- 1 [9] Dumeaux V, Fjukstad B, Fjosne HE, Frantzen JO, Holmen MM, Rodegerdts E, Schlichting E,
2 Borresen-Dale AL, Bongo LA, Lund E and Hallett M. Interactions between the tumor and the blood
3 systemic response of breast cancer patients. *PLoS Comput Biol* 2017; 13: e1005680.
- 4 [10] Dumeaux V, Olsen KS, Nuel G, Paulssen RH, Borresen-Dale AL and Lund E. Deciphering normal
5 blood gene expression variation--The NOWAC postgenome study. *PLoS Genet* 2010; 6: e1000873.
- 6 [11] Breast cancer paper
- 7 [12] Lund E and Dumeaux V. Systems epidemiology in cancer. *Cancer Epidemiol Biomarkers Prev* 2008;
8 17: 2954-2957.
- 9 [13] Lund E, Holden L, Bovelstad H, Plancade S, Mode N, Gunther CC, Nuel G, Thalabard JC and Holden
10 M. A new statistical method for curve group analysis of longitudinal gene expression data illustrated for
11 breast cancer in the NOWAC postgenome cohort as a proof of principle. *BMC Med Res Methodol*
12 2016; 16: 28.
- 13 [14] Lund E, Dumeaux V, Braaten T, Hjartaker A, Engeset D, Skeie G and Kumle M. Cohort profile: The
14 Norwegian Women and Cancer Study--NOWAC--Kvinner og kreft. *Int J Epidemiol* 2008; 37: 36-41.
- 15 [15] Lund E, Kumle M, Braaten T, Hjartaker A, Bakken K, Eggen E and Gram TI. External validity in a
16 population-based national prospective study--the Norwegian Women and Cancer Study (NOWAC).
17 *Cancer Causes Control* 2003; 14: 1001-1008.
- 18 [16] Gavriluyk O, Braaten T, Weiderpass E, Licaj I and Lund E. Lifetime number of years of menstruation
19 as a risk index for postmenopausal endometrial cancer in the Norwegian Women and Cancer Study.
20 *Acta Obstet Gynecol Scand* 2018.
- 21 [17] Bovelstad H, Holsbo E, Bongo L and Lund E. A Standard Operating Procedure For Outlier Removal In
22 Large-Sample Epidemiological Transcriptomics Datasets. *bioRxiv* 2017;
23 <https://doi.org/10.1101/144519>.
- 24 [18] Shi W, Oshlack A and Smyth GK. Optimizing the noise versus bias trade-off for Illumina whole
25 genome expression BeadChips. *Nucleic Acids Res* 2010; 38: e204.
- 26 [19] Lin SM, Du P, Huber W and Kibbe WA. Model-based variance-stabilizing transformation for Illumina
27 microarray data. *Nucleic Acids Res* 2008; 36: e11.
- 28 [20] Du P, Feng G, Kibbe W and Lin S. lumiHumanIDMapping: Illumina Identifier mapping for Human. R
29 package version 1.10.1. 2016.
- 30 [21] Smyth GK. Linear Models and Empirical Bayes Methods for Assessing Differential Expression in
31 Microarray Experiments. *Statistical Applications in Genetics and Molecular Biology* 2004; 3: 1-25.
- 32 [22] Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A, Pomeroy
33 SL, Golub TR, Lander ES and Mesirov JP. Gene set enrichment analysis: a knowledge-based approach
34 for interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A* 2005; 102: 15545-15550.
- 35 [23] Hänzelmann S, Castelo R and Guinney J. GSEA: gene set variation analysis for microarray and RNA-
36 Seq data. *BMC Bioinformatics* 2013; 14: 7.
- 37 [24] Reiner A, Yekutieli D and Benjamini Y. Identifying differentially expressed genes using false discovery
38 rate controlling procedures. *Bioinformatics* 2003; 19: 368-375.
- 39 [28] Onstad MA, Schmandt RE and Lu KH. Addressing the Role of Obesity in Endometrial Cancer Risk,
40 Prevention, and Treatment. *J Clin Oncol* 2016; 34: 4225-4230.
- 41 [29] Calle EE, Rodriguez C, Walker-Thurmond K and Thun MJ. Overweight, obesity, and mortality from
42 cancer in a prospectively studied cohort of U.S. adults. *N Engl J Med* 2003; 348: 1625-1638.
- 43 [30] Li Q. Transforming growth factor beta signaling in uterine development and function. *J Anim Sci*
44 *Biotechnol* 2014; 5: 52.
- 45 [31] Piestrzeniewicz-Ulanska D, Brys M, Semczuk A, Rechberger T, Jakowicki JA and Krajewska WM.
46 TGF-beta signaling is disrupted in endometrioid-type endometrial carcinomas. *Gynecol Oncol* 2004;
47 95: 173-180.
- 48 [32] Kvale G, Heuch I and Ursin G. Reproductive factors and risk of cancer of the uterine corpus: a
49 prospective study. *Cancer Res* 1988; 48: 6217-6221.
- 50 [33] Dossus L, Allen N, Kaaks R, Bakken K, Lund E, Tjonneland A, Olsen A, Overvad K, Clavel-Chapelon
51 F, Fournier A, Chabbert-Buffet N, Boeing H, Schutze M, Trichopoulou A, Trichopoulos D, Lagiou P,
52 Palli D, Krogh V, Tumino R, Vineis P, Mattiello A, Bueno-de-Mesquita HB, Onland-Moret NC,
53 Peeters PH, Dumeaux V, Redondo ML, Duell E, Sanchez-Cantalejo E, Arriola L, Chirlaque MD,
54 Ardanaz E, Manjer J, Borgquist S, Lukanova A, Lundin E, Khaw KT, Wareham N, Key T, Chajes V,
55 Rinaldi S, Slimani N, Mouw T, Gallo V and Riboli E. Reproductive risk factors and endometrial
56 cancer: the European Prospective Investigation into Cancer and Nutrition. *Int J Cancer* 2010; 127: 442-
57 451.
- 58 [34] Schonfeld SJ, Hartge P, Pfeiffer RM, Freedman DM, Greenlee RT, Linet MS, Park Y, Schairer C,
59 Visvanathan K and Lacey JV, Jr. An aggregated analysis of hormonal factors and endometrial cancer
60 risk by parity. *Cancer* 2013; 119: 1393-1401.

- 1 [35] Chen Q, Tong M, Guo F, Lau S and Zhao M. Parity Correlates with the Timing of Developing
2 Endometrial Cancer, But Not Subtype of Endometrial Cancer. *J Cancer* 2015; 6: 1087-1092.
- 3 [36] Masuda K, Banno K, Yanokura M, Kobayashi Y, Kisu I, Ueki A, Ono A, Asahara N, Nomura H,
4 Hirasawa A, Susumu N and Aoki D. Relationship between DNA Mismatch Repair Deficiency and
5 Endometrial Cancer. *Mol Biol Int* 2011; 2011: 256063.
- 6 [37] Gotthardt D, Putz EM, Grundschober E, Prchal-Murphy M, Straka E, Kudweis P, Heller G, Bago-
7 Horvath Z, Witalisz-Siepracka A, Cumaraswamy AA, Gunning PT, Strobl B, Muller M, Moriggl R,
8 Stockmann C and Sexl V. STAT5 Is a Key Regulator in NK Cells and Acts as a Molecular Switch from
9 Tumor Surveillance to Tumor Promotion. *Cancer Discov* 2016; 6: 414-429.
- 10 [38] Paiva P, Van Damme MP, Tellbach M, Jones RL, Jobling T and Salamonsen LA. Expression patterns
11 of hyaluronan, hyaluronan synthases and hyaluronidases indicate a role for hyaluronan in the
12 progression of endometrial cancer. *Gynecol Oncol* 2005; 98: 193-202.
- 13 [39] Yabushita H, Kishida T, Fusano K, Kanyama K, Zhuo L, Itano N, Kimata K and Noguchi M. Role of
14 hyaluronan and hyaluronan synthase in endometrial cancer. *Oncol Rep* 2005; 13: 1101-1105.

15

16

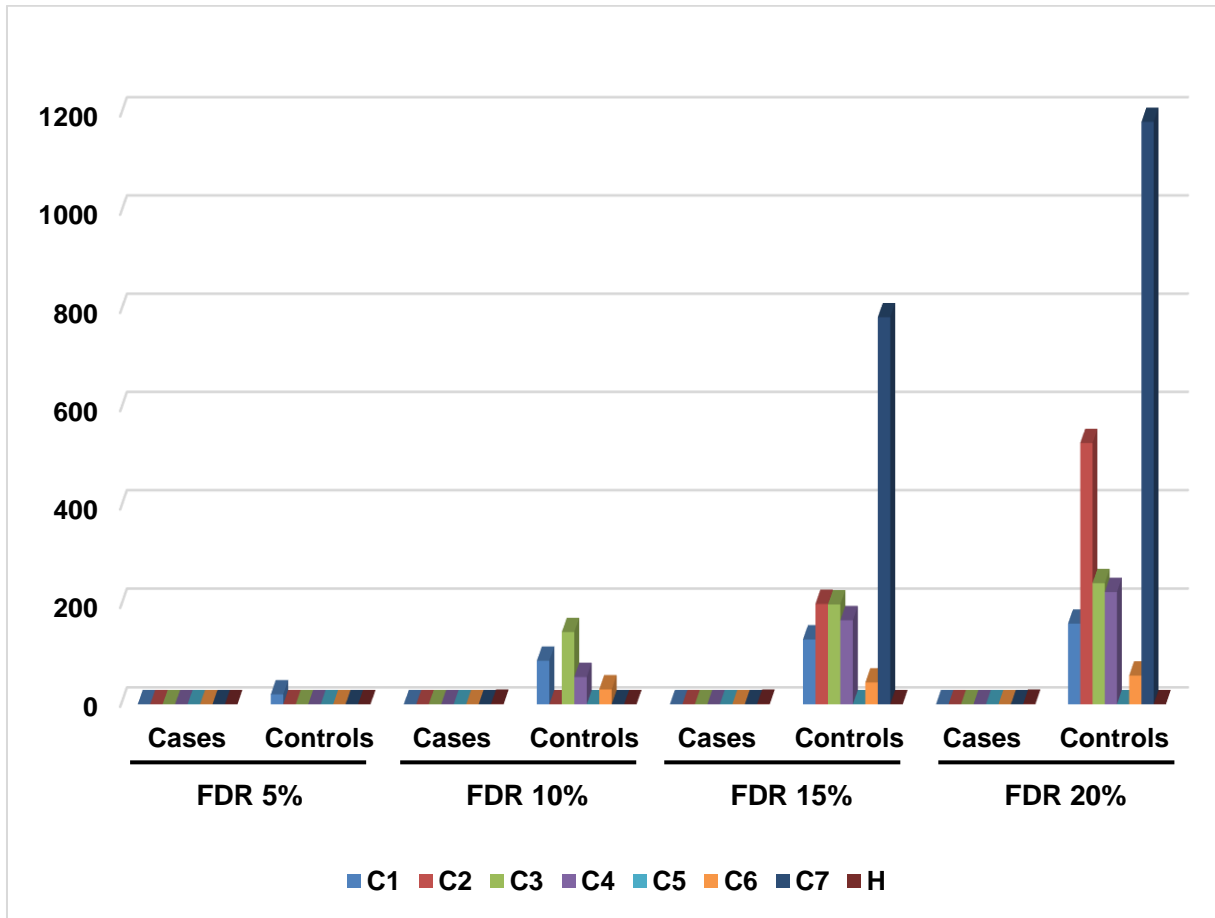
1 **Figure 1. Study population.**



2

3 **Notes:** a. Relative risk estimates of EC in the Norwegian Women and Cancer Study calculated
 4 using multivariable model adjusted for BMI, use of HRT, use of OC, smoking and alcohol
 5 consumption

1 **Figure 2. Number of significantly enriched gene sets in association with increasing**
 2 **parity.**



3

4

1 **Table 1. Baseline characteristics of the study (N = 79 case/control pairs)**

Characteristics	EC cases	Controls
Age (mean, \pm SE)	49.3 (6.3)	49.6 (6.5)
Menopausal status		
<i>Premenopausal</i>	38 (56 %)	41 (46 %)
<i>Postmenopausal</i>	30 (44 %)	49 (54 %)
Age at menopause (mean, \pm SE)	49.9 (4.2)	47.5 (5.5)
Age at menarche (mean, \pm SE)	12.9 (1.4)	13.2 (1.3)
Parity		
0	8	8
1	11	9
2	38	35
3	17	17
4	5	9
5	0	1
Cumulative duration of breastfeeding (mean, \pm SE)	13.5 (17.0)	10.4 (7.2)
Ever consumption of oral contraceptives (%)	43%	30.7%
LNYM		
<25	10.7%	15.4%
25-29	10.7%	14.1%
30-34	36.0%	41.0%
35-39	34.7%	26.9%
40+	8.0%	2.6%
Body mass index (mean, \pm SE)	27.5(5.4)	25.8 (4.9)
< 25 (%)	35.4%	58.2%
> 25 (%)	64.6%	41.8%
Ever coffee consumption (%)	92.4%	93.7%

2

1 **Table 2. Top 10 differentially expressed genes associated with parity and BMI.**

Parity				
Gene	logFC^a	p-value	q-value^b	Function
NUDT22	-0.081	8.15E-05	0.163	UDP-glucose and UDP-galactose hydrolase
SH2B2	-0.104	8.34E-05	0.163	Regulator of tyrosine kinase receptor activity
NUP188	-0.056	0.000123829	0.163	Nuclear pore complex involved in the flow of various substances between the cytoplasm and nucleoplasm
TRIP12	0.134	0.000157622	0.163	E3 ubiquitin-protein ligase, involved in regulation of DNA repair
APBA3	-0.094	0.000204851	0.163	Involved in signal transduction and synaptic transmission
CYB5A	0.055	0.000221835	0.163	Electron carrier, regulates hemoglobin metabolism
CEP250	-0.053	0.000252363	0.163	Required for interphase progression of the cell cycle
NRM	-0.073	0.000258652	0.163	Encodes protein residing within the inner nuclear membrane. May be involved in apoptosis
PLRG1	0.097	0.000348695	0.163	Regulator of alternative splicing
TWF2	-0.083	0.000479402	0.163	involved in motile processes and endocytosis regulation
BMI				
Gene	logFC^a	p-value	q-value^b	Function
ALS2	0.016	2.69E-05	0.099	GTPase regulator. Involved in the development of spinal neurons
TAOK1	0.026	3.11E-05	0.099	Involved in p38 MAPK signaling, apoptosis regulation and cytoskeleton maintenance
ZZEF1	-0.019	7.34E-05	0.099	Involved in calcium ion binding.
DNAJB1	-0.020	7.35E-05	0.099	Stimulates ATP hydrolysis and promotes folding and unfolding of proteins
PROSC	-0.017	7.87E-05	0.099	Involved in homeostasis regulation of pyridoxal 5-phosphate (active form of B6 vitamin)
H2AFY	-0.023	8.35E-05	0.099	Histone-coding gene that represses transcription and inactivates X chromosome

EDEM1	-0.021	0.000109136	0.099	Involved in protein processic in endoplasmic reticulum.
ZBTB44	0.024	0.000111261	0.099	Zinc finger protein realted to nucleic acid binding
SFRS9	-0.019	0.000132804	0.105	Regulates mRNA maturation.
ANKRD11	-0.020	0.000156181	0.110	Inhibits ligand-dependent activation of transcription

Notes: ^aLogFC is the estimated log-fold change in gene expression when the parity increases continuously. ^bq-value is an FDR adjusted p-value.

1 **Supplementary Table 1. Gene set enrichment analysis.**

Number of significant gene sets		
(FDR 5% – 10% – 15% – 20%)		
Parity	Cases	Controls
C1	0 - 0 - 0 - 0	20 - 88 - 131 - 163
C2	0 - 0 - 0 - 0	0 - 0 - 203 - 530
C3	0 - 0 - 0 - 0	0 - 146 - 202 - 245
C4	0 - 0 - 0 - 0	0 - 55 - 170 - 227
C5	0 - 0 - 0 - 0	0 - 0 - 0 - 0
C6	0 - 0 - 0 - 0	0 - 30 - 44 - 58
C7	0 - 0 - 0 - 0	0 - 0 - 786 - 1184
H	0 - 1 - 1 - 1	0 - 0 - 0 - 0
Coffee		
C1	0 - 0 - 0 - 0	0 - 0 - 0 - 0
C2	0 - 0 - 0 - 0	0 - 0 - 0 - 0
C3	0 - 0 - 0 - 0	0 - 0 - 0 - 0
C4	0 - 0 - 0 - 0	0 - 0 - 0 - 0
C5	0 - 0 - 0 - 0	0 - 0 - 0 - 0
C6	0 - 0 - 0 - 0	0 - 0 - 0 - 0
C7	0 - 0 - 0 - 0	0 - 0 - 0 - 0
H	0 - 0 - 0 - 0	0 - 0 - 0 - 1
BMI		
C1	0 - 0 - 0 - 0	0 - 0 - 0 - 0
C2	1 - 1 - 1 - 1	0 - 0 - 2 - 2
C3	0 - 0 - 0 - 0	0 - 0 - 0 - 0
C4	0 - 0 - 0 - 0	0 - 0 - 0 - 0
C5	0 - 0 - 0 - 0	0 - 0 - 0 - 0
C6	0 - 0 - 0 - 0	0 - 0 - 0 - 0
C7	0 - 0 - 0 - 0	0 - 0 - 0 - 0
H	0 - 0 - 0 - 0	1 - 1 - 1 - 1
LN YM		

C1	0-0-0-0	0-1-1-2
C2	0-0-0-0	0-0-0-0
C3	0-0-0-0	0-0-0-0
C4	0-0-0-0	0-0-0-0
C5	0-0-0-0	0-0-0-0
C6	0-0-0-0	0-0-0-0
C7	0-0-0-0	0-0-0-0
H	0-0-0-0	0-0-1-1

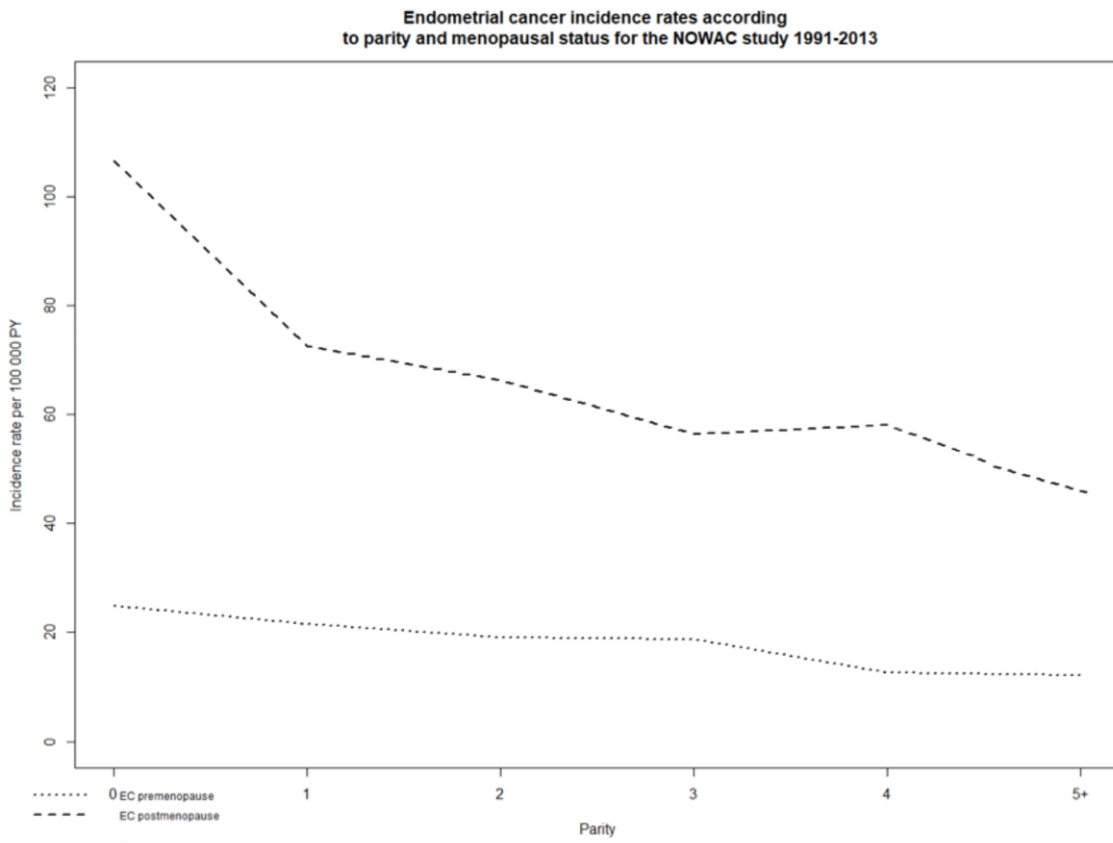
Age of menopause

C1	0-0-0-0	0-0-0-0
C2	0-0-0-0	0-0-0-0
C3	0-0-0-0	0-0-0-0
C4	0-0-0-0	0-0-0-0
C5	0-0-0-0	0-0-0-0
C6	0-0-0-0	0-0-0-0
C7	0-0-0-0	0-0-0-0
H	0-0-0-0	0-0-0-0

OC use

C1	0-0-0-0	0-0-0-0
C2	0-0-0-1	0-0-0-0
C3	0-0-0-0	0-0-0-0
C4	0-0-0-0	1-1-1-1
C5	0-0-0-0	0-0-1-1
C6	0-0-0-0	0-0-0-0
C7	0-0-0-0	0-0-0-0
H	0-0-0-0	0-0-0-0

1 **Supplementary Figure 1. EC risk according to parity**



2

3