

## Five endometrial cancer risk loci identified through genome-wide association analysis

Timothy HT Cheng<sup>1\*</sup>, Deborah J Thompson<sup>2\*</sup>, Tracy A O'Mara<sup>3</sup>, Jodie N Painter<sup>3</sup>, Dylan M Glubb<sup>3</sup>, Susanne Flach<sup>1</sup>, Anabelle Lewis<sup>1</sup>, Juliet D French<sup>3</sup>, Luke Freeman-Mills<sup>1</sup>, David Church<sup>1</sup>, Maggie Gorman<sup>1</sup>, Lynn Martin<sup>1</sup>, National Study of Endometrial Cancer Genetics Group (NSECg)<sup>1</sup>, Shirley Hodgson<sup>4</sup>, Penelope M Webb<sup>3</sup>, The Australian National Endometrial Cancer Study Group (ANECs)<sup>3</sup>, John Attia<sup>5,6</sup>, Elizabeth G Holliday<sup>5,6</sup>, Mark McEvoy<sup>6</sup>, Rodney J Scott<sup>5,7-9</sup>, Anjali K Henders<sup>3</sup>, Nicholas G Martin<sup>3</sup>, Grant W Montgomery<sup>3</sup>, Dale R Nyholt<sup>3,10</sup>, Shahana Ahmed<sup>11</sup>, Catherine S Healey<sup>11</sup>, Mitul Shah<sup>11</sup>, Joe Dennis<sup>2</sup>, Peter A Fasching<sup>12,13</sup>, Matthias W Beckmann<sup>13</sup>, Alexander Hein<sup>13</sup>, Arif B Ekici<sup>14</sup>, Per Hall<sup>15</sup>, Kamila Czene<sup>15</sup>, Hatef Darabi<sup>15</sup>, Jingmei Li<sup>15</sup>, Thilo Dörk<sup>16</sup>, Matthias Dürst<sup>17</sup>, Peter Hillemanns<sup>18</sup>, Ingo Runnebaum<sup>17</sup>, Frederic Amant<sup>19</sup>, Stefanie Schrauwen<sup>19</sup>, Hui Zhao<sup>20,21</sup>, Diether Lambrechts<sup>20,21</sup>, Jeroen Depreeuw<sup>19-21</sup>, Sean C Dowdy<sup>22</sup>, Ellen L Goode<sup>23</sup>, Brooke L Fridley<sup>24</sup>, Stacey J Winham<sup>23</sup>, Tormund S Njølstad<sup>25,26</sup>, Helga B Salvesen<sup>25,26</sup>, Jone Trovik<sup>25,26</sup>, Henrica MJ Werner<sup>25,26</sup>, Katie Ashton<sup>5,8,9</sup>, Geoffrey Otton<sup>27</sup>, Tony Proietto<sup>27</sup>, Tao Liu<sup>28</sup>, Miriam Mints<sup>29</sup>, Emma Tham<sup>28</sup>, RENDOCAS<sup>28</sup>, CHIBCHA Consortium<sup>1</sup>, Mulin Jun Li<sup>30</sup>, Shun Yip<sup>30</sup>, Junwen Wang<sup>30</sup>, Manjeet K Bolla<sup>2</sup>, Kyriaki Michailidou<sup>2</sup>, Qin Wang<sup>2</sup>, Jonathan P Tyrer<sup>11</sup>, Malcolm Dunlop<sup>31,32</sup>, Richard Houlston<sup>33</sup>, Claire Palles<sup>1</sup>, John L Hopper<sup>34</sup>, AOCs Group<sup>3,35</sup>, Julian Peto<sup>36</sup>, Anthony J Swerdlow<sup>33,37</sup>, Barbara Burwinkel<sup>38,39</sup>, Hermann Brenner<sup>40-42</sup>, Alfons Meindl<sup>43</sup>, Hiltrud Brauch<sup>42,44,45</sup>, Annika Lindblom<sup>28</sup>, Jenny Chang-Claude<sup>46</sup>, Fergus J Couch<sup>23,47</sup>, Graham G Giles<sup>34,48,49</sup>, Vessela N Kristensen<sup>50-52</sup>, Angela Cox<sup>53</sup>, Julie M Cunningham<sup>23,47</sup>, Paul D P Pharoah<sup>11</sup>, Alison M Dunning<sup>11</sup>, Stacey L Edwards<sup>3</sup>, Douglas F Easton<sup>2,11+</sup>, Ian Tomlinson<sup>1+</sup>, Amanda B Spurdle<sup>3+</sup>

\* contributed equally to this work

+ to whom correspondence should be addressed

<sup>1</sup> Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford, OX3 7BN, UK.

<sup>2</sup> Centre for Cancer Genetic Epidemiology, Department of Public Health and Primary Care, University of Cambridge, Cambridge, CB1 8RN, UK.

<sup>3</sup> Department of Genetics and Computational Biology, QIMR Berghofer Medical Research Institute, Brisbane, QLD, 4006, Australia.

<sup>4</sup> Department of Clinical Genetics, St George's, University of London, London, SW17 0RE, UK.

<sup>5</sup> Hunter Medical Research Institute, John Hunter Hospital, Newcastle, NSW, 2305, Australia.

<sup>6</sup> Centre for Clinical Epidemiology and Biostatistics, School of Medicine and Public Health, University of Newcastle, NSW, 2305, Australia.

<sup>7</sup> Hunter Area Pathology Service, John Hunter Hospital, Newcastle, NSW, 2305, Australia.

<sup>8</sup> Centre for Information Based Medicine, University of Newcastle, NSW, 2308, Australia.

<sup>9</sup> School of Biomedical Sciences and Pharmacy, University of Newcastle, Newcastle, NSW, 2308, Australia.

<sup>10</sup> Institute of Health and Biomedical Innovation, Queensland University of Technology, Brisbane, 4006, Australia.

<sup>11</sup> Centre for Cancer Genetic Epidemiology, Department of Oncology, University of Cambridge, Cambridge, CB1 8RN, UK.

<sup>12</sup> University of California at Los Angeles, Department of Medicine, Division of Hematology/Oncology, David Geffen School of Medicine, Los Angeles, CA, 90095, USA.

<sup>13</sup> Department of Gynecology and Obstetrics, University Hospital Erlangen, Friedrich-Alexander University Erlangen-Nuremberg, Erlangen, 91054, Germany.

<sup>14</sup> Institute of Human Genetics, University Hospital Erlangen, Friedrich-Alexander-University Erlangen-Nuremberg, Erlangen, 91054, Germany.

- <sup>15</sup> Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, SE-171 77, Sweden.
- <sup>16</sup> Hannover Medical School, Gynaecology Research Unit, Hannover, 30625, Germany.
- <sup>17</sup> Department of Gynaecology, Jena University Hospital - Friedrich Schiller University, Jena, 07743, Germany.
- <sup>18</sup> Hannover Medical School, Clinics of Gynaecology and Obstetrics, Hannover, 30625, Germany.
- <sup>19</sup> Department of Obstetrics and Gynecology, Division of Gynecologic Oncology, University Hospitals, KU Leuven - University of Leuven, 3000, Belgium.
- <sup>20</sup> Vesalius Research Center, Leuven, 3000, Belgium.
- <sup>21</sup> Laboratory for Translational Genetics, Department of Oncology, University Hospitals Leuven, Leuven, 3000, Belgium.
- <sup>22</sup> Department of Obstetrics and Gynecology, Division of Gynecologic Oncology, Mayo Clinic, Rochester, MN, 55905, USA.
- <sup>23</sup> Department of Health Sciences Research, Mayo Clinic, Rochester, MN, 55905, USA.
- <sup>24</sup> Department of Biostatistics, University of Kansas Medical Center, Kansas City, KS, 66160, USA.
- <sup>25</sup> Centre for Cancerbiomarkers, Department of Clinical Science, The University of Bergen, 5020, Norway.
- <sup>26</sup> Department of Obstetrics and Gynecology, Haukeland University Hospital, Bergen, 5021, Norway.
- <sup>27</sup> School of Medicine and Public Health, University of Newcastle, Newcastle, NSW, 2308, Australia.
- <sup>28</sup> Department of Molecular Medicine and Surgery, Karolinska Institutet, Stockholm, SE-171 77, Sweden.
- <sup>29</sup> Department of Women's and Children's Health, Karolinska Institutet, Karolinska University Hospital, Stockholm, SE-171 77, Sweden.
- <sup>30</sup> Centre for Genomic Sciences, Li Ka Shing Faculty of Medicine, The University of Hong Kong, Hong Kong SAR, China.
- <sup>31</sup> Colon Cancer Genetics Group, Institute of Genetics and Molecular Medicine, University of Edinburgh, Edinburgh, EH4 2XU, UK.
- <sup>32</sup> MRC Human Genetics Unit, Western General Hospital Edinburgh, Edinburgh, EH4 2XU, UK.
- <sup>33</sup> Division of Genetics and Epidemiology, Institute of Cancer Research, London, SM2 5NG, UK.
- <sup>34</sup> Centre for Epidemiology and Biostatistics, Melbourne School of Population and Global Health, The University of Melbourne, Vic, 3010, Australia.
- <sup>35</sup> Peter MacCallum Cancer Center, The University of Melbourne, Melbourne, 3002, Australia.
- <sup>36</sup> London School of Hygiene and Tropical Medicine, London, WC1E 7HT, UK.
- <sup>37</sup> Division of Breast Cancer Research, Institute of Cancer Research, London, SM2 5NG, UK.
- <sup>38</sup> Molecular Biology of Breast Cancer, Department of Gynecology and Obstetrics, University of Heidelberg, Heidelberg, 69120, Germany.
- <sup>39</sup> Molecular Epidemiology Group, German Cancer Research Center, DKFZ, Heidelberg, 69120, Germany.
- <sup>40</sup> Division of Clinical Epidemiology and Aging Research, German Cancer Research Center (DKFZ), Heidelberg, 69120, Germany.
- <sup>41</sup> Division of Preventive Oncology, German Cancer Research Center (DKFZ) and National Center for Tumor Diseases (NCT), Heidelberg, 69120, Germany
- <sup>42</sup> German Cancer Consortium (DKTK), German Cancer Research Center (DKFZ), Heidelberg, 69120, Germany
- <sup>43</sup> Department of Obstetrics and Gynecology, Division of Tumor Genetics, Technical University of Munich, Munich, 80333, Germany.
- <sup>44</sup> Dr. Margarete Fischer-Bosch-Institute of Clinical Pharmacology, Stuttgart, 70376, Germany.
- <sup>45</sup> University of Tübingen, Tübingen, 72074, Germany
- <sup>46</sup> Division of Cancer Epidemiology, German Cancer Research Center (DKFZ), Heidelberg, 69120, Germany.
- <sup>47</sup> Department of Laboratory Medicine and Pathology, Mayo Clinic, Rochester, MN, 55905, USA.

<sup>48</sup> Cancer Epidemiology Centre, Cancer Council Victoria, Melbourne, Vic, 3004, Australia.

<sup>49</sup> Department of Epidemiology and Preventive Medicine, Monash University, Melbourne, Vic, 3004, Australia.

<sup>50</sup> Department of Genetics, Institute for Cancer Research, The Norwegian Radium Hospital, Oslo, 0310, Norway.

<sup>51</sup> The K.G. Jebsen Center for Breast Cancer Research, Institute for Clinical Medicine, Faculty of Medicine, University of Oslo, Oslo, 0316, Norway .

<sup>52</sup> Department of Clinical Molecular Oncology, Division of Medicine, Akershus University Hospital, Lørenskog, 1478, Norway.

<sup>53</sup> Sheffield Cancer Research, Department of Oncology, University of Sheffield, Sheffield, S10 2RX, UK.

**Abbreviations:** EC, endometrial cancer; CI, confidence interval; GWAS, genome-wide association study; LD, linkage disequilibrium; OR, odds ratio; kb, kilobase; Mb, megabase; PCA, principal components analysis; DHS, DNase1 hypersensitivity site.

## Abstract

We conducted a meta-analysis of three endometrial cancer (EC) GWAS and two replication phases totaling 7,737 EC cases and 37,144 controls of European ancestry. Genome-wide imputation and meta-analysis identified five novel risk loci at genome-wide significance at likely regulatory regions on chromosomes 13q22.1 (rs11841589, near *KLF5*), 6q22.31 (rs13328298, in *LOC643623* and near *HEY2* and *NCOA7*), 8q24.21 (rs4733613, telomeric to *MYC*), 15q15.1 (rs937213, in *EIF2AK4*, near *BMF*) and 14q32.33 (rs2498796, in *AKT1* near *SIVA1*). A second independent 8q24.21 signal (rs17232730) was found. Functional studies of the intergenic 13q22.1 locus showed that rs9600103 (pairwise  $r^2=0.98$  with rs11841589) is located in a region of active chromatin that interacts with the *KLF5* promoter region. The rs9600103-T EC protective allele suppressed gene expression *in vitro* suggesting that the regulation of *KLF5* expression, a gene linked to uterine development, is implicated in tumorigenesis. These findings provide enhanced insight into the genetic and biological basis of EC.

Endometrial cancer (EC) is the fourth most common cancer in women in the United States<sup>1</sup> and Europe<sup>2</sup>, and the most common cancer of the female reproductive system. The familial relative risk is  $\sim 2^{3,4}$ , but highly penetrant germline mutations in mismatch repair genes<sup>5</sup>, and DNA polymerases<sup>6,7</sup> account for only a small proportion of the familial aggregation. Our previous GWAS and subsequent fine-mapping identified the only two reported genome-wide significant EC risk loci, tagged by rs11263763 in *HNF1B* intron 1<sup>8</sup> and rs727479 in *CYP19A1* intron 4<sup>9</sup>.

To identify additional EC risk loci, we re-analysed data from our previous GWAS (ANECS, SEARCH datasets<sup>10</sup>) and conducted a meta-analysis with two further studies (**Supplementary Figure 1**). The first was an independent GWAS; the National Study of Endometrial Cancer (NSECG), including 925 EC cases genotyped using the Illumina 660W array, 1,286 cancer-free controls from the CORGI/SP1 GWAS<sup>11,12</sup> and 2,674 controls from the 1958 Birth Cohort<sup>13</sup>. The second study comprised 4,330 EC cases and 26,849 controls from Europe, the United States and Australia, genotyped using a custom array designed by the Collaborative Oncological Gene-environment Study (COGS) initiative<sup>14–17</sup> (**Supplementary Table 1, Supplementary Note**).

We first performed genome-wide imputation using 1000 Genomes Project data, allowing us to assess up to 8.6 million variants with allele frequency  $\geq 1\%$  across the different studies. Per-allele odds ratios and P-values for all SNPs in the GWAS and iCOGS were obtained using a logistic regression model. There was little evidence of systematic overdispersion of the test statistic ( $\lambda_{GC}=1.002-1.038$ , **Supplementary Figure 3**). A fixed-effects meta-analysis was conducted for all 2.3 million typed and well-imputed SNPs (info score  $> 0.90$ ) in a total of 6,542 EC cases and 36,393 controls. The strongest associations were with SNPs in LD with previously identified EC risk SNPs in *HNF1B*<sup>10,8,18</sup> and *CYP19A1*<sup>9,19</sup> (**Figure 1, Table 1**). For fourteen 1.5Mb regions containing at least one novel SNP with  $P_{meta} < 10^{-5}$ , we performed regional imputation using an additional reference panel that comprised 196 high-coverage whole genome-sequenced UK individuals (**Supplementary Table 2**).

Five novel regions containing at least one EC risk SNP with  $P_{meta} < 10^{-7}$  were identified and the most strongly associated SNP in each region was genotyped in an additional 1,195 NSECG EC cases and 751 controls using competitive allele-specific PCR (KASPar, KBiosciences) and the Fluidigm BioMark System (**Supplementary Table 3**). Duplicate samples displayed concordance  $> 98.5\%$  between different genotyping platforms (**Supplementary Table 4**). All five SNPs were associated with EC at genome-wide significance ( $P < 5 \times 10^{-8}$ , **Table 1, Figure 2**), and these associations remained highly

significant when analysis was restricted to cases with endometrioid subtype only. Endometrioid-only analysis did not reveal any additional risk loci. eQTL analysis (**Online Methods**) in normal uterine tissue,<sup>20</sup> and EC tumour and adjacent normal tissue<sup>21</sup> did not yield any SNPs robustly associated with the expression of nearby genes at the EC risk loci (**Supplementary Table 7**). However, for each risk locus, bioinformatic analysis including cell-type-specific expression and histone modification data identified correlated SNPs within 500kb in likely enhancers and multiple potential regulatory targets (**Supplementary Table 6, Supplementary Figure 5**). The most compelling candidates for future functional analysis are described below.

rs13328298 (OR=1.13, 95%CI:1.09–1.18,  $P=3.73\times10^{-10}$ ) on 6q22.31 lies in the long non-coding RNA *LOC643623*, 54kb upstream of *HEY2* and 86kb upstream of *NCOA7*. *HEY2* is a helix-loop-helix transcriptional repressor in the Notch pathway, which maintains stem cells, and dysregulation has been associated with different cancers<sup>22</sup>. *NCOA7* modulates the activity of the estrogen receptor via direct binding<sup>23</sup>.

The second locus (rs4733613, OR=0.84, 95%CI:0.80–0.89,  $P=3.09\times10^{-9}$ ) is at 8q24.21. Stepwise conditional logistic regression identified another independent signal in this region, rs17232730 (pairwise  $r^2=0.02$ ,  $P_{\text{cond}}=1.29\times10^{-5}$ , **Table 2**). Both EC SNPs lie further from *MYC* (784-846kb telomeric) than most of the other cancer SNPs in the region, including those for cancers of the bladder<sup>24,25</sup>, breast<sup>26,16</sup>, colorectum<sup>11,27</sup>, ovary<sup>28</sup> and prostate<sup>29,30</sup>. rs17232730 is in moderate LD with the ovarian cancer SNP rs10088218 ( $r^2=0.43$ ), with both cancers sharing the same risk allele, but rs4733613 is not in LD ( $r^2\leq0.02$ ) with any other cancer SNP in the region (**Supplementary Figure 5**). A role in tumorigenesis is implicated for several miRNAs in the region<sup>31</sup>. Of these, miR-1207-5p is reported to repress *TERT*, a locus also implicated in EC risk<sup>32</sup>.

The lead SNP at 15q15 (rs937213; OR=0.90, 95%CI:0.86–0.93,  $P=1.77\times10^{-8}$ ) lies within an intron of *EIF2AK4*. *EIF2AK4* encodes a kinase that phosphorylates EIF2 $\alpha$  and downregulates protein synthesis during cellular stress<sup>33</sup>. Another nearby gene, *BMF*, encodes an apoptotic regulator moderately to highly expressed in glandular endometrial tissue<sup>34</sup>.

At 14q42, the lead SNP rs2498796 (OR=0.89, 95%CI:0.85–0.93,  $P=3.55\times10^{-8}$ ) lies in intron 3 of oncogene *AKT1*, which is highly expressed in the endometrium<sup>34</sup>. Several SNPs in LD with rs2498796 are bioinformatically linked with regulation of *AKT1* and four other nearby genes (*SIVA1*, *ZBTB42*, *ADSSL1* and *INF2*; **Supplementary Table 6, Supplementary**

**Figure 5).** *AKT1* acts in the PI3K/AKT/MTOR intracellular signaling pathway, which affects cell survival and proliferation<sup>35</sup> and is activated in endometrial tumors<sup>36</sup>, especially aggressive disease<sup>37,38,39</sup>. *SIVA1* encodes an apoptosis regulatory protein that inhibits p53 activity<sup>40,41</sup> and enhances epithelial–mesenchymal transition to promote motility and invasiveness of epithelial cells<sup>42</sup>. *INF2* expression is reported to act as a promigratory signal in gastric cancer cells treated with mycophenolic acid<sup>43</sup>.

The final novel EC SNP was rs11841589 (OR=1.15, 95%CI:1.11–1.21,  $P=4.83\times10^{-11}$ ) on chromosome 13q22.1, 163kb and 445kb downstream from Kruppel-like factors *KLF5* and *KLF12*, respectively. *KLF5* is a transcription factor associated with cell cycle regulation, and it plays a role in uterine development, homoeostasis and tumorigenesis<sup>44–47</sup>. Elevated *KLF5* levels are strongly correlated with activating *KRAS* mutations<sup>48</sup> and *KLF5* is targeted for degradation by the tumor suppressor *FBXW7*. Both *FBXW7* and *KRAS* are commonly mutated in EC<sup>49</sup>. rs11841589 was one of a group of five highly correlated SNPs ( $r^2\geq0.98$ ) surpassing genome wide significance in a 3kb LD block bounded by rs9600103 ( $P=8.70\times10^{-11}$ ) and rs11841589 (**Figure 4a**). There was no residual association signal at this locus ( $P_{\text{cond}} > 0.05$ ) after conditioning for rs11841589. Bioinformatic analysis suggested that the causal variant at the intergenic 13q22.1 locus may affect a regulatory element that modifies *KLF5* expression (**Supplementary Figure 5**); rs9600103 overlaps a vertebrate conservation peak, and a DNaseI hypersensitivity site (DHS) in estrogen and tamoxifen-treated ENCODE<sup>50</sup> Ishikawa cells (**Figure 4a**). In addition, in a Hi-C chromatin capture experiment in Hela S3 cells<sup>51</sup>, an interaction loop was observed between a segment containing the *KLF5* promoter and the rs11841589/rs9600103 locus ( $P=0.004$ , **Supplementary Figure 6**).

We further investigated the epigenetic landscape of a 16kb region around rs11841589 and rs9600103 that contained the SNPs most strongly associated with EC, by analysis of three EC cell lines: Ishikawa is homozygous for the rs9600103-A and rs11841589-G high-risk alleles, and provided a comparison with the ENCODE data; ARK-2 is homozygous for the low-risk T alleles at both SNPs; and AN3CA is a non-*KLF5* expressing line that is homozygous for the high-risk alleles (**Supplementary Figure 7**). We conducted formaldehyde-assisted identification of regulatory elements (FAIRE, to identify regions of open chromatin), and chromatin immunoprecipitation (ChIP) using antibodies against H3K4Me2 (marker of transcription factor binding<sup>52</sup>) and panH4Ac (marker of active chromatin). Although the anti-H4Ac ChIP did not display a consistent signal in the region, signals from FAIRE and anti-H3K4Me2 ChIP were specifically present in the *KLF5*-expressing lines and were co-located with the conservation peak and DHS from the

ENCODE data at rs9600103, providing strong evidence for open chromatin and transcription factor binding here (**Figure 4a**). We then conducted chromatin conformation capture experiments for the *KLF5*-expressing Ishikawa endometrial cancer cells and we found a significant interaction between the *NcoI* restriction fragment containing the rs11841589/rs9600103 risk loci SNPs and the promoter region of *KLF5* (**Figure 4b**).

The regulatory nature of the region around rs9600103 and rs11841589 was investigated using allele-specific luciferase enhancer reporter assays in Ishikawa cells (**Figure 4c**). Paired t-tests were used to compare the relationships between fragments containing the rs11841589 and rs9600103 alleles, and the pGL3-Promoter reporter vector (no insert) control (**Supplementary Table 8**). Fragments containing the rs9600103-T, rs11841589-T and rs11841589-G alleles had activity significantly lower than that of the pGL3-Promoter control ( $P \leq 0.014$ ). In contrast, the construct containing the rs9600103-A risk allele had luciferase expression similar to the pGL3-Promoter control ( $P = 0.23$ ) and significantly higher than that of rs9600103-T ( $P = 0.02$ ), rs11841589-T ( $P = 0.05$ ) and rs11841589-G ( $P = 0.04$ ). These results suggest that the EC risk tagged by rs11841589 is at least partly due to a regulatory element containing rs9600103, which interacts with the *KLF5* promoter region, and the risk rs9600103-A allele is likely associated with increased gene expression.

In summary, this meta-analysis identified five novel EC risk loci at genome-wide significance, bringing the total number of common EC risk loci identified by GWAS to seven (**Figure 1**). Together with other risk SNPs reaching study-wide significance<sup>32,53,54</sup>, these explain ~1.6% of the EC familial relative risk. Novel EC risk SNPs lie in likely enhancers predicted to regulate genes or miRNAs with known or suspected roles in tumorigenesis, and we specifically showed that a functional SNP at 13q22.1 may sit within a transcriptional repressor of *KLF5*. Our findings further clarify the genetic etiology of EC, provide regions for functional follow-up, and add key information for future risk stratification models.

## Methods

Methods and any associated references are available in the online version of the paper at <http://www.nature.com/naturegenetics/>.



## Acknowledgments

The authors thank the many individuals who participated in this study and the numerous institutions and their staff who supported recruitment, detailed in full in the Supplementary Text.

The iCOGS endometrial cancer analysis was supported by NHMRC project grant [ID#1031333] to ABS, DFE and AMD. ABS, PW, GWM, and DRN are supported by the NHMRC Fellowship scheme. AMD was supported by the Joseph Mitchell Trust. IT is supported by Cancer Research UK and the Oxford Comprehensive Biomedical Research Centre. THTC is supported by the Rhodes Trust and the Nuffield Department of Medicine. Funding for the iCOGS infrastructure came from: the European Community's Seventh Framework Programme under grant agreement no 223175 [HEALTH-F2-2009-223175] [COGS], Cancer Research UK [C1287/A10118, C1287/A 10710, C12292/A11174, C1281/A12014, C5047/A8384, C5047/A15007, C5047/A10692, C8197/A16565], the National Institutes of Health [CA128978] and Post-Cancer GWAS initiative [1U19 CA148537, 1U19 CA148065 and 1U19 CA148112 - the GAME-ON initiative], the Department of Defence [W81XWH-10-1-0341], the Canadian Institutes of Health Research [CIHR] for the CIHR Team in Familial Risks of Breast Cancer, Komen Foundation for the Cure, the Breast Cancer Research Foundation, and the Ovarian Cancer Research Fund.

ANECs recruitment was supported by project grants from the NHMRC [ID#339435], The Cancer Council Queensland [ID#4196615] and Cancer Council Tasmania [ID#403031 and ID#457636]. SEARCH recruitment was funded by a programme grant from Cancer Research UK [C490/A10124]. Stage 1 and stage 2 case genotyping was supported by the NHMRC [ID#552402, ID#1031333]. Control data were generated by the Wellcome Trust Case Control Consortium (WTCCC), and a full list of the investigators who contributed to the generation of the data is available from the WTCCC website. We acknowledge use of DNA from the British 1958 Birth Cohort collection, funded by the Medical Research Council grant G0000934 and the Wellcome Trust grant 068545/Z/02 - funding for this project was provided by the Wellcome Trust under award 085475. NSECG was supported by the EU FP7 CHIBCHA grant, Wellcome Trust Centre for Human Genetics Core Grant 090532/Z/09Z, and CORGI was funded by Cancer Research UK. Recruitment of the QIMR Berghofer controls was supported by the NHMRC. The University of Newcastle, the Gladys M Brawn Senior Research Fellowship scheme, The Vincent Fairfax Family Foundation, the Hunter Medical Research Institute and the Hunter Area Pathology Service all contributed towards the costs of establishing the Hunter Community Study.

The Bavarian Endometrial Cancer Study (BECS) was partly funded by the ELAN fund of the University of Erlangen. The Hannover-Jena Endometrial Cancer Study was partly supported by the Rudolf Bartling Foundation. The Leuven Endometrium Study (LES) was supported by the Verelst Foundation for endometrial cancer. The Mayo Endometrial Cancer Study (MECS) and Mayo controls (MAY) were supported by grants from the National Cancer Institute of United States Public Health Service [R01 CA122443, P30 CA15083, P50 CA136393, and GAME-ON the NCI Cancer Post-GWAS Initiative U19 CA148112], the Fred C and Katherine B Andersen Foundation, the Mayo Foundation, and the Ovarian Cancer Research Fund with support of the Smith family, in memory of Kathryn Sladek Smith. MoMaTEC received financial support from a Helse Vest Grant, the University of Bergen, Melzer Foundation, The Norwegian Cancer Society (Harald Andersens legat), The Research Council of Norway and Haukeland University Hospital. The Newcastle Endometrial Cancer Study (NECS) acknowledges contributions from the University of Newcastle, The NBN Children's Cancer Research Group, Ms Jennie Thomas and the Hunter Medical Research Institute. RENDOCAS was supported through the regional agreement on medical training and clinical research (ALF) between Stockholm County Council and Karolinska Institutet [numbers: 20110222, 20110483, 20110141 and DF 07015], The Swedish Labor Market Insurance [number 100069] and The Swedish Cancer Society [number 11 0439]. The Cancer Hormone Replacement Epidemiology in Sweden Study (CAHRES, formerly called The Singapore and Swedish Breast/Endometrial Cancer Study; SASBAC) was supported by funding from the Agency for Science, Technology and Research of Singapore (A\*STAR), the US National Institutes of Health and the Susan G. Komen Breast Cancer Foundation.

The Breast Cancer Association Consortium (BCAC) is funded by Cancer Research UK [C1287/A10118, C1287/A12014]. The Ovarian Cancer Association Consortium (OCAC) is supported by a grant from the Ovarian Cancer Research Fund thanks to donations by the family and friends of Kathryn Sladek Smith [PPD/RPCI.07], and the UK National Institute for Health Research Biomedical Research Centres at the University of Cambridge. Additional funding for individual control groups is detailed in the Supplementary Information

**Table 1: Risk loci associated with EC at  $P < 5 \times 10^{-8}$  in the meta-analysis.**

Locus	SNP	Position	Nearby gene(s)	EA	OA	EAF	All histologies Allelic OR (95%CI)	$P$	$I^2$	Endometrioid histology Allelic OR (95%CI)	$P$	$I^2$
<b>Novel GWAS loci</b>												
13q22.1	rs11841589	73,814,891	<i>KLF5, KLF12</i>	G	T	0.74	1.15 (1.11-1.21)	$4.83 \times 10^{-11}$	0.19	1.16 (1.10-1.21)	$6.01 \times 10^{-10}$	0.00
6q22.31	rs13328298	126,016,580	<i>HEY2, NCOA7</i>	G	A	0.58	1.13 (1.09-1.18)	$3.73 \times 10^{-10}$	0.00	1.15 (1.11-1.20)	$1.02 \times 10^{-11}$	0.00
8q24.21	rs4733613	129,599,278	<i>MYC</i>	G	C	0.87	0.84 (0.80-0.89)	$3.09 \times 10^{-9}$	0.00	0.84 (0.79-0.89)	$7.70 \times 10^{-9}$	0.09
15q15.1	rs937213	40,322,124	<i>EIF2AK, BMF</i>	T	C	0.58	0.90 (0.86-0.93)	$1.77 \times 10^{-8}$	0.36	0.90 (0.86-0.94)	$2.22 \times 10^{-7}$	0.30
14q32.33	rs2498796	105,243,220	<i>AKT1, SIVA1</i>	G	A	0.70	0.89 (0.85-0.93)	$3.55 \times 10^{-8}$	0.00	0.88 (0.85-0.92)	$4.22 \times 10^{-8}$	0.00
<b>Previously reported GWAS loci</b>												
17q12	rs11263763	36,103,565	<i>HNF1B</i>	A	G	0.54	1.20 (1.15-1.25)	$2.78 \times 10^{-19}$	0.37	1.20 (1.15-1.25)	$6.51 \times 10^{-17}$	0.52
15q21	rs2414098	51,537,806	<i>CYP19A1</i>	C	T	0.62	1.17 (1.13-1.23)	$4.51 \times 10^{-13}$	0.00	1.18 (1.13-1.23)	$2.48 \times 10^{-13}$	0.00

Positions in build 37; EA, Effect allele; OA, Other allele; EAF, effect allele frequency;  $I^2$ , heterogeneity  $I^2$  statistic<sup>55</sup>. For all novel loci, the lead SNP was either directly genotyped or imputed with an information score of more than 0.9. *HNF1B* and *CYP19A1* have been previously reported by Painter *et al.*<sup>8</sup> and Thompson *et al.*<sup>9</sup>.

**Table 2: Conditional analysis of 8q24 locus showing two independent association signals.**

SNP	Position	EA	OA	EAF	Pairwise $r^2$ with		All histology meta-analysis		Conditioning on rs4733613		Conditioning on rs17232730	
					rs4733613	rs17232730	Allelic OR (95%CI)	$P$	Allelic OR (95%CI)	$P$	Allelic OR (95%CI)	$P$
rs4733613	129,599,278	G	C	0.87	-	0.02	0.84 (0.79-0.89)	$5.64 \times 10^{-9}$	-	-	0.86 (0.81-0.91)	$2.32 \times 10^{-7}$
rs17232730	129,537,746	G	C	0.88	0.02	-	1.17 (1.10-1.24)	$4.46 \times 10^{-7}$	1.14 (1.08-1.22)	$1.29 \times 10^{-5}$	-	-
rs10088218*	129,543,949	G	A	0.87	0.02	0.43	1.14 (1.07-1.20)	$1.65 \times 10^{-5}$	1.12 (1.05-1.18)	$2.92 \times 10^{-4}$	1.01 (0.91-1.12)	0.818

Positions in build 37; EA, Effect allele; OA, Other allele; EAF, effect allele frequency.

\*rs10088218 is associated with ovarian cancer (all subtypes), with the association being more significant for cancers of serous histology.

rs10088218-G is the risk allele for both EC and ovarian cancer.

## References

1. Siegel, R., Ma, J., Zou, Z. & Jemal, A. Cancer statistics, 2014. *CA. Cancer J. Clin.* **64**, 9–29 (2014).
2. Ferlay, J. *et al.* Cancer incidence and mortality patterns in Europe: estimates for 40 countries in 2012. *Eur. J. Cancer Oxf. Engl. 1990* **49**, 1374–1403 (2013).
3. Gruber, S. B. & Thompson, W. D. A population-based study of endometrial cancer and familial risk in younger women. Cancer and Steroid Hormone Study Group. *Cancer Epidemiol. Biomark. Prev. Publ. Am. Assoc. Cancer Res. Cosponsored Am. Soc. Prev. Oncol.* **5**, 411–417 (1996).
4. Win, A. K., Reece, J. C. & Ryan, S. Family history and risk of endometrial cancer: a systematic review and meta-analysis. *Obstet. Gynecol.* **125**, 89–98 (2015).
5. Barrow, E., Hill, J. & Evans, D. G. Cancer risk in Lynch Syndrome. *Fam. Cancer* **12**, 229–240 (2013).
6. Church, D. N. *et al.* DNA polymerase  $\epsilon$  and  $\delta$  exonuclease domain mutations in endometrial cancer. *Hum. Mol. Genet.* **22**, 2820–2828 (2013).
7. Palles, C. *et al.* Germline mutations affecting the proofreading domains of POLE and POLD1 predispose to colorectal adenomas and carcinomas. *Nat. Genet.* **45**, 136–144 (2013).
8. Painter, J. N. *et al.* Fine-mapping of the HNF1B multicancer locus identifies candidate variants that mediate endometrial cancer risk. *Hum. Mol. Genet.* (2014).  
doi:10.1093/hmg/ddu552
9. Thompson, D. J. *et al.* CYP19A1 fine-mapping and Mendelian randomisation: estradiol is causal for endometrial cancer. *Endocr. Relat. Cancer* (2015). doi:10.1530/ERC-15-0386
10. Spurdle, A. B. *et al.* Genome-wide association study identifies a common variant associated with risk of endometrial cancer. *Nat. Genet.* **43**, 451–454 (2011).
11. Tomlinson, I. *et al.* A genome-wide association scan of tag SNPs identifies a susceptibility variant for colorectal cancer at 8q24.21. *Nat. Genet.* **39**, 984–988 (2007).

12. Tenesa, A. *et al.* Genome-wide association scan identifies a colorectal cancer susceptibility locus on 11q23 and replicates risk loci at 8q24 and 18q21. *Nat. Genet.* **40**, 631–637 (2008).
13. Wellcome Trust Case Control Consortium. Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature* **447**, 661–678 (2007).
14. Pharoah, P. D. P. *et al.* GWAS meta-analysis and replication identifies three new susceptibility loci for ovarian cancer. *Nat. Genet.* **45**, 362–370, 370e1–2 (2013).
15. Sakoda, L. C., Jorgenson, E. & Witte, J. S. Turning of COGS moves forward findings for hormonally mediated cancers. *Nat. Genet.* **45**, 345–348 (2013).
16. Michailidou, K. *et al.* Large-scale genotyping identifies 41 new loci associated with breast cancer risk. *Nat. Genet.* **45**, 353–361 (2013).
17. Eeles, R. A. *et al.* Identification of 23 new prostate cancer susceptibility loci using the iCOGS custom genotyping array. *Nat. Genet.* **45**, 385–391, 391e1–2 (2013).
18. De Vivo, I. *et al.* Genome-wide association study of endometrial cancer in E2C2. *Hum. Genet.* **133**, 211–224 (2014).
19. Setiawan, V. W. *et al.* Two estrogen-related variants in CYP19A1 and endometrial cancer risk: a pooled analysis in the Epidemiology of Endometrial Cancer Consortium. *Cancer Epidemiol. Biomark. Prev. Publ. Am. Assoc. Cancer Res. Cosponsored Am. Soc. Prev. Oncol.* **18**, 242–247 (2009).
20. GTEx Consortium. The Genotype-Tissue Expression (GTEx) project. *Nat. Genet.* **45**, 580–585 (2013).
21. Cancer Genome Atlas Research Network *et al.* Integrated genomic characterization of endometrial carcinoma. *Nature* **497**, 67–73 (2013).
22. Katoh, M. & Katoh, M. Integrative genomic analyses on HES/HEY family: Notch-independent HES1, HES3 transcription in undifferentiated ES cells, and Notch-dependent HES1, HES5, HEY1, HEY2, HEYL transcription in fetal tissues, adult tissues, or cancer. *Int. J. Oncol.* **31**, 461–466 (2007).

23. Shao, W., Halachmi, S. & Brown, M. ERAP140, a conserved tissue-specific nuclear receptor coactivator. *Mol. Cell. Biol.* **22**, 3358–3372 (2002).
24. Rothman, N. *et al.* A multi-stage genome-wide association study of bladder cancer identifies multiple susceptibility loci. *Nat. Genet.* **42**, 978–984 (2010).
25. Kiemeny, L. A. *et al.* Sequence variant on 8q24 confers susceptibility to urinary bladder cancer. *Nat. Genet.* **40**, 1307–1312 (2008).
26. Easton, D. F. *et al.* Genome-wide association study identifies novel breast cancer susceptibility loci. *Nature* **447**, 1087–1093 (2007).
27. Whiffin, N. *et al.* Identification of susceptibility loci for colorectal cancer in a genome-wide meta-analysis. *Hum. Mol. Genet.* **23**, 4729–4737 (2014).
28. Goode, E. L. *et al.* A genome-wide association study identifies susceptibility loci for ovarian cancer at 2q31 and 8q24. *Nat. Genet.* **42**, 874–879 (2010).
29. Eeles, R. A. *et al.* Identification of seven new prostate cancer susceptibility loci through a genome-wide association study. *Nat. Genet.* **41**, 1116–1121 (2009).
30. Gudmundsson, J. *et al.* Genome-wide association and replication studies identify four variants associated with prostate cancer susceptibility. *Nat. Genet.* **41**, 1122–1126 (2009).
31. Huppi, K., Pitt, J. J., Wahlberg, B. M. & Caplen, N. J. The 8q24 gene desert: an oasis of non-coding transcriptional activity. *Front. Genet.* **3**, 69 (2012).
32. Carvajal-Carmona, L. G. *et al.* Candidate locus analysis of the TERT-CLPTM1L cancer risk region on chromosome 5p15 identifies multiple independent variants associated with endometrial cancer risk. *Hum. Genet.* (2014). doi:10.1007/s00439-014-1515-4
33. Berlanga, J. J., Santoyo, J. & De Haro, C. Characterization of a mammalian homolog of the GCN2 eukaryotic initiation factor 2alpha kinase. *Eur. J. Biochem. FEBS* **265**, 754–762 (1999).
34. Uhlén, M. *et al.* Proteomics. Tissue-based map of the human proteome. *Science* **347**, 1260419 (2015).
35. Cantley, L. C. The phosphoinositide 3-kinase pathway. *Science* **296**, 1655–1657 (2002).

36. Slomovitz, B. M. & Coleman, R. L. The PI3K/AKT/mTOR pathway as a therapeutic target in endometrial cancer. *Clin. Cancer Res. Off. J. Am. Assoc. Cancer Res.* **18**, 5856–5864 (2012).
37. Salvesen, H. B. *et al.* Integrated genomic profiling of endometrial carcinoma associates aggressive tumors with indicators of PI3 kinase activation. *Proc. Natl. Acad. Sci. U. S. A.* **106**, 4834–4839 (2009).
38. Shoji, K. *et al.* The oncogenic mutation in the pleckstrin homology domain of AKT1 in endometrial carcinomas. *Br. J. Cancer* **101**, 145–148 (2009).
39. Cohen, Y. *et al.* AKT1 pleckstrin homology domain E17K activating mutation in endometrial carcinoma. *Gynecol. Oncol.* **116**, 88–91 (2010).
40. Du, W. *et al.* Suppression of p53 activity by Siva1. *Cell Death Differ.* **16**, 1493–1504 (2009).
41. Wang, X. *et al.* Siva1 inhibits p53 function by acting as an ARF E3 ubiquitin ligase. *Nat. Commun.* **4**, 1551 (2013).
42. Li, N. *et al.* Siva1 suppresses epithelial-mesenchymal transition and metastasis of tumor cells by inhibiting stathmin and stabilizing microtubules. *Proc. Natl. Acad. Sci. U. S. A.* **108**, 12851–12856 (2011).
43. Dun, B. *et al.* Mycophenolic acid inhibits migration and invasion of gastric cancer cells via multiple molecular pathways. *PloS One* **8**, e81702 (2013).
44. Shi, H., Zhang, Z., Wang, X., Liu, S. & Teng, C. T. Isolation and characterization of a gene encoding human Kruppel-like factor 5 (IKLF): binding to the CAAT/GT box of the mouse lactoferrin gene promoter. *Nucleic Acids Res.* **27**, 4807–4815 (1999).
45. Simmen, R. C. M. *et al.* The emerging role of Krüppel-like factors in endocrine-responsive cancers of female reproductive tissues. *J. Endocrinol.* **204**, 223–231 (2010).
46. Mutter, G. L. *et al.* Global expression changes of constitutive and hormonally regulated genes during endometrial neoplastic transformation. *Gynecol. Oncol.* **83**, 177–185 (2001).



47. Davis, H. *et al.* FBXW7 mutations typically found in human cancers are distinct from null alleles and disrupt lung development. *J. Pathol.* **224**, 180–189 (2011).
48. Nandan, M. O. *et al.* Krüppel-like factor 5 mediates cellular transformation during oncogenic KRAS-induced intestinal tumorigenesis. *Gastroenterology* **134**, 120–130 (2008).
49. Forbes, S. A. *et al.* The Catalogue of Somatic Mutations in Cancer (COSMIC). *Curr. Protoc. Hum. Genet. Editor. Board Jonathan Haines* **AI Chapter 10**, Unit 10.11 (2008).
50. ENCODE Project Consortium *et al.* Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project. *Nature* **447**, 799–816 (2007).
51. Rao, S. S. P. *et al.* A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell* **159**, 1665–1680 (2014).
52. Wang, Y., Li, X. & Hu, H. H3K4me2 reliably defines transcription factor binding regions in different cells. *Genomics* **103**, 222–228 (2014).
53. O'Mara, T. A. *et al.* Comprehensive genetic assessment of the ESR1 locus identifies a risk region for endometrial cancer. *Endocr. Relat. Cancer* **22**, 851–861 (2015).
54. Cheng, T. H. *et al.* Meta-analysis of genome-wide association studies identifies common susceptibility polymorphisms for colorectal and endometrial cancer near SH2B3 and TSHZ1. *Sci. Rep.* **5**, 17369 (2015).
55. Higgins, J. P. T. & Thompson, S. G. Quantifying heterogeneity in a meta-analysis. *Stat. Med.* **21**, 1539–1558 (2002).

## Figure legends

### Figure 1: EC meta-analysis Manhattan plot

Manhattan plot of  $-\log_{10}$ -transformed P-values from meta-analysis of 22 autosomes. There are seven loci surpassing genome wide significance including two known loci: 15q21 (*CYP19A1*) and 17q12 (*HNF1B*) and five novel loci: 6q22 (*NCOA7*, *HEY2*), 8q24 (*MYC*), 13q22 (*KLF5*), 14q32 (*AKT1*, *SIVA1*), 15q15 (*EIF2AK4*, *BMF*).

### Figure 2: Forest plots of novel EC risk loci

The odds ratio and 95% confidence intervals of each study of the meta-analysis are listed and shown in the adjacent plot. The  $I^2$  heterogeneity scores (all  $<0.4$ ) suggest that there is no marked difference in effects between studies. The SNPs represented are: a) rs11841589 (13q22), b) rs13328298 (6q22), c) rs4733613 (8q24), d) rs17232730 (8q24, pairwise  $r^2$  0.02 with rs4733613), e) rs937213 (15q15) and f) rs2498796 (14q32).

### Figure 3: Regional association plots for the five novel loci associated with EC.

The  $-\log_{10}$  P-values from the meta-analysis and regional imputation for three GWAS and eight iCOGS groups are shown for SNPs at: a) 13q22.1, b) 6q22, c) & d) 8q24, e) 15q15 and f) 14q32.33. The SNP with the lowest P-value at each locus is labeled and marked as a purple diamond, and the dot color represents the LD with the top SNP. The blue line shows recombination rates in cM/Mb. All plotted SNPs are either genotyped or have an IMPUTE info score of more than 0.9 in all datasets. **Supplementary Figure 4** displays similar regional association plots with a larger number of SNPs using a less stringent info score cut-off.

### Figure 4: The 13q22.1 EC susceptibility locus

a) Diagram showing the 16kb region (position 73,804,930- 73,820,618) around rs11841589, rs9600103 and correlated SNPs rs7981863, rs7988505 and rs7989799 (black marks).

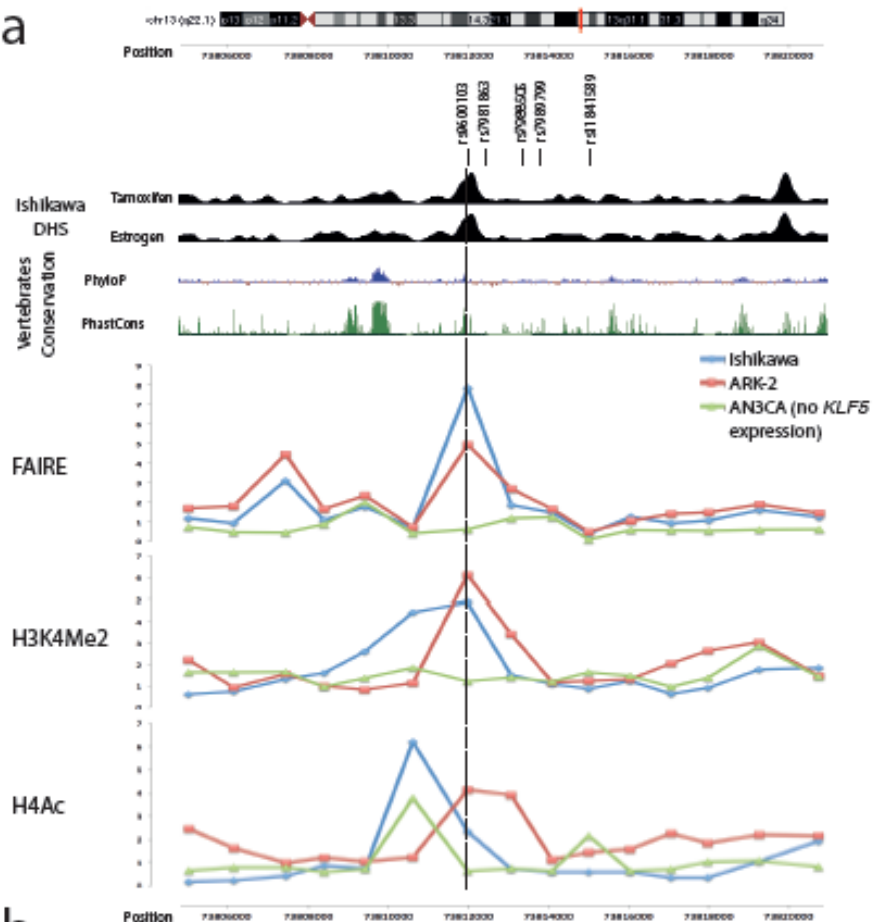
FAIRE and ChIP assays with anti-H3K4Me2 and anti-H4Ac antibodies for three EC cell lines ARK-2 (rs9600103-TT), Ishikawa (rs9600103-AA) and AN3CA (rs9600103-AA) are shown, with the y-axis displaying enrichment normalized to non-crosslinked genomic DNA/sonicated input DNA, relative to the *Rhodopsin* promoter as a negative control using the  $\Delta\Delta C_t$  method. DNaseI hypersensitivity site (DHS) density signal in ENCODE EC Ishikawa cells (**Supplementary Note**) are shown, from experiments with cell lines treated with estrogen and tamoxifen. 100 vertebrates conservation is also displayed. Vertical dotted line represents the position of rs9600103.

b) *3C experiment for KLF5-expressing Ishikawa cells.* Relative interaction frequencies between an *NcoI* restriction fragment containing risk SNPs rs9600103 and rs11841589 (bait fragment) with *NcoI* fragments across the region were calculated using qPCR with normalization to the signal from a control BAC 3C library and a non-interacting chromosomal region, using the  $\Delta\Delta C_t$  method. The graph shows the frequencies plotted against the fragment position on chromosome 13. A significant interaction is seen with the fragment containing a *KLF5* transcriptional start site in Ishikawa cells.

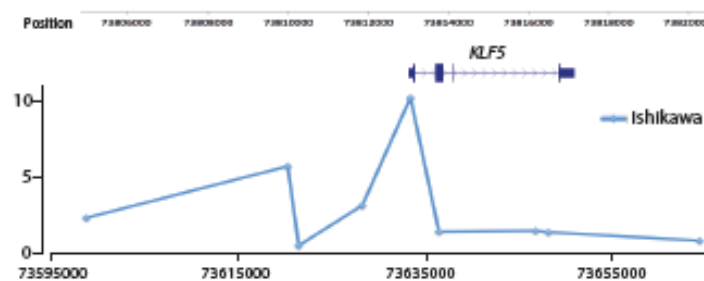
c) *Luciferase reporter assay to analyze the activity of 3kb fragments containing either rs9600103 or rs11841589 using the pGL3 promoter vector in Ishikawa cells.* Green arrows represent the low-risk alleles, and red arrows the high-risk alleles. Error bars represent the standard error of the mean. Data were normalized by subtraction of background luminescence and normalized to pGL4 Renilla activity. Luciferase activity in the rs9600103-A risk allele was more than double than that of the rs9600103-T protective allele ( $P=0.018$ ). Paired t-tests between the different fragments also showed that the rs9600103-A high-risk allele has significantly higher expression compared with both rs11841589 alleles (0.045, 0.039) (**Supplementary Table 8**). Schematic diagram displays position on chromosome 13 of the fragment sequences and the arrows represents the position of the two SNPs.

Fig. 4

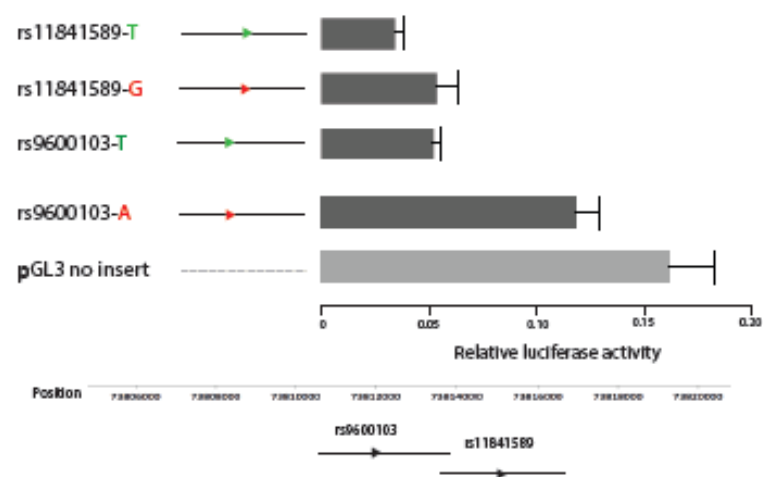
a



b



c



## ONLINE METHODS

Cases and controls were matched as summarized in **Supplementary Table 1** and a detailed description of each sample set can be found in the **Supplementary Note**. **Supplementary Figure 1** is a flow diagram that illustrates the overall study design.

### Additional EC GWAS

The National Study of Endometrial Cancer Genetics (NSECG) consisted of 925 histologically confirmed endometrial cancer (EC) cases from the UK. 86% of these cases had endometrioid-only histology and genotyping was done using Illumina 660W Quad arrays.

These cases were matched with 1,286 cancer-free controls from the UK1/CORGI<sup>1</sup> and SP1<sup>2</sup> colorectal studies with genotyping conducted on Illumina Hap550, Illumina Hap300 and Illumina Hap240S arrays. Additionally, publically available 1958 Birth Cohort<sup>3</sup> controls from the Wellcome Trust Case Control Consortium (WTCCC2)<sup>4</sup> genotyped on Illumina Infinium 1.2M arrays were included.

### Original EC GWAS

As described previously, cases with confirmed endometrioid histology were selected from two population studies of endometrial cancer; the UK Studies of Epidemiology and Risk factors in Cancer Heredity (SEARCH, n=681) and the Australian National Endometrial Cancer Study (ANECS, n=606), and genotypes generated using Illumina Infinium 610K arrays<sup>5</sup>. Compared with our previous study<sup>5</sup>, ANECS and SEARCH were analysed as two groups and additional controls<sup>6,7</sup> were used for this meta-analysis. SEARCH cases were compared with 2,501 controls from the National Blood Service (NBS) part of the WTCCC2 controls<sup>4</sup>. ANECS cases were compared to controls recruited as part of the Hunter Community Study<sup>6</sup> or Brisbane Adolescent Twin Study<sup>8</sup>, both genotyped using the Illumina Infinium 610K arrays.

### Phase 1 iCOGS replication

For the iCOGS genotyping stage of the study, 4,330 women with a confirmed diagnosis of endometrial cancer and European ancestry were recruited via 11 separate studies in Western Europe, North America and Australia, collectively called the Endometrial Cancer Association Consortium (ECAC).

Healthy female controls with European ancestry and known age at sampling were selected from controls genotyped by the Breast Cancer Association Consortium (BCAC)<sup>9</sup> iCOGS project, or the Ovarian Cancer Association Consortium (OCAC)<sup>10</sup> iCOGS project. The eight case-control groups were matched based on geographical location, and principal components analysis (PCA) was conducted such that individuals who clustered outside the main centroid in pairwise plots of the first four PCs were excluded (**Supplementary Figure 2**).

Cases and controls were genotyped on a custom Illumina Infinium iSelect array with 211,155 SNPs, designed by the Collaborative Oncological Gene-environment Study (iCOGS), a collaborative project involving four consortia. SNPs were included on this array based on promising regions of interest in previous breast, ovarian and prostate<sup>11</sup> studies, and also the 1,483 top SNPs from our previous EC GWAS<sup>5</sup> analysis. Cases and MoMaTEC

controls were genotyped by Genome Quebec Innovation Center. BCAC and OCAC control samples were genotyped at four centres. Raw intensity data files for all consortia were sent to the COGS data co-ordination centre at the University of Cambridge for centralized genotype calling and QC, so that all case and control genotypes were called using the same procedure.

## SNP genotyping arrays quality control

Genotype calling was done using Illumina's proprietary Gencall algorithm and Illuminus<sup>12</sup>. Duplicate samples displayed >99% concordance. Standard quality control measures applied to genotyping arrays are described in our original GWAS<sup>5</sup> and these include genotypic call rate <0.95, deviation from Hardy-Weinberg Equilibrium (HWE) at  $P < 10^{-6}$  and visual inspection of cluster plots for most significant SNPs. For iCOGS, all EC cases and MoMaTEC controls were genotyped by Genome Quebec Innovation Center. BCAC and OCAC control samples were genotyped at four centres. Raw intensity data files for all consortia were sent to the COGS data co-ordination centre at the University of Cambridge for centralized genotype calling and QC, so that all case and control genotypes were called using the same procedure. Duplicate samples for quality showed a concordance of >99%. Samples were excluded based on the following measures: missingness >5%, heterozygosity rates  $((N-O)/N) > 5$  S.D from the mean, X chromosome heterozygosity rate (PLINK F-score) >0.2, and pairwise identity by descent (IBD) >0.1875 (cut-off for second-degree relatives). Principal components analysis (PCA) was conducted using Eigenstrat<sup>13</sup> software. Analysis was conducted using PLINK<sup>14</sup>, and R packages GenABEL and SNPMatrix<sup>15,16</sup>.

## Phase 2 NSECG replication

The second replication phase consisted of directly genotyping five SNPs with  $P < 10^{-7}$  and IMPUTE info scores of >0.94 from the NSECG/ANECs/SEARCH/iCOGS meta-analysis. Genotyping was done in NSECG samples that had not previously been used in the NSECG GWAS or the NSECG iCOGS. Genotyping was conducted using competitive allele-specific PCR (KASPar, KBiosciences) and the Fluidigm BioMark<sup>TM</sup> HD System, using standard protocols. The genotyping call rate was >0.98 and there was a >0.985 concordance between different genotyping platforms (**Supplementary Table 4**). There was no significant deviation from HWE ( $P > 0.05$ ). Primers used for genotyping are listed in **Supplementary Table 5**.

## Genome-wide and regional imputation

Genome-wide imputation for all SNP array generated data was conducted using IMPUTE v2<sup>17</sup> using 1000 Genomes project (2012 release) as a reference panel. For the first-pass genome-wide analysis we pre-phased chromosomes using SHAPEIT<sup>18</sup> to improve the computational speed. Imputation was carried out separately for each of the three GWAS studies (for each GWAS study the cases and controls were imputed together as a single dataset, using only SNPs which passed QC in both cases and controls) and for the iCOGS study (all studies within iCOGS were imputed together). SNPs with  $MAF < 0.1\%$  were removed from all studies prior to imputation. Genome-wide imputation produced 9,594,066 SNPs with  $MAF \geq 1\%$  and  $info \geq 0.4$  in at least one of the three GWAS and eight iCOGS groups. Of these, 8,308,423 SNPs met these criteria in all studies. The iCOGS genotyping array (~200,000 SNPs) is aimed at capturing previously prioritised cancer SNPs and not genome-wide coverage, but nonetheless 8,631,871 SNPs met  $MAF \geq 1\%$  and  $info \geq 0.4$  criteria, of which 5,437,135 had  $info \geq 0.7$  and 2,333,040 had  $info \geq 0.9$ .

Regional imputation of regions of interest (1.5Mb region around SNPs with meta-analysis  $P < 10^{-5}$ ) used both 1000 Genomes 2012 release and 196 high-coverage, whole genome-sequenced UK individuals as reference panels as a means to improve imputation accuracy<sup>19</sup>. All SNPs reported in this study had an info score of more than 0.9 in all datasets.

## Association testing

Association testing was done using SNPTEST v2<sup>20</sup> employing frequentist tests with a logistic regression model for each of the 11 groups as matched in **Supplementary Table 1**. There was little evidence of systematic over-dispersion of the test statistic from the quantile-quantile plots (**Supplementary Figure 3**) and the genomic inflation  $\lambda_{GC}$ , calculated using all genotyped SNPs passing QC for the three GWAS. For iCOGS, 105,000 SNPs after LD-pruning ( $r^2 < 0.2$ ) and  $> 500$ kb from the 1,483 EC prioritized SNPs on the iCOGS were used.  $\lambda_{GC}$  was between 1.002 and 1.038 for each study. Conditional logistic regression analysis was conducted for each locus of genome-wide significance using SNPTEST to look for the presence of multiple independent association signals. This was done in a stepwise manner, first conditioning for the most significant SNP and subsequently for any SNPs that remained significant at  $P_{cond} < 10^{-4}$ . Regional association plots (**Figure 1, Supplementary Figure 4**) were created using LocusZoom<sup>21</sup>.

## Meta-analysis

Inverse variance, fixed effects meta-analysis of the 11 groups (three GWAS, eight iCOGS groups) was conducted using GWAMA<sup>22</sup>. The per allele effect size of each SNP in a particular study is represented by  $\beta$  (the log-odds ratio) and its standard error. Inter-study differences are represented by the  $I^2$  heterogeneity score<sup>23,24</sup>. Forest plots of the genome-wide significant loci (**Figure 2**) provided a visual representation of risk effects across different studies and these were made using the rmeta package (<http://cran.r-project.org/web/packages/rmeta/>). A random-effects meta-analysis was also performed for SNPs with  $I^2 > 0.3$ . The results of the second replication phase (NSEC replication) were meta-analyzed in a 12-way meta-analysis for the top 5 SNPs yielding a total of 7,737 EC cases and 37,144 controls. 6,635 (86%) of the EC cases had endometrioid-only histology and association testing and meta-analysis were also conducted with just these samples.

## Bioinformatic analysis and functional annotation of genome-wide significant EC risk loci

The five novel genome-wide significant loci and SNPs in LD ( $r^2 > 0.7$  in European 1000 Genomes) were annotated using HaploregV2<sup>25</sup>, RegulomeDB<sup>26</sup> and data from ENCODE<sup>27</sup> in **Supplementary Table 6**. This includes information such as promoter and enhancer histone marks, DHS, bound proteins, altered motifs, GENCODE and dbSNP annotations, RegulomeDB score and PhastCons conservation scores.

Bioinformatic analysis in **Supplementary Figure 6** used datasets described by Hnisz *et al.*<sup>28</sup> and Corradin *et al.*<sup>29</sup> in order to identify likely enhancers in a cell-specific context for the EC risk loci. Enhancer-gene interactions are predicted by identifying 'super-enhancers' (regions containing neighbouring H3K27Ac modifications) from 86 cell and tissue types and then the expressed transcript with transcription start site closest to the centre of the super-enhancer was assigned as the target gene. PresTIGE pairs cell-type specific H3K4Me1 and gene expression data from 13 cell types to identify likely enhancer-gene interactions.

## **Endometrial-tissue expression quantitative trait loci (eQTL) analysis for associated SNPs using GTEx and TCGA data**

Publicly available data generated by the Genotype-Tissue Expression Project (GTEx)<sup>30</sup> and The Cancer Genome Atlas (TCGA; [www.cancergenome.nih.gov](http://www.cancergenome.nih.gov)) were accessed to examine tissue-specific eQTLs. For GTEx, expression and genotype data were generated from 70 normal uteri from post-mortem biopsies, using an Affymetrix Expression array and Illumina Omni 5M SNP array. GTEx provided processed results, evaluating association between genotype and expression data. The expression levels are represented as a rank normalized score. TCGA genotype and copy number variation (CNV) data were derived from Affymetrix 6.0 SNP arrays. Expression data were from RNAseq arrays (Illumina HiSeq and Illumina GA) for 458 endometrial cancer tissues and 30 adjacent normal endometrial tissues. Association analyses for TCGA datasets were performed as follows. Genes within 500kb flanking our SNPs of interest were selected for analysis. Since there may be significant variation in tumour tissue copy number, somatic CNVs were taken into account by regressing gene expression to average copy number spanning the gene. Residual unexplained variance in gene expression was then regressed on the genotype of the lead SNP at each locus, using genotyped or imputed data. Statistical comparisons were subject to Bonferroni correction for number of tests (number of sample sets, and number of genes assessed).

## **DNA and RNA extraction from EC cell lines**

Cells were snap frozen with dry ice after centrifugation, and DNA and RNA were extracted using DNeasy DNA extraction (Qiagen) and RNeasy minikit (Qiagen) according to manufacturer's instructions. Nucleic acids were then quantified using Nanodrop 2000 (ThermoScientific) spectrophotometry.

## **Quantification of *KLF5* expression in 11 EC cell lines**

Extracted RNA was treated with DNase 1 and complimentary DNA (cDNA) was reverse transcribed from RNA using the High Capacity cDNA Reverse Transcription Kit (Applied Biosystems), according to the manufacturer's protocol. TaqMan Gene Expression Assays were used for *KLF5* and *GAPDH* (details available from authors). The absolute expression of *KLF5* was quantified using qRT-PCR using the ABI 7900HT cycler (Applied Biosystems), and the critical threshold was manually set at 0.2. Standard protocols were applied to calculate relative expression using the  $\Delta\Delta CT$  method as described by Livak and Schmittgen<sup>31</sup> and *GAPDH* was used as an endogenous control.

## **Formaldehyde-assisted identification of regulatory elements (FAIRE)**

Formaldehyde-Assisted Isolation of Regulatory Elements (FAIRE) was conducted using the method adapted from Giresi et al<sup>32</sup>. Briefly, cross-linking was done on a rocker at room temperature. 1% formaldehyde was added to approximately  $10^8$  cells for 5 minutes, after which 115 mM glycine was added to inhibit the cross-linking. For each cell line, a non-crosslinked control was prepared in parallel for all of the remaining steps. After two rinses with 4°C phosphate buffered saline solution (PBS), the cells were suspended in successive lysis buffers: Lysis buffer I (50 mM HEPES-KOH, 140 mM NaCl, 1 mM EDTA, 10% glycerol, 0.5% NP-40, 0.25% tritonX-100); lysis buffer II (10 mM tris-HCl, 200 mM NaCl, 1 mM EDTA, 0.5 mM EGTA); and lysis buffer III (10 mM tris-HCl, 2100 mM NaCl, 1 mM EDTA, 0.1% sodium deoxycholate, 0.5%N-lauroylsarcosine). Cells were incubated on a rocker at 4°C for



10 minutes in each lysis buffer after which they were spun down at 1300 g for 5 minutes so that the supernatant could be removed. The cells were then sonicated using the Bioruptor in seven to fifteen 30-second cycles to generate fragments 100-1000 bp in size. Gel electrophoresis in 1% agarose was used to confirm the size of the DNA fragments. The DNA was extracted with a standard phenol/chloroform method and ethanol-precipitated. 50ng of DNA from paired crosslinked and non-crosslinked cells was analyzed in duplicate by SYBR-green quantitative PCR (qPCR using primers at roughly 1kb intervals in the 13q22.1 region downstream of *KLF5* (**Supplementary Table 7**). The  $\Delta\Delta C_t$  method<sup>31</sup> was used to normalize the results to the input DNA from the non-crosslinked cells and then expressed relative to the Rhodopsin promoter as a negative control. For each experiment there were two replicates for the crosslinked cells and non-crosslinked controls, each performed on two occasions.

### Cross-linked Chromatin immunoprecipitation (ChIP)

About  $10^8$  cells were cross-linked using 1% formaldehyde for 10 minutes. Glycine was used to stop the cross-linking, cells were then rinsed twice in PBS, and cell scrapers were used to detach cells adhered to the Petri dish surface. The cells were then resuspended in lysis buffer (1% sodium dodecyl sulfate (SDS), 10 mM EDTA (Ambion), 50mM Tris-HCl (Ambion)) and incubated for 10 minutes. The cells were then sonicated using the Bioruptor (Diagenode) in 7 to 15 30-second cycles to generate fragments 1000-1500 bp in size. Gel electrophoresis in 1% agarose confirmed the size of the DNA fragments. The fragmented DNA was then diluted ten times to the immuno-precipitation dilution buffer (1% tritonX-100, 2 nM EDTA, 20 mM Tris-HCl, 150 mM sodium chloride and each cell line was separated into four tubes: input chromatin, no-antibody-control and one tube for each antibody. 5 ul of anti-dimethyl-histone H3 Lys4 (Millipore 07-030) and anti-acetyl-histone H4 (Millipore 06-866) were added to the antibody tubes and, along with the no-antibody-control, incubated overnight at 4°C for immunoprecipitation. The input chromatin was kept refrigerated at 4°C until the reverse cross-linking of day 2. Phenylmethylsulfonyl fluoride and protease inhibitor was added to the lysis buffer and IP dilution buffer to deactivate proteases, while sodium butyrate was added to these solutions to inhibit histone deacetylases. 5 ul of protein A Dynabeads was added to each tube and incubated for 4 hours. A series of washes were done using Tris/Sucrose/EDTA (TSE) I (1% tritonX-100, 2 mM EDTA, 20 mM Tris-HCl, 150 mM NaCl, 0.1 %SDS), TSE II (1% tritonX-100, 2 mM EDTA, 20 mM Tris-HCl, 500 mM NaCl, 0.1% SDS), Buffer III (0.25 M lithium chloride, 1 mM EDTA, 10 mM Tris-HCl, 1% tergitol-type NP-40, 1% sodium deoxycholate) and tris-EDTA (Tris-EDTA 1X). 300 ul of extraction solution (1% SDS 0.1 M sodium bicarbonate) was added and the Dynabeads were removed after a 30 minute incubation. Then 0.7 M NaCl was added and reverse cross-linking occurred overnight at 65°C. DNA was then purified using the QIAquick PCR purification kit (Qiagen) according to the manufacturer's protocol. 1ul of DNA was analyzed in duplicate or triplicate by SYBR green qPCR as above and the  $\Delta\Delta C_t$  method was used to identify areas with enrichment. For each experiment there were two replicates for each anti-body along with the input and no-antibody control, each performed on two occasions. Primers used are listed in **Supplementary Table 7**.

### Chromatin conformation capture (3C)

Experiments were performed as described in Ghoussaini *et al.*<sup>33</sup> Briefly, *KLF5*-expressing Ishikawa endometrial cancer cell lines were crosslinked with 1% formaldehyde for 10 mins, quenched with 125mM glycine, washed with PBS and collected by scraping. Cells were lysed for 30 min on ice in 10 mM Tris-HCl, pH 7.5, 10 mM NaCl, 0.2% Igepal with protease inhibitors and homogenized in a Dounce homogenizer. Nuclei were pelleted and resuspended in 1ml 1.2X restriction buffer (NEB 3.1) with 0.3% SDS for 1h at 37°C. 2% Triton X-100 was added then 1000U NcoI was added 3 times over 24h at 37°C with shaking.

The enzyme was inactivated and digested DNA was diluted 8X before ligation with 4000U of T4 DNA ligase overnight at 16°C. Crosslinks were reversed by proteinase K digestion at 65°C overnight, and then the DNA was purified by phenol–chloroform extraction and ethanol precipitation. The final DNA pellet was dissolved in 10 mM Tris (pH 7.5) and purified through Amicon Ultra 0.5 ml columns (Millipore). 3C interactions were quantified by SYTO9 qPCR (performed on a RotorGene 6000) using primers designed to amplify across ligated NcoI restriction fragments with one constant primer within the risk fragment (including rs11841589 and rs9600103) and a series of test primers within NcoI fragments spanning 76 kb of the *KLF5* promoter region. BAC clones (RP11-81D9 and RP11-179I20) covering the region were digested with NcoI, ligated with T4 ligase and then used to determine PCR efficiency. 3C analyses were performed on three independent 3C libraries, with each data point in duplicate. Data were normalized to the signal from the BAC clone library and from a non-interacting chromosomal region using the  $\Delta\Delta C_t$  method with incorporated individual primer pair efficiencies.

### Luciferase reporter assays

For luciferase reporter assays, the regions chr:13 73,810,509-73,813,452 around rs9600103 and chr13:73,813,268-73,816,290 around rs11841589 were cloned into the pGL3-Promoter vector (Promega) to test for regulatory effects in Ishikawa cells. Ishikawa cells were selected because they express *KLF5*, showed evidence of a DHS, FAIRE and H3K4Me2 enrichment at rs9600103 and were readily transfectable. Site-directed mutagenesis was used so that both the high- and low-risk alleles of rs9600103 and rs11841589 were tested. After sequencing to verify the correct insert sequences, cells were transiently co-transfected using lipofectamine with the appropriate pGL3-Promoter constructs, and with the Renilla luciferase pGL4.75 vector (Promega) as a control for transfection efficiency. After 48 hours, luciferase activity was measured (Dual-Glo Luciferase Assay System, Promega), and after subtracting background from lipofectamine-only controls, firefly luciferase activity from the putative enhancer regions was normalized to the Renilla luciferase values for each sample. Levels of firefly luciferase activity were compared with a control plasmid consisting of an empty pGL3 and also a noncoding 2.2-kb stretch of plasmid sequence (taken from the pENTR1A plasmid, Invitrogen) cloned into the pGL3-Promoter vector that we had previously used as a length of DNA with no regulatory activity<sup>34</sup>. Luciferase activity experiments had three or four replicates, each performed on three occasions (total of 11 assays). Primers used in these experiments are listed in **Supplementary Table 5**.

ANOVA found significant differences in luciferase levels ( $P < 0.0001$ ,  $F: 11.6$ ) but no significant differences between replicates conducted on different days ( $P = 0.91$ ,  $F: 0.09$ ). There were no significant differences between the pENTR1A control and the empty pGL3 vector ( $P = 0.085$ ) and pGL3 no insert is used as the control. We then conducted paired t-tests for all comparisons using the average of biological repeats, between the pGL3 no insert, rs9600103-A, rs9600103-T, rs11841589-G and rs11841589-T fragments (**Supplementary Table 8**).

### References

1. Tomlinson, I. *et al.* A genome-wide association scan of tag SNPs identifies a susceptibility variant for colorectal cancer at 8q24.21. *Nat. Genet.* **39**, 984–988 (2007).

2. Tenesa, A. *et al.* Genome-wide association scan identifies a colorectal cancer susceptibility locus on 11q23 and replicates risk loci at 8q24 and 18q21. *Nat. Genet.* **40**, 631–637 (2008).
3. Power, C. & Elliott, J. Cohort profile: 1958 British birth cohort (National Child Development Study). *Int. J. Epidemiol.* **35**, 34–41 (2006).
4. Wellcome Trust Case Control Consortium. Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature* **447**, 661–678 (2007).
5. Spurdle, A. B. *et al.* Genome-wide association study identifies a common variant associated with risk of endometrial cancer. *Nat. Genet.* **43**, 451–454 (2011).
6. McEvoy, M. *et al.* Cohort profile: The Hunter Community Study. *Int. J. Epidemiol.* **39**, 1452–1463 (2010).
7. Painter, J. N. *et al.* Genome-wide association study identifies a locus at 7p15.2 associated with endometriosis. *Nat. Genet.* **43**, 51–54 (2011).
8. McGregor, B. *et al.* Genetic and environmental contributions to size, color, shape, and other characteristics of melanocytic naevi in a sample of adolescent twins. *Genet. Epidemiol.* **16**, 40–53 (1999).
9. Michailidou, K. *et al.* Large-scale genotyping identifies 41 new loci associated with breast cancer risk. *Nat. Genet.* **45**, 353–361 (2013).
10. Pharoah, P. D. P. *et al.* GWAS meta-analysis and replication identifies three new susceptibility loci for ovarian cancer. *Nat. Genet.* **45**, 362–370, 370e1–2 (2013).
11. Eeles, R. A. *et al.* Identification of 23 new prostate cancer susceptibility loci using the iCOGS custom genotyping array. *Nat. Genet.* **45**, 385–391, 391e1–2 (2013).
12. Teo, Y. Y. *et al.* A genotype calling algorithm for the Illumina BeadArray platform. *Bioinforma. Oxf. Engl.* **23**, 2741–2746 (2007).
13. Price, A. L. *et al.* Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet.* **38**, 904–909 (2006).

14. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
15. Aulchenko, Y. S., Ripke, S., Isaacs, A. & van Duijn, C. M. GenABEL: an R library for genome-wide association analysis. *Bioinforma. Oxf. Engl.* **23**, 1294–1296 (2007).
16. Clayton, D. & Leung, H.-T. An R package for analysis of whole-genome association studies. *Hum. Hered.* **64**, 45–51 (2007).
17. Howie, B., Marchini, J. & Stephens, M. Genotype imputation with thousands of genomes. *G3 Bethesda Md* **1**, 457–470 (2011).
18. Delaneau, O., Zagury, J.-F. & Marchini, J. Improved whole-chromosome phasing for disease and population genetic studies. *Nat. Methods* **10**, 5–6 (2013).
19. Timpson, N. J. *et al.* A rare variant in APOC3 is associated with plasma triglyceride and VLDL levels in Europeans. *Nat. Commun.* **5**, 4871 (2014).
20. Marchini, J., Howie, B., Myers, S., McVean, G. & Donnelly, P. A new multipoint method for genome-wide association studies by imputation of genotypes. *Nat. Genet.* **39**, 906–913 (2007).
21. Pruim, R. J. *et al.* LocusZoom: regional visualization of genome-wide association scan results. *Bioinforma. Oxf. Engl.* **26**, 2336–2337 (2010).
22. Mägi, R. & Morris, A. P. GWAMA: software for genome-wide association meta-analysis. *BMC Bioinformatics* **11**, 288 (2010).
23. Higgins, J. P. T. & Thompson, S. G. Quantifying heterogeneity in a meta-analysis. *Stat. Med.* **21**, 1539–1558 (2002).
24. Huedo-Medina, T. B., Sánchez-Meca, J., Marín-Martínez, F. & Botella, J. Assessing heterogeneity in meta-analysis: Q statistic or I<sup>2</sup> index? *Psychol. Methods* **11**, 193–206 (2006).
25. Ward, L. D. & Kellis, M. HaploReg: a resource for exploring chromatin states, conservation, and regulatory motif alterations within sets of genetically linked variants. *Nucleic Acids Res.* **40**, D930–934 (2012).

26. Boyle, A. P. *et al.* Annotation of functional variation in personal genomes using RegulomeDB. *Genome Res.* **22**, 1790–1797 (2012).
27. ENCODE Project Consortium *et al.* Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project. *Nature* **447**, 799–816 (2007).
28. Hnisz, D. *et al.* Super-enhancers in the control of cell identity and disease. *Cell* **155**, 934–947 (2013).
29. Corradin, O. *et al.* Combinatorial effects of multiple enhancer variants in linkage disequilibrium dictate levels of gene expression to confer susceptibility to common traits. *Genome Res.* **24**, 1–13 (2014).
30. GTEx Consortium. The Genotype-Tissue Expression (GTEx) project. *Nat. Genet.* **45**, 580–585 (2013).
31. Livak, K. J. & Schmittgen, T. D. Analysis of relative gene expression data using real-time quantitative PCR and the 2(-Delta Delta C(T)) Method. *Methods San Diego Calif* **25**, 402–408 (2001).
32. Giresi, P. G., Kim, J., McDaniell, R. M., Iyer, V. R. & Lieb, J. D. FAIRE (Formaldehyde-Assisted Isolation of Regulatory Elements) isolates active regulatory elements from human chromatin. *Genome Res.* **17**, 877–885 (2007).
33. Ghoussaini, M. *et al.* Evidence that breast cancer risk at the 2q35 locus is mediated through IGFBP5 regulation. *Nat. Commun.* **4**, 4999 (2014).
34. Lewis, A. *et al.* A polymorphic enhancer near GREM1 influences bowel cancer risk through differential CDX2 and TCF7L2 binding. *Cell Rep.* **8**, 983–990 (2014).

## SUPPLEMENTARY NOTE

### Detailed Description of the Case and Control Sample Sets

A summary of the studies included in the GWAS and both replication phases is shown in **Supplementary Table 1**, with additional details provided below. **Supplementary Figure 1** provides a flow diagram of the overall study design. All studies were of women of European ancestry. All studies have the relevant IRB approval in each country in accordance with the principles embodied in the Declaration of Helsinki, and informed consent was obtained from all participants. A total of 7,737 cases and 37,144 controls were included in this analysis. Cases and controls were matched based on geographical location and case-control clustering in principal components analysis (PCA) (**Supplementary Figure 1**).

#### ***EC and control GWAS Sample Sets:***

Quality control (QC) was applied to all GWAS sets, following standard QC approaches detailed in Spurdle et al<sup>1</sup>. Also see online methods.

#### **NSECG**

National Study of the Genetics of Endometrial Cancer (NSECG) cases were identified from collaborating clinicians throughout the UK from 2008 to 2013, taking care not to recruit from centres involved in SEARCH. Inclusion criteria were adenocarcinomas of the uterus presenting at 70 years of age or younger. Almost all cases were incident and sampled within 6 months of diagnosis. Peripheral blood was collected from each participant and DNA extracted using standard methods and the participants completed the associated questionnaire. Tumour histology was confirmed from routine hospital reports and further details of histopathology and other tumor pathology characteristic was abstracted from these clinical pathology reports. 925 samples were genotyped using the Illumina 660W Quads in the GWAS scan, 965 samples were genotyped in the phase 1 replication using iCOGS arrays, and a further 1195 were genotyped using KASPar and Fluidigm genotyping for the second replication phase. There was no overlap in samples used and all cases were of European ancestry.

#### **ANECs**

The Australian National Endometrial Cancer Study (ANECs) is an Australian population-based case-control family study of cancer of the uterine corpus<sup>2</sup>. Women aged 18-79, newly diagnosed with histologically confirmed primary cancer of the endometrium between July 2005 and December 2007 were identified through major hospitals nationally, and also from state-based cancer registries. Excluding women who could not be contacted (mostly due to death, illness or failure to contact), case participation rate was 63%. Participants completed a detailed questionnaire providing clinical and epidemiological information, including ethnicity of all four grandparents. Information on tumor pathology characteristics was abstracted in standardized format from clinical pathology reports for all patients. 606 ANECs samples all of endometrioid-only histology were used for the original EC GWAS and a further 538 were genotyped using iCOGS for the first replication phase.

#### **SEARCH**

The Studies of Epidemiology and Risk factors in Cancer Heredity (SEARCH) is an ongoing population-based study with cases ascertained through the Eastern Cancer Registration and Information Centre (<http://www.ecric.org.uk>). All women diagnosed with endometrial cancer between the ages of 18-69 years (average age of diagnosis 58 years) from August 2001 to September 2007 were eligible for inclusion. Approximately 54% of eligible patients have enrolled in the study. Women taking part in the study were asked to provide a 20ml blood

sample for DNA analysis, and to complete a comprehensive epidemiological questionnaire. Controls were also drawn from SEARCH (<http://ccge.medschl.cam.ac.uk/search/>), but had no prior history of cancer at the time of recruitment. They were female, also between the ages of 18-69 at the time of recruitment and matched to cases in geographical profile. Approximately 35% of eligible controls enrolled in the study. All participants reported Caucasian ethnicity. Information on tumor pathology characteristics was provided by the Eastern Cancer Registration and Information Centre and was derived from clinical pathology reports for all patients. 681 SEARCH samples with endometroid-only histology were used in the original GWAS and a further 773 non-overlapping cases were used in the iCOGS analysis.

### **UK1/CORGI**

The UK1 Colorectal Tumour Gene Identification (CoRGI) is a GWAS for colorectal neoplasia<sup>3</sup>. The 894 controls matched with the NSECG cases were spouses or partners unaffected by cancer and without a personal family history (to second degree relative level) of colorectal neoplasia. Known dominant polyposis syndromes, HNPCC/Lynch syndrome or bi-allelic MUTYH mutation carriers were excluded. All cases and controls were of white UK ethnic origin. Genotyping was done on the Illumina Hap550 arrays.

### **Scotland Phase 1**

Scotland Phase1 is a colorectal cancer GWAS<sup>4</sup> with 1012 cancer-free population controls. Known dominant polyposis syndromes, HNPCC/Lynch syndrome or bi-allelic MUTYH mutation carriers were excluded. Control subjects were sampled from the Scottish population NHS registers, matched by age ( $\pm 5$  years), gender and area of residence within Scotland. A subset of 392 controls from this dataset were matched with the NSECG GWAS cases and these were chosen based on case-control clustering on PCA. Genotyping was done on the Illumina Hap300 and Hap 240S arrays.

### **QIMR**

The Queensland Institute of Medical Research (QIMR) control sample is a subsection of subjects recruited as part of the Brisbane Adolescent Twin Study<sup>5,6</sup>. Twins were recruited from schools in Brisbane, Australia and surrounding areas of southeast Queensland and were examined close to their 12th birthday. Blood was obtained from all twins and most parents. Parents were asked the ancestry of all eight great-grandparents of the twins. More than 95% of great-grandparents were identified as being of northern European ancestry, mainly from Britain and Ireland. This analysis used genotype data from parents and siblings only, extracted from an existing Illumina 610K BeadChip genome-wide association scan<sup>7</sup> and recalled using the Illuminus algorithm. After QC, 1846 QIMR controls were available for inclusion in the analysis.

### **HCS**

The Hunter Community Study (HCS) is a population-based cohort study consisting of men and women aged 55-85 years of age in Newcastle, New South Wales, Australia<sup>8</sup>. Participants were randomly selected from the NSW State electoral roll (listing on the electoral roll is compulsory in Australia) and contacted between December 2004 and December 2007. Non-English speaking persons and those living in a residential aged-care facility were ineligible for participation in the study. Participants were asked to complete five self-report questionnaires as well as attend the HCS data collection centre so clinical measures could be obtained. In total, 44.5% of eligible controls agreed to participate in this study. Genotype data for this study were extracted from an existing Illumina 610K BeadChip genome-wide association study scan and recalled using the Illuminus algorithm. After QC, 1237 HCS controls were available for inclusion in the analysis.

### **WTCCC**

Controls utilized were genotyped as part of the Wellcome Trust Case Control Consortium (WTCCC2)<sup>9</sup>. These controls are drawn from two sources: 2,674 controls from the 1958 Birth Cohort (1958BC), a population-based study in the United Kingdom of individuals born in 1 week in 1958<sup>10</sup>; and 2,501 controls identified through the UK National Blood Service (NBS)<sup>9</sup>. 1958BC controls were matched with NSECG cases and the NBS controls were matched with SEARCH cases.

### ***Phase 1 iCOGS replication Case Sample Sets:***

All samples in the first replication phase were genotyped as part of the Collaborative Oncological Gene-environment Study (iCOGS) initiative on a custom Illumina Infinium iSelect array. Cases from ANECS and SEARCH and NSECG were recruited as detailed above, and are non-overlapping.

### **BECS**

The Bavarian Endometrial Cancer Cases and Controls Study (BECS) is a single-center case-control study, conducted between 2002 and 2008, with the aim of investigating genetic and epidemiological risk factors for endometrial cancer. Cases were either incident cases referred to the University Hospital Erlangen by surrounding practitioners (66% of the case sample set), or prevalent cases that were outpatients in follow-up care approached within 6.2 ( $\pm 4.6$  SD) years after treatment for primary endometrial cancer in the same hospital (34% of the case sample set). Epidemiological information was collected by a structured questionnaire completed during an interview and clinical data for the cases was obtained from clinical health records.

### **CAHRES**

Details of the population selection process have been published previously for the Cancer Hormone Replacement Epidemiology Study (CAHRES)<sup>11</sup>. Formerly known as the Singapore and Sweden Breast/Endometrial Cancer Study (SASBAC), this population based case-control study was conducted among Swedish women aged 50-74 years, who were residing in Sweden between January 1<sup>st</sup> 1994 and December 31<sup>st</sup> 1995. Endometrial cancer cases were identified through the nation-wide cancer registries in Sweden. All participants provided detailed questionnaire information. For endometrial cancer, histological specimens were reviewed and re-classified by the study pathologist. All participants reported Caucasian ethnicity.

### **HJECS**

The Hannover-Jena Endometrial Cancer Study (HJECS), a hospital-based case-control study, included 250 German women, aged 31-89 years, who were recruited either at the Friedrich Schiller University of Jena or at Hannover Medical School after having been diagnosed with histologically confirmed primary incident endometrial carcinoma between 2004 and 2010. Epidemiological data were obtained from questionnaires, and information on tumor stage and histology was obtained from pathology and clinical reports. Over 98% were of German descent. Interviews were conducted at either the Friedrich Schiller University of Jena or at Hannover Medical School, and peripheral blood was collected for the extraction of DNA from white blood cells.

### **LES**

The Leuven Endometrial Study (LES) is a hospital based case-control study. Eligible cases, identified by active surveillance of electronic patient files at the Leuven University Hospital, were white women aged 27-80 years diagnosed with endometrial cancer. Clinical data for endometrial cancer patients were recorded during interview at the time of diagnosis, and from pathology reports. All medical records were reviewed by trained abstractors and pathology reports compatible with primary, invasive, epithelial endometrial adenocarcinoma



of all stages (I –IV) and all grades were consulted. Participation rates exceeded 95% for cases.

### **MECS**

The Mayo Endometrial Cancer Study (MECS) includes a clinic-based prospective collection of primary endometrial cases diagnosed from 2008 to 2011 and seen at Mayo Clinic Rochester with primary endometrial cancer diagnosed at age 18 and older. DNA was isolated from white blood cells using qiagen isolation kit. DNA concentration was measured with picogreen. Clinical data were abstracted from electronic medical records and supplemented with a risk factor questionnaire. Control data were obtained from Mayo Clinic OCAC controls (MAY) and BCAC controls (MCBCS).

### **MoMaTEC**

Molecular Markers in Treatment of Endometrial Cancer (MoMaTEC) cases were recruited from an unselected patient population primarily treated for endometrial carcinoma at Haukeland University Hospital, Bergen during 2001-2009. This is the referral hospital for Hordaland county; the area is demographically well defined, with about 450,000 inhabitants, representing approximately 10% of the Norwegian population and with a similar incidence rate and prognosis as the total Norwegian population of endometrial cancers<sup>12–14</sup>. Clinical Information for cases regarding age, FIGO stage, histologic subtype, grade and prognosis was extracted from medical records. DNA was extracted from peripheral blood samples.

### **NECS**

The Newcastle Endometrial Cancer Study (NECS) includes histologically confirmed endometrial cancer cases consecutively recruited from 1992 up to 2005 at the Hunter Centre for Gynaecological Cancer, John Hunter Hospital, Newcastle, New South Wales, Australia<sup>15</sup>. The final analysis included 194 endometrial cancer patients. Data on reproductive and environmental risk factors including ethnicity was collected using self reported questionnaires. Information regarding recurrence, stage, grade and histology of endometrial cancer was collected from medical records. Patients presenting at this hospital-based site were captured by ANECS recruitment from 2005 onwards.

### **RENDOCAS**

The Registry of Endometrial Cancer in Sweden (RENDOCAS) is a hospital based case-control study. Patients (n=520) who underwent surgery for endometrial cancer at Karolinska University hospital Solna, Sweden between 2008 and 2011 were included in the study. For each patient, the following was collected: blood and tumor samples; detailed family history and formulation of a pedigree where all suspected cancer cases were verified in medical records/pathology report if possible; questionnaire covering relevant environmental factors underlying endometrial cancer.

### ***iCOGS Control Sample Sets:***

As indicated in **Supplementary Table 1**, iCOGS endometrial cancer case sample sets were matched with controls from the same countries and also clustered with cases in PCA (**Supplementary Figure 2**). Controls were genotyped using the same iCOGS array and data were largely drawn from healthy controls participating in the Breast Cancer Association Consortium (BCAC)<sup>16</sup> part of the iCOGS initiative. Additional controls were from the Mayo Clinic via the Ovarian Cancer Association Consortium (OCAC)<sup>17</sup>, and Norwegian female controls recruited in Bergen for use in the MoMaTEC case-control genotyping studies.

### ***Endometrioid and non-endometrioid histology analysis***

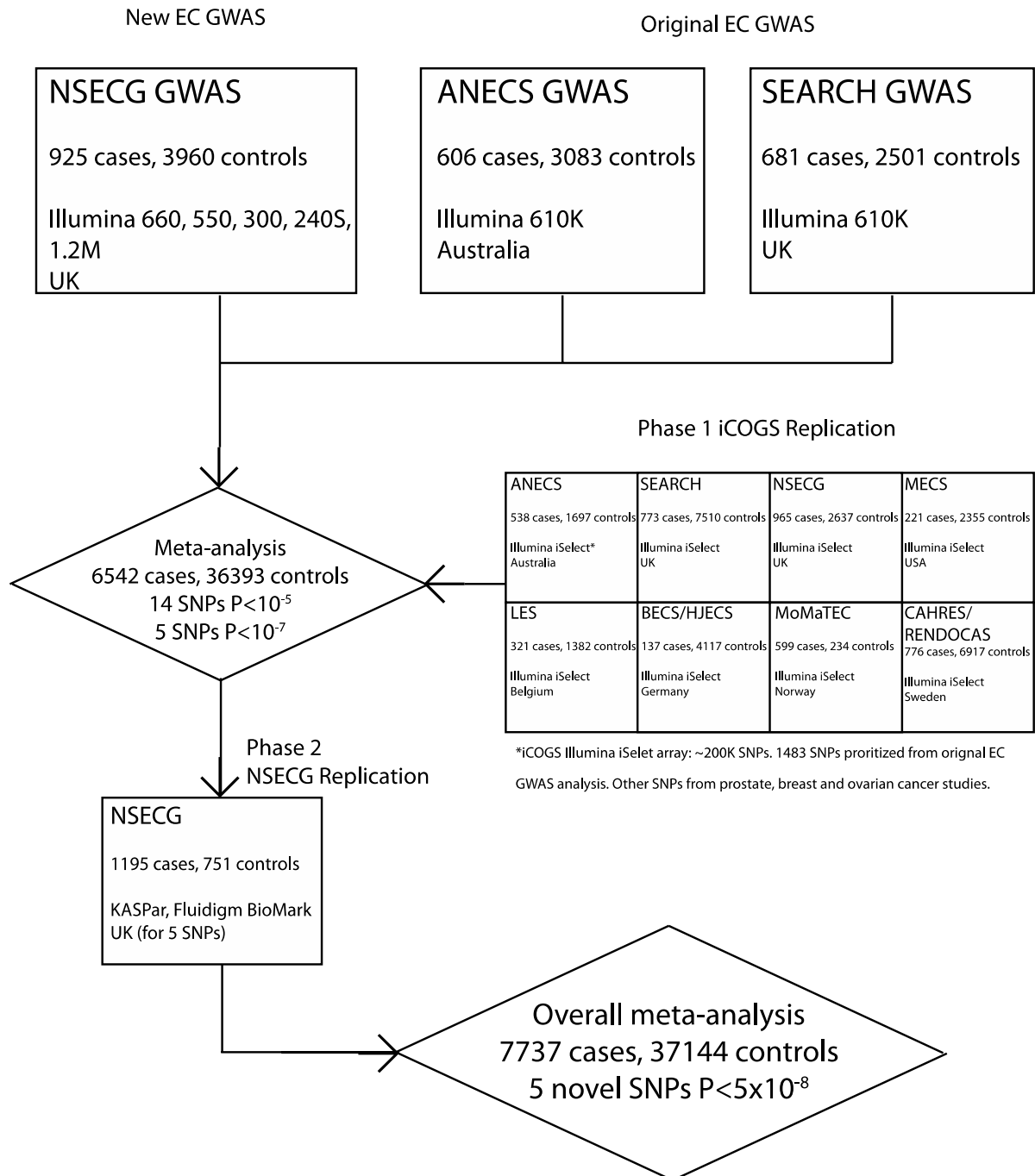
Cases were defined as having endometrioid subtype based on pathology report of endometrioid histology only. Non-endometrioid subtypes included carcinosarcoma, clear cell, serous, mucinous, and tumours of mixed histology (any combination). 6,635 (86%) of the 7,737 EC cases displayed endometrioid-only histology and association testing and meta-analysis was also conducted using endometrioid-only histology cases. The results of this analysis for the novel risk loci are shown in **Table 1**. Endometrioid-only phase 1 meta-analysis (n=5,590) found only novel risk loci that were identified in the all histologies analysis (**Table 1**, **Supplementary Table 2**). Analysis of 952 EC cases that displayed non-endometrioid histology found no SNPs near genome-wide significance and this is expected given the limited statistical power.

### ***Ishikawa and ECC-1 cells***

Results from Ishikawa and ECC-1 cells are listed separately in publicly available ENCODE<sup>18</sup> tier 3 data but STR-profiling has shown that these two cell lines are very similar<sup>19</sup>. Based on the results presented by Korch *et al.*, the International Cell Line Authentication Committee (ICLAC) recommended in the 2013 Database of Cross-Contaminated or Misidentified Cell Lines that ECC-1 be re-identified as Ishikawa cells. Our *in vitro* functional analysis for the 13q22 locus made use of our supply of Ishikawa cells for FAIRE and ChIP experiments and ECC-1 cells for luciferase reporter assays. Both cell lines displayed identical genotypes for rs9600103 and rs11841589, similar *KLF5* expression levels, and 20x sequencing using the Ion AmpliSeq™ Comprehensive Cancer Panel confirmed that variants in Ishikawa and ECC-1 are 90% concordant (based on 3,004 exonic SNVs in 409 cancer-related genes). In line with these findings and ICLAC recommendations, we have presented functional work on these cells as Ishikawa cells.

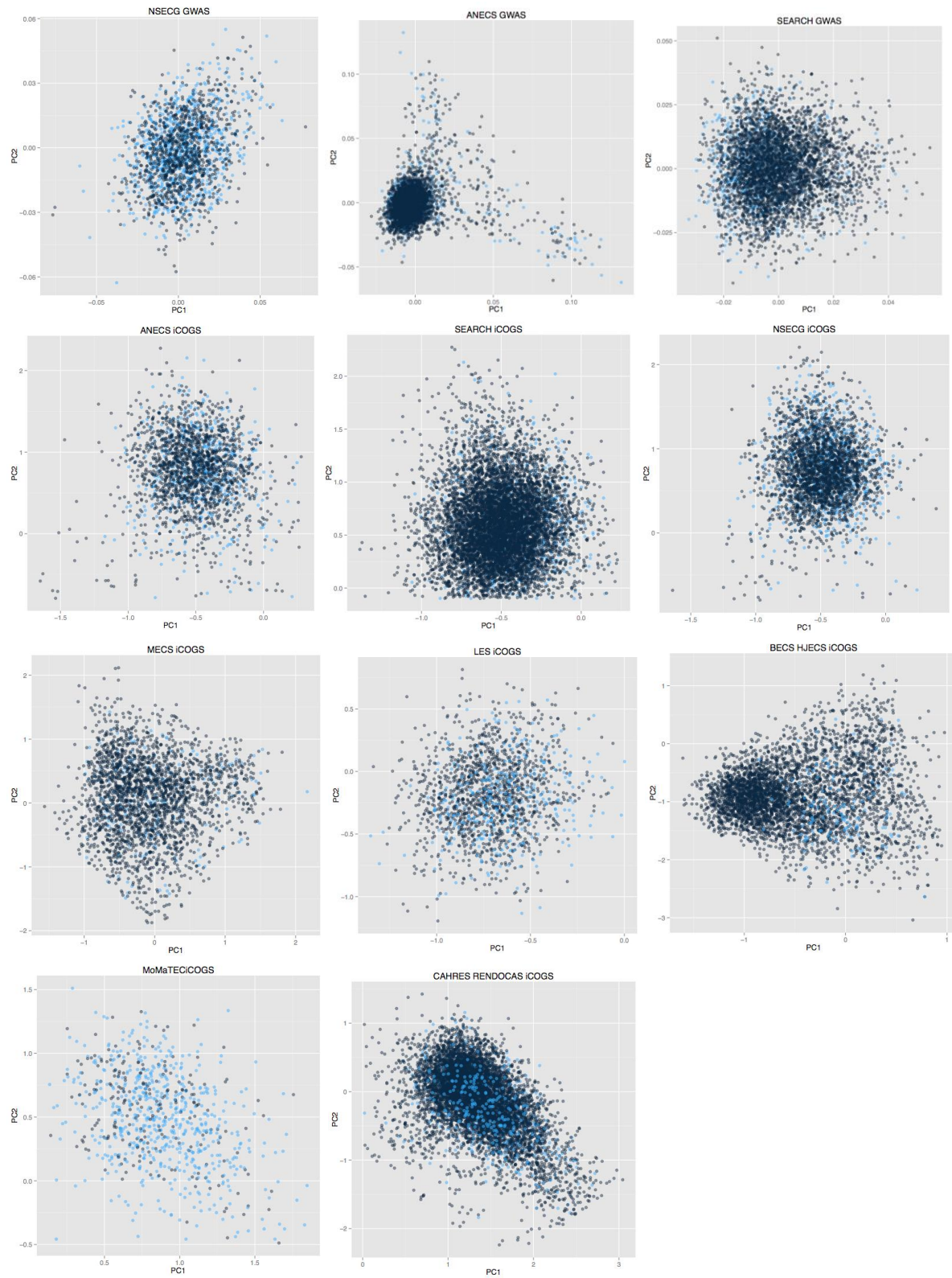
### Supplementary Figure 1: EC meta-analysis flow diagram.

This schematic figure illustrates the EC meta-analysis study design. The new NSECG GWAS was meta-analysed with a re-analysis of the original EC GWAS (ANECS and SEARCH) and phase 1 iCOGS replication (eight groups). This meta-analysis of 6,542 cases and 36,393 controls yielded 14 regions with SNPs  $P < 10^{-5}$ , of which five regions had SNPs  $P < 10^{-7}$ . These five SNPs were brought forward to the phase 2 NSECG replication and were confirmed as novel genome-wide significant risk loci ( $P < 5 \times 10^{-8}$ ) in the overall meta-analysis of 7,737 cases and 37,144 controls.



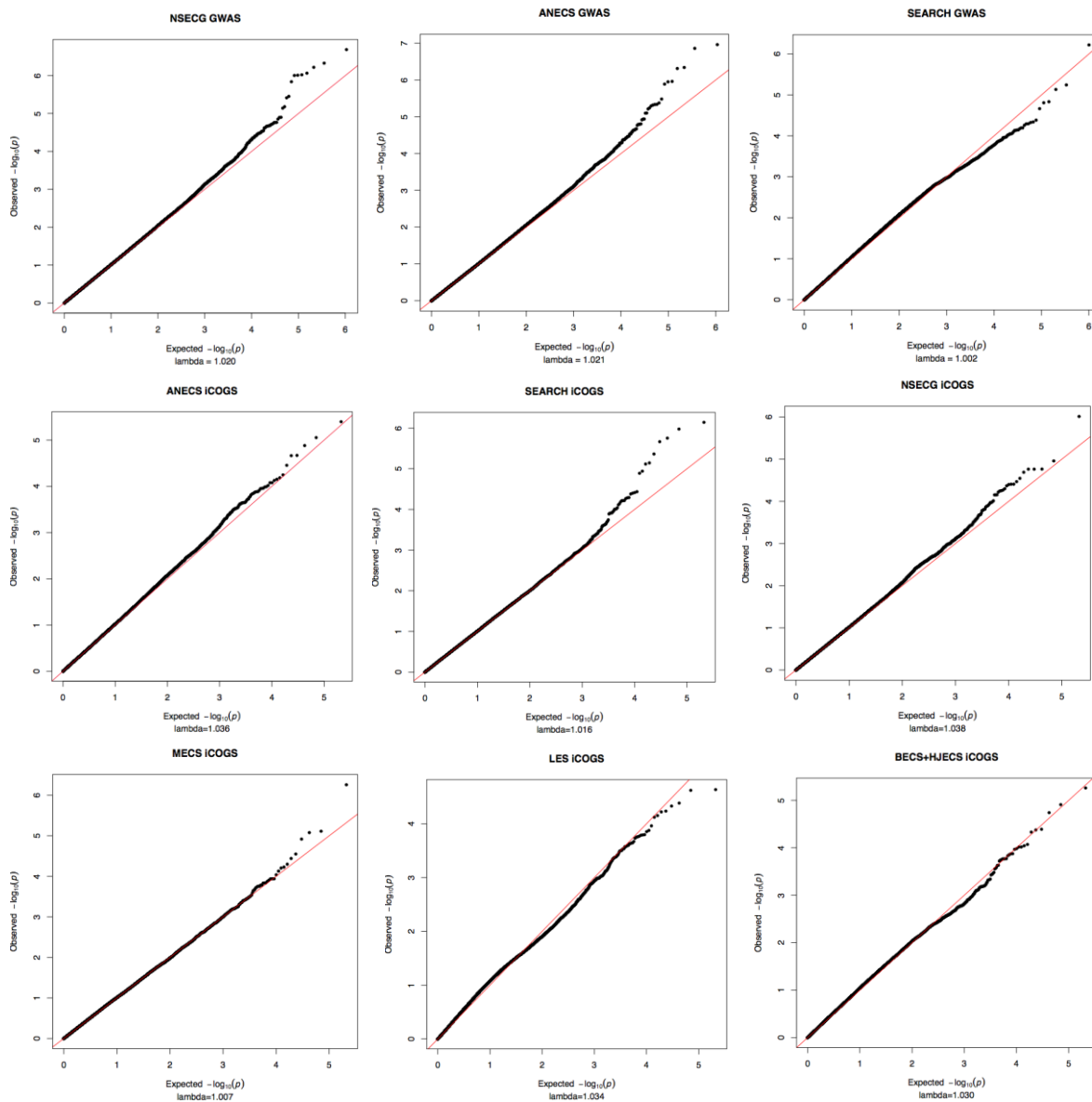
## Supplementary Figure 2: Principal Components Analysis (PCA) of three GWAS and eight iCOGS studies.

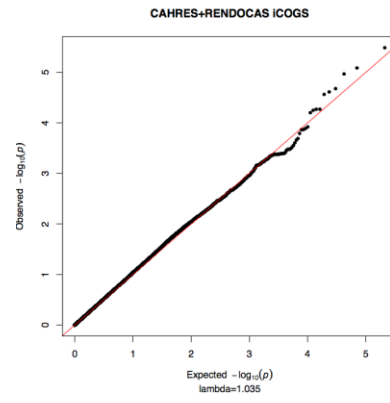
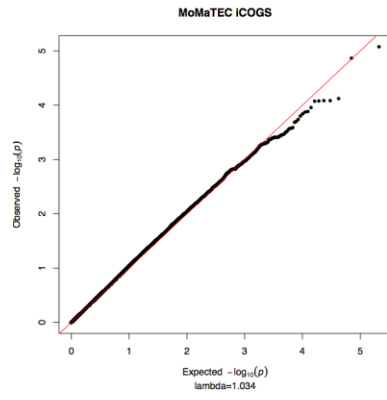
Plots of the first two principal components (PCs) in each study. EC cases are represented by blue dots, whereas controls are in black. Samples were excluded if they clustered away from the centroid in the first four PCs and these are plots of the samples used in the analysis.



### Supplementary Figure 3: Quantile-quantile plots of the ranked trend test statistics for three GWAS and eight iCOGS groups.

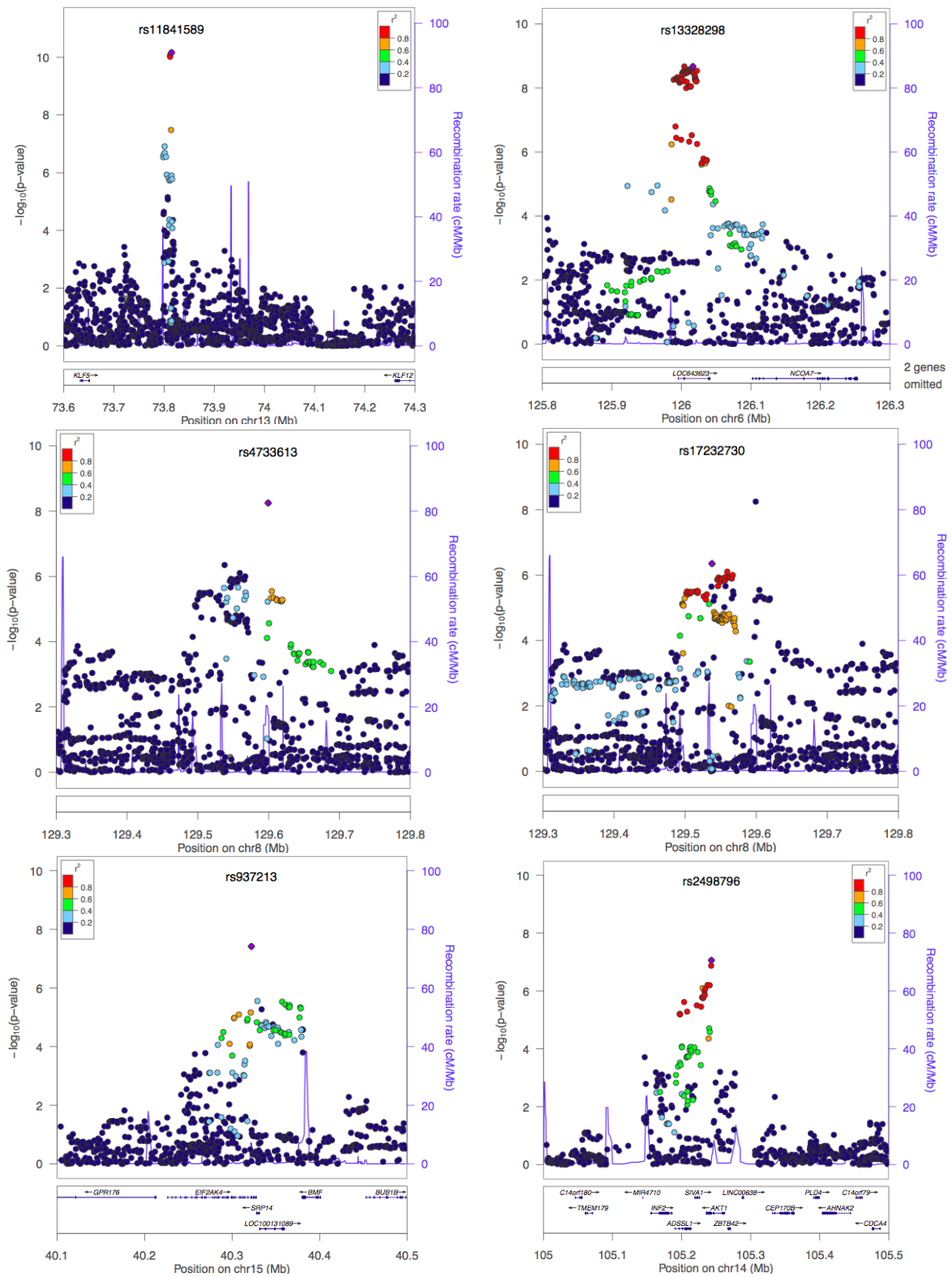
The  $-\log_{10}$  transformed observed P-values (y-axis) were plotted against the expected P-values under the null hypothesis (x-axis). The red line denotes the expectation under no deviation from the null hypothesis. The QQ-plots show little evidence of genomic inflation and the  $\lambda_{GC}$  for each study are: NSECG GWAS 1.020, ANECS GWAS 1.021, SEARCH GWAS 1.002, ANECS iCOGS 1.036, SEARCH iCOGS 1.016, NSECG iCOGS 1.038, MECS iCOGS 1.007, LES iCOGS 1.034, BECS iCOGS 1.030, MoMaTEC iCOGS 1.034, CAHRES RENDOCAS iCOGS 1.037. For the three GWAS, all genotyped SNPs passing QC are displayed. For iCOGS, 105,000 SNPs after LD-pruning ( $r^2 < 0.2$ ) and  $> 500\text{kb}$  from the 1,483 EC prioritized SNPs on the iCOGS are displayed.





**Supplementary Figure 4: Regional association plots for the five novel loci associated with endometrial cancer.**

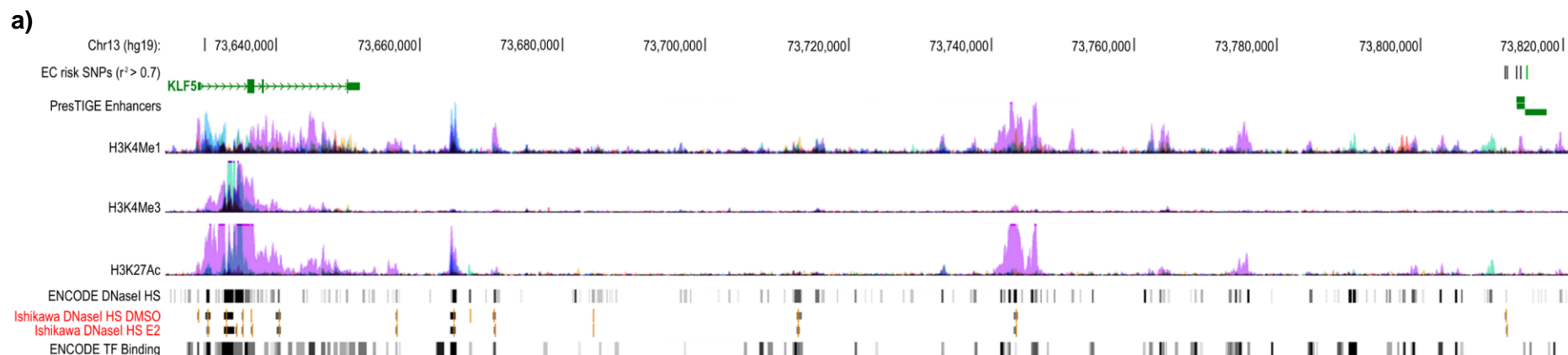
The  $-\log_{10} P$  values from the meta-analysis and regional imputation for NSECG, ANECS, SEARCH, and eight iCOGS groups are shown for SNPs at: a) 13q22.1, b) 6q22, c) & d) 8q24, e) 15q15 and f) 14q32.33. The SNP with the lowest P value at each locus is labelled and marked as a purple diamond, and the dot color represents the LD with the top SNP. The blue line shows recombination rates in cM/Mb. Compared with **Figure 2**, more SNPs are displayed in these plots. SNPs with info scores of more than 0.6 in iCOGS and more than 0.9 in NSECG, ANECS, and SEARCH are included.



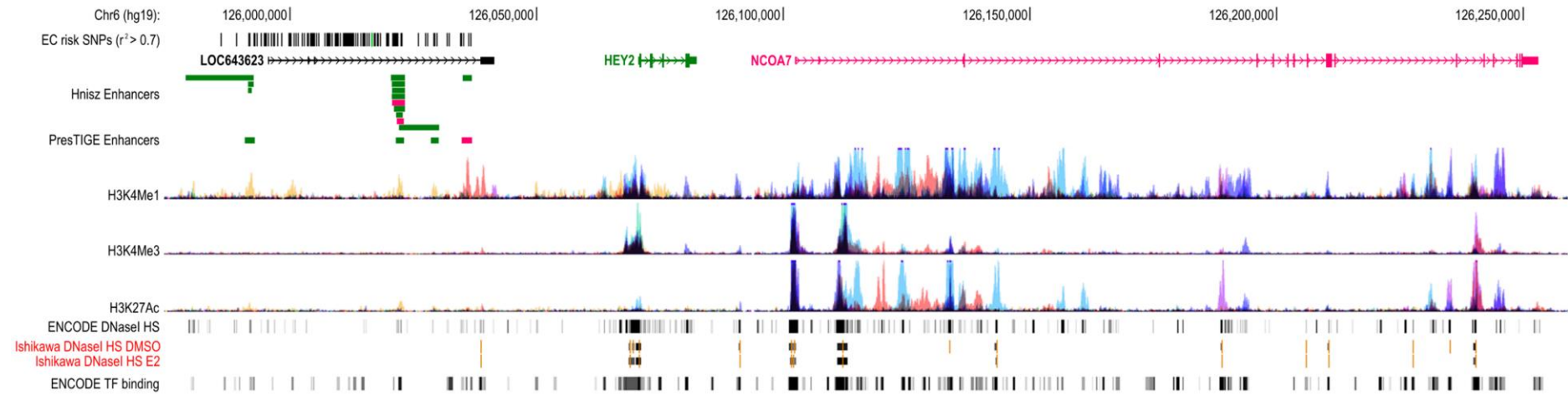


### Supplementary Figure 5: Genetic landscape of novel endometrial cancer associated regions.

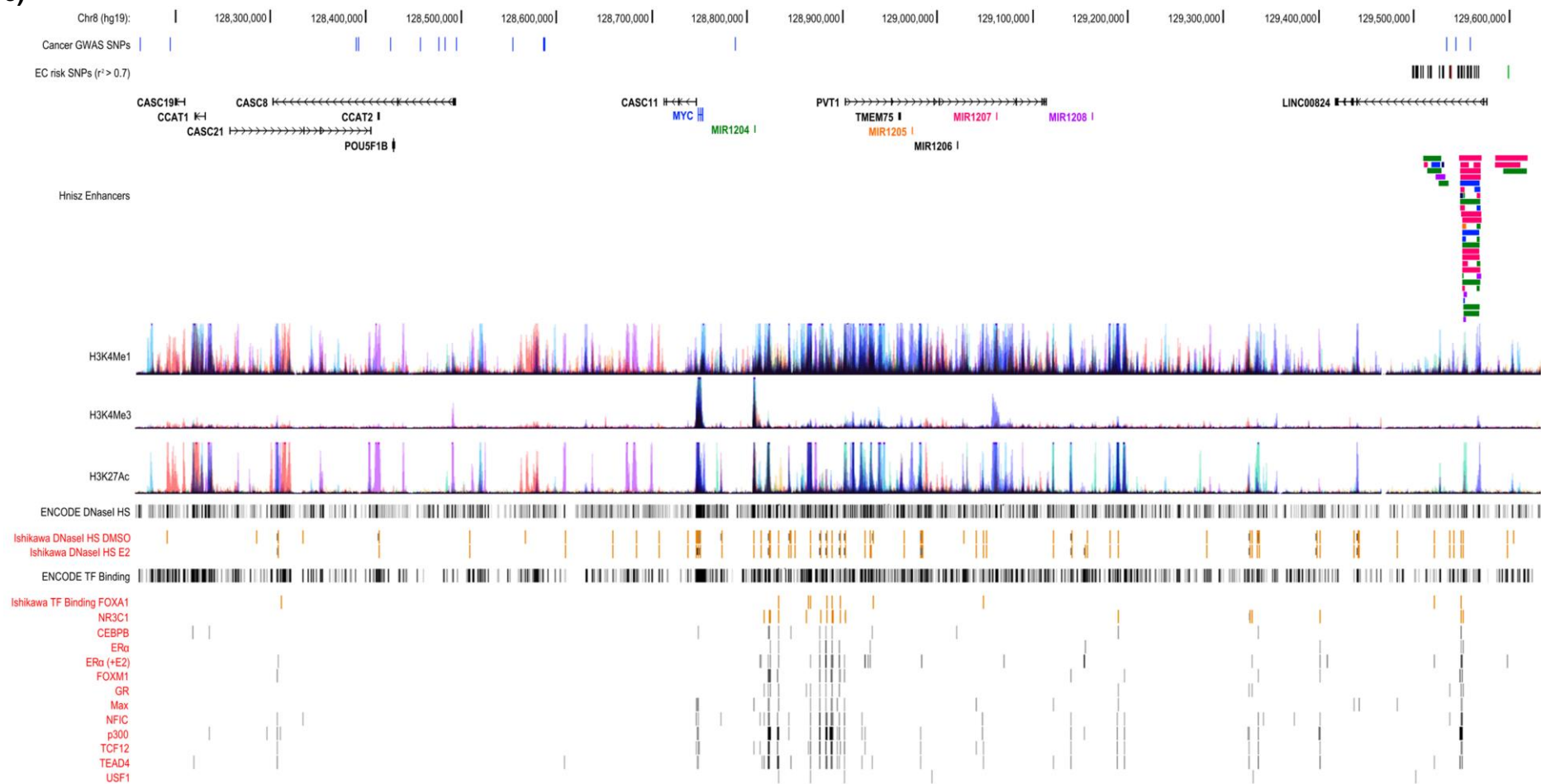
Plots for novel risk loci at a) 13q22.1, b) 6q22.31, c) 8q24.21, d) 15q15 and e) 14q32. SNPs in strong LD ( $r^2 > 0.7$ ) with the lead EC-risk SNP have been plotted for each region and the lead SNP denoted in green. The second, independently associated SNP found at 8q24.21 after conditioning on the lead SNP is denoted in red (c). Previously reported cancer risk SNPs identified by GWAS at 8q24 are shown in blue (c), none of which are in LD ( $r^2 \leq 0.02$ ) with EC risk SNPs. Likely enhancers identified by Hnisz *et al.*<sup>20</sup> and PresTIGE<sup>21</sup> that overlap EC-risk associated SNPs are depicted as colored bars, where the color of the likely enhancer matches the schematic of its predicted target gene, as determined by correlations with gene expression. As described in Online Methods Hnisz *et al.* predicted enhancer-gene interactions by identifying 'super-enhancers' (regions containing neighboring H3K27Ac modifications) from 86 cell and tissue types and then the expressed transcript with transcription start site closest to the centre of the super-enhancer was assigned as the target gene. PresTIGE pairs cell-type specific H3K4Me1 and gene expression data from 13 cell types to identify likely enhancer-gene interactions. Additional tracks include: Histone modifications associated with promoters (H3K4Me3) and enhancers (H3K4Me1 and H3K27Ac) from seven ENCODE Project cell types; DNaseI hypersensitivity sites (DHS) and transcription factor (TF) binding identified in 125 and 91 ENCODE Project cell types, respectively; DHS identified in Ishikawa endometrial cancer cells using DMSO vehicle and under estrogen (E2) stimulation are shown; transcription factor binding regions in Ishikawa cells that encompass EC-risk SNP loci are also displayed. For all risk loci, EC-risk associated SNPs co-locate with at least one enhancer predicted by cell-type specific analysis, implicating the following genes/transcripts as worthy of investigation: a) *KLF5*; b) *HEY2*, *NCOA7*; c) *MYC*, *MIR1204*, *MIR1205*, *MIR1207*, *MIR1208*; d) *BMF*, *GPR176*, *SRP14*, *LOC100131089*; e) *AKT*, *ADSSL1*, *INF2*, *ZBTB42*, *SIVA1*. For four loci (a, c, d and e), likely enhancers overlap with at least one region displaying evidence of regulatory activity (DHS and/or TF binding) in the single endometrial cancer cell line (Ishikawa) assayed by ENCODE.



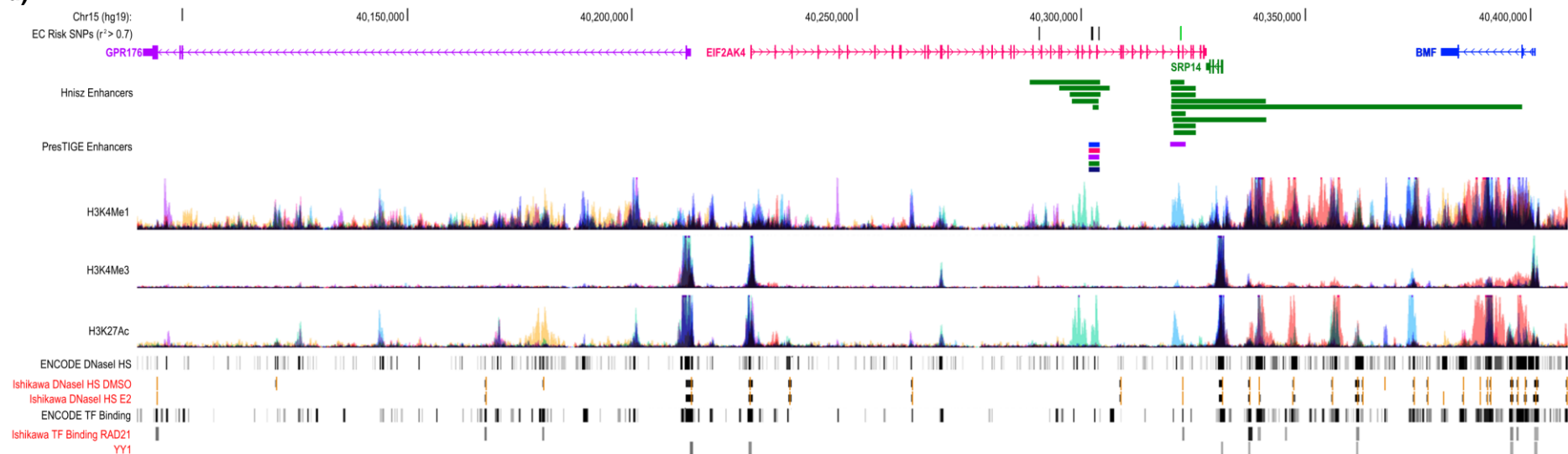
b)

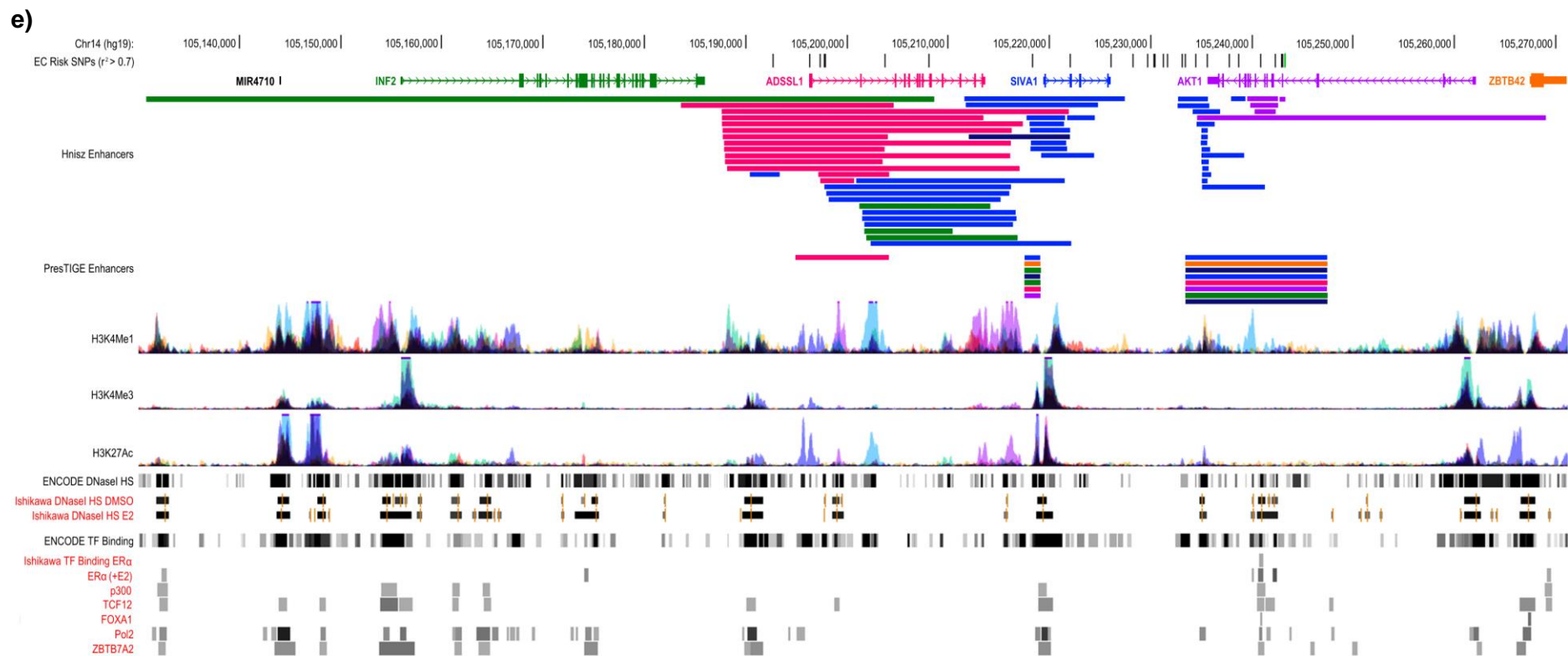


c)

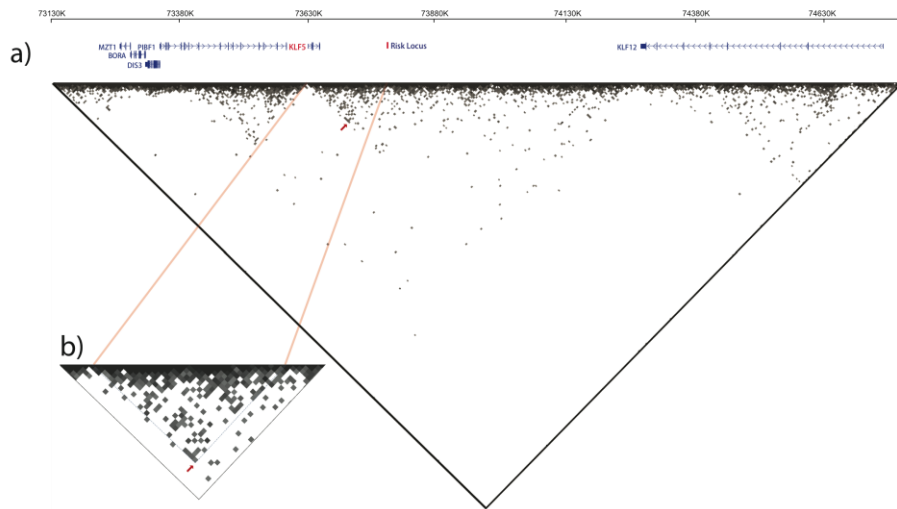


d)



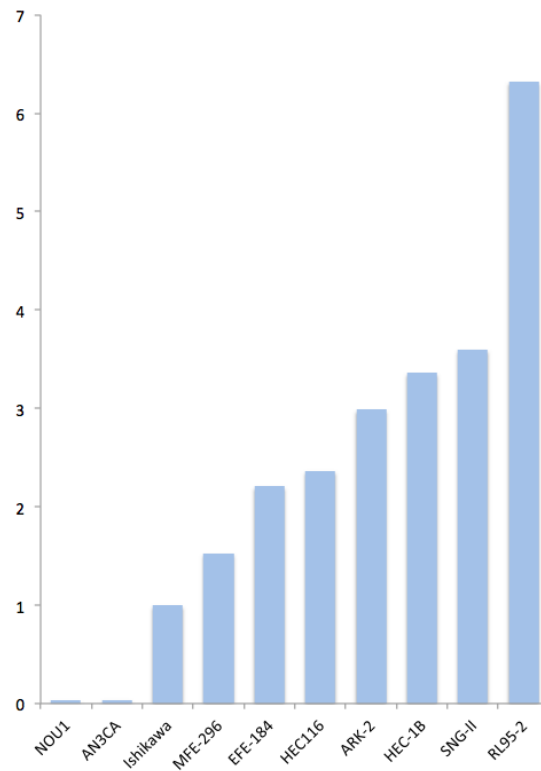


**Supplementary Figure 6: Hi-C chromatin capture of 13q22 locus in HeLa S3 cells.** a) 5Kb KR normalized contact matrix in Hi-C experiment for HeLa S3 cells was used to represent the interaction pattern between *KLF5* and risk locus rs11841589/rs9600103<sup>22</sup>. A loop was anchored at the *KLF5* promoter and the risk locus (see b) for the small topologically associated domain and the red arrow for loop anchor), which indicated distal *cis*-regulatory element with in the risk locus. The interaction between the rs11841589/rs9600103 risk locus with the *KLF5* promoter was the strongest interaction observed out of the 262 protein-coding genes on chromosome 13 (P=0.004). The color scheme in the contact matrix is KR normalized score with the black indicating a strong interaction.



### Supplementary Figure 7: Quantification of *KLF5* expression in EC cell lines

Expression of *KLF5* in 11 EC cell lines as described in **Online methods** using qRT-PCR, expression levels on the x-axis are relative to *KLF5* expression in Ishikawa cells using the ddCT method.



## References

1. Spurdle, A. B. *et al.* Genome-wide association study identifies a common variant associated with risk of endometrial cancer. *Nat. Genet.* **43**, 451–454 (2011).
2. Spurdle, A., Webb, P. & Australian National Endometrial Cancer Study. Re: Excess of early onset multiple myeloma in endometrial cancer probands and their relatives suggests common susceptibility. *Gynecol. Oncol.* **109**, 153; author reply 154 (2008).
3. Tomlinson, I. *et al.* A genome-wide association scan of tag SNPs identifies a susceptibility variant for colorectal cancer at 8q24.21. *Nat. Genet.* **39**, 984–988 (2007).
4. Tenesa, A. *et al.* Genome-wide association scan identifies a colorectal cancer susceptibility locus on 11q23 and replicates risk loci at 8q24 and 18q21. *Nat. Genet.* **40**, 631–637 (2008).
5. McGregor, B. *et al.* Genetic and environmental contributions to size, color, shape, and other characteristics of melanocytic naevi in a sample of adolescent twins. *Genet. Epidemiol.* **16**, 40–53 (1999).
6. Zhu, G. *et al.* A major quantitative-trait locus for mole density is linked to the familial melanoma gene CDKN2A: a maximum-likelihood combined linkage and association analysis in twins and their sibs. *Am. J. Hum. Genet.* **65**, 483–492 (1999).
7. Painter, J. N. *et al.* Genome-wide association study identifies a locus at 7p15.2 associated with endometriosis. *Nat. Genet.* **43**, 51–54 (2011).
8. McEvoy, M. *et al.* Cohort profile: The Hunter Community Study. *Int. J. Epidemiol.* **39**, 1452–1463 (2010).
9. Wellcome Trust Case Control Consortium. Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature* **447**, 661–678 (2007).
10. Power, C. & Elliott, J. Cohort profile: 1958 British birth cohort (National Child Development Study). *Int. J. Epidemiol.* **35**, 34–41 (2006).



11. Weiderpass, E. *et al.* Risk of endometrial cancer following estrogen replacement with and without progestins. *J. Natl. Cancer Inst.* **91**, 1131–1137 (1999).
12. Wik, E. *et al.* Deoxyribonucleic acid ploidy in endometrial carcinoma: a reproducible and valid prognostic marker in a routine diagnostic setting. *Am. J. Obstet. Gynecol.* **201**, 603.e1–7 (2009).
13. Salvesen, H. B., Iversen, O. E. & Akslen, L. A. Prognostic significance of angiogenesis and Ki-67, p53, and p21 expression: a population-based endometrial carcinoma study. *J. Clin. Oncol. Off. J. Am. Soc. Clin. Oncol.* **17**, 1382–1390 (1999).
14. Salvesen, H. B. *et al.* Integrated genomic profiling of endometrial carcinoma associates aggressive tumors with indicators of PI3 kinase activation. *Proc. Natl. Acad. Sci. U. S. A.* **106**, 4834–4839 (2009).
15. Ashton, K. A. *et al.* The influence of the Cyclin D1 870 G>A polymorphism as an endometrial cancer risk factor. *BMC Cancer* **8**, 272 (2008).
16. Michailidou, K. *et al.* Large-scale genotyping identifies 41 new loci associated with breast cancer risk. *Nat. Genet.* **45**, 353–361 (2013).
17. Pharoah, P. D. P. *et al.* GWAS meta-analysis and replication identifies three new susceptibility loci for ovarian cancer. *Nat. Genet.* **45**, 362–370, 370e1–2 (2013).
18. ENCODE Project Consortium *et al.* Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project. *Nature* **447**, 799–816 (2007).
19. Korch, C. *et al.* DNA profiling analysis of endometrial and ovarian cell lines reveals misidentification, redundancy and contamination. *Gynecol. Oncol.* **127**, 241–248 (2012).
20. Hnisz, D. *et al.* Super-enhancers in the control of cell identity and disease. *Cell* **155**, 934–947 (2013).
21. Corradin, O. *et al.* Combinatorial effects of multiple enhancer variants in linkage disequilibrium dictate levels of gene expression to confer susceptibility to common traits. *Genome Res.* **24**, 1–13 (2014).
22. Rao, S. S. P. *et al.* A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell* **159**, 1665–1680 (2014).

## ECAC Study Collaborators

**The ANECS Group comprises:** AB Spurdle, PM Webb, J Young (QIMR Berghofer Medical Research Institute); Consumer representative: L McQuire; Clinical Collaborators: NSW: S Baron-Hay, D Bell, A Bonaventura, A Brand, S Braye, J Carter, F Chan, C Dalrymple, A Ferrier (deceased), G Gard, N Hacker, R Hogg, R Houghton, D Marsden, K McIlroy, G Otton, S Pather, A Proietto, G Robertson, J Scurry, R Sharma, G Wain, F Wong; Qld: J Armes, A Crandon, M Cummings, R Land, J Nicklin, L Perrin, A Obermair, B Ward; SA: M Davy, T Dodd, J Miller, M Oehler, S Paramasivum, J Pierides, F Whitehead; Tas: P Blomfield, D Challis; Vic: D Neesham, J Pyman, M Quinn, R Rome, M Weitzer; WA: B Brennan, I Hammond, Y Leung, A McCartney (deceased), C Stewart, J Thompson; Project Managers: S O'Brien, S Moore; Laboratory Manager: K Ferguson; Pathology Support: M Walsh; Admin Support: R Cicero, L Green, J Griffith, L Jackman, B Ranieri; Laboratory Assistants: M O'Brien, P Schultz; Research Nurses: B Alexander, C Baxter, H Croy, A Fitzgerald, E Herron, C Hill, M Jones, J Maidens, A Marshall, K Martin, J Mayhew, E Minehan, D Roffe, H Shirley, H Steane, A Stenlake, A Ward, S Webb, J White.

**CHIBCHA (study of hereditary cancer in Europe and Latin America) collaborators include:** Ma. Magdalena Echeverry de Polanco, Mabel Elena Bohórquez, Rodrigo Prieto, Angel Criollo, Carolina Ramírez, Ana Patricia Estrada, Jhon Jairo Suárez (Grupo de Citogenética Filogenia y Evolución de Poblaciones, Universidad del Tolima, Colombia); Augusto Rojas Martinez (Center for Research and Development in Health Sciences, Universidad Autónoma de Nuevo León, Monterrey, Mexico); Silvia Rogatto, Samuel Aguiar Jnr, Ericka Maria Monteiro Santos (Department of Urology, School of Medicine, UNESP - São Paulo State University, Botucatu, Brazil); Monica Sans, Valentina Colistro, Pedro C. Hidalgo, Patricia Mut (Department of Biological Anthropology, College of Humanities and Educational Sciences, University of the Republic, Magallanes, Montevideo, Uruguay); Angel Carracedo, Clara Ruiz Ponte, Ines Quntela Garcia (Fundacion Publica Galega de Medicina Xenomica, CIBERER, Genomic Medicine Group-University of Santiago de Compostela, Hospital Clinico, Santiago de Compostela, Galicia, Spain); Sergi Castellvi-Bel (Department of Gastroenterology, Institut de Malalties Digestives i Metabòliques, Hospital Clínic, Centro de Investigación Biomédica en Red de Enfermedades Hepáticas y Digestivas, IDIBAPS, University of Barcelona, Barcelona, Catalonia, Spain); Manuel Teixeira (Department of Genetics, Portuguese Oncology Institute, Rua Dr. António Bernardino de Almeida, Porto, Portugal).

**The NSECG Group comprises:** Ian Tomlinson (Oxford University); M Adams, A Al-Samarraie, S Anwar, R Athavale, S Awad, A Bali, A Barnes, G Cawdell, S Chan, K Chin, P Cornes, M Crawford, J Cullimore, S Ghaem-Maghami, R Gornall, J Green, M Hall, M Harvey, J Hawe, A Head, J Herod, M Hingorani, M Hocking, C Holland, T Hollingsworth, J Hollingworth, T Ind, R Irvine, C Irwin, M Katesmark, S Kehoe, G Kheng-Chew, K Lankester, A Linder, D Luesley, C B-Lynch, V McFarlane, R Naik, N Nicholas, D Nugent, S Oates, A Oladipo, A Papadopoulos, S Pearson, D Radstone, S Raju, A Rathmell, C Redman, M Rymer, P Sarhanis, G Sparrow, N Stuart, S Sundar, A Thompson, S Tinkler, S Trent, A Tristram, N Walji, R Woolas.

**RENDOCAS investigators include:** Annika Lindblom, Gerasimos Tzortzatos, Miriam Mints, Emma Tham, Ofra Castro, Kristina Gemzell-Danielsson.

**SEARCH collaborators include:** Helen Baker, Caroline Baynes, Don Conroy, Bridget Curzon, Patricia Harrington, Sue Irvine, Craig Luccarini, Rebecca Mayes, Hannah Munday, Barbara Perkins, Daisy Pharoah, Radka Platte, Anne Stafford and Judy West.

### BCAC and OCAC Study Collaborators (for control samples):

***The Australian Ovarian Cancer Study Group comprises:*** R Stuart-Harris; NSW- F Kirsten, J Rutovitz, P Clingan, A Glasgow, A Proietto, S Braye, G Otton, J Shannon, T Bonaventura, J Stewart, S Begbie, M Friedlander, D Bell, S Baron-Hay, A Ferrier (deceased), G Gard, D Nevell, N Pavlakis, S Valmadre, B Young, C Camaris, R Crouch, L Edwards, N Hacker, D Marsden, G Robertson, P Beale, J Beith, J Carter, C Dalrymple, R Houghton, P Russell, L Anderson, M Links, J Grygiel, J Hill, A Brand, K Byth, R Jaworski, P Harnett, R Sharma, G Wain; QLD- D Purdie, D Whiteman, B Ward, D Papadimos, A Crandon, M Cummings, K Horwood. A Obermair, L Perrin, D Wyld, J Nicklin; SA- M Davy, MK Oehler, C Hall, T Dodd, T Healy, K Pittman, D Henderson, J Miller, J Pierdes, A Achan; TAS- P Blomfield, D Challis, R McIntosh, A Parker; VIC- B Brown, R Rome, D Allen, P Grant, S Hyde, R Laurie M Robbie, D Healy, T Jobling, T Manolitsas, J McNealage, P Rogers, B Susil, E Sumithran, I Simpson, I Haviv, K Phillips, D Rischin, S Fox, D Johnson, S Lade, P Waring, M Loughrey, N O'Callaghan, B Murray, L Mileshekin, P Allan; V Billson, J Pyman, D Neesham, M Quinn, A Hamilton, C Underhill, R Bell, LF Ng, R Blum, V Ganju; WA- I Hammond, A McCartney (deceased), C Stewart, Y Leung, M Buck, N Zeps (WARTN); AOCS Management Group- DDL Bowtell, AC Green, G Chenevix-Trench, A deFazio, D Gertig, PM Webb.

***BSUCH collaborator:*** Peter Bugert

***ESTHER collaborators:*** Volker Arndt, Heiko Müller, Christa Stegmaier

***GENICA Network collaborators:*** Wing-Yee Lo, Christina Justenhoven, Ute Hamann, Thomas Brüning, Beate Pesch, Yon-Dschun Ko, Sylvia Rabstein, Anne Lotz, Christina Baisch, Hans-Peter Fischer, Volker Harth.

## Supplementary Acknowledgements

The authors thank the many individuals who participated in this study and the numerous institutions and their staff who have supported recruitment.

ANECs thanks members of the Molecular Cancer Epidemiology and Cancer Genetic laboratories at QIMR Berghofer Medical Research Institute for technical assistance, and the ANECs research team for assistance with the collection of risk factor information and blood samples. ANECs also gratefully acknowledges the cooperation of the following institutions: NSW: John Hunter Hospital, Liverpool Hospital, Mater Misericordiae Hospital (Sydney), Mater Misericordiae Hospital (Newcastle), Newcastle Private Hospital, North Shore Private Hospital, Royal Hospital for Women, Royal Prince Alfred Hospital, Royal North Shore Hospital, Royal Prince Alfred Hospital, St George Hospital; Westmead Hospital, Westmead Private Hospital; Qld: Brisbane Private Hospital, Greenslopes Hospital, Mater Misericordiae Hospitals, Royal Brisbane and Women's Hospital, Wesley Hospital, Queensland Cancer Registry; SA: Adelaide Pathology Partners, Burnside Hospital, Calvary Hospital, Flinders Medical Centre, Queen Elizabeth Hospital, Royal Adelaide Hospital, South Australian Cancer Registry; Tas: Launceston Hospital, North West Regional Hospitals, Royal Hobart Hospital; Vic: Freemasons Hospital, Melbourne Pathology Services, Mercy Hospital for Women, Royal Women's Hospital, Victorian Cancer Registry; WA: King Edward Memorial Hospital, St John of God Hospitals Subiaco & Murdoch, Western Australian Cancer Registry. SEARCH thanks the SEARCH research team for recruitment, and also acknowledges the assistance of the Eastern Cancer Registration and Information Centre for subject recruitment.

BECS thanks Reiner Strick, Silke Landrith and Sonja Oeser for their logistic support during the study.

CAHRES (formerly known as SASBAC) thanks Li Yuqing from the Genome Institute of Singapore for contributions to this study, and also acknowledges previous input to SASBAC resource creation by Anna Christensson, Boel Bissmarck, Kirsimari Aaltonen, Karl von Smitten, Nina Puolakka, Christer Halldén, Lim Siew Lan and Irene Chen, Lena U.

Rosenberg, Mattias Hammarström, and Eija Flygare.

HJECS thanks Wen Zheng, Hermann Hertel, and Tjong-Won Park-Simon at Hannover Medical School for their contribution to sample recruitment.

LES gratefully acknowledges Helena Soenen, Gilian Peuteman and Dominiek Smeets for their technical assistance.

MECS thanks Tom Sellers, Catherine Phelan, Andrew Berchuck, and Kimberly Kalli, Amanda von Bismarck, Luisa Freyer and Lisa Rogmann.

MoMaTEC thanks Britt Edvardsen, Ingjerd Bergo and Mari Kyllsø Halle for technical assistance and Inger Marie Aksnes and Tor Audun Hervig at the Blood bank, Haukeland University Hospital for assistance with control recruitment.

NECS thanks staff at the University of Newcastle and the Hunter Medical Research Institute. NSECG thank Ella Barclay and Lynn Martin for their contribution, and acknowledge the invaluable help of the National Cancer Research Network with the collection of study participants.

QIMR Berghofer thanks Margie Wright, Lisa Bowdler, Sara Smith, Megan Campbell and Scott Gordon for control sample collection and data processing, Kerenaftali Klein for statistical advice and Brendan Ryan for assistance with the figures.

REDOCAS thanks Berith Wejderot, Sigrid Sahlen, Tao Liu, Margareta Ström, Maria Karlsson, and Birgitta Byström for their contribution to the study.

BSUCH thanks the Medical Faculty, Mannheim, the Diemtmar Hopp Foundation and the German Cancer Research Center.

MCCS was made possible by the contribution of many people, including the original investigators and the diligent team who recruited the participants and who continue working on follow up. We would like to express our gratitude to the many thousands of Melbourne residents who continue to participate in the study.

The UKBGS thank Breakthrough Breast Cancer and the Institute of Cancer Research for support and funding, and the Study participants, Study staff, and the doctors, nurses and other health care staff and data providers who have contributed to the Study. The ICR acknowledges NHS funding to the NIHR Biomedical Research Centre.

In addition, the iCOGS study would not have been possible without the contributions of: Andrew Berchuck (OCAC), Rosalind A. Eeles, Ali Amin Al Olama, Zsofia Kote-Jarai, Sara Benlloch (PRACTICAL), Georgia Chenevix-Trench, Antonis Antoniou, Lesley McGuffog, Fergus Couch and Ken Offit (CIMBA), Andrew Lee, and Ed Dicks, Craig Luccarini and the staff of the Centre for Cancer Genetic Epidemiology Laboratory (Cambridge), Javier Benitez, Anna Gonzalez-Neira and the staff of the CNIO genotyping unit, Jacques Simard and Daniel C. Tessier, Francois Bacot, Daniel Vincent, Sylvie LaBoissière and Frederic Robidoux and the staff of the McGill University and Génome Québec Innovation Centre, Stig E. Bojesen, Sune F. Nielsen, Borge G. Nordestgaard, and the staff of the Copenhagen DNA laboratory, Sharon A. Windebank, Christopher A. Hilker, Jeffrey Meyer and the staff of Mayo Clinic Genotyping Core Facility.

### Acknowledgments of Funding to BCAC/OCAC control groups

The ESTHER study was funded by the Baden-Württemberg state Ministry of Science, Research and Arts (Stuttgart, Germany), the Federal Ministry of Education and Research (Berlin, Germany) and the Federal Ministry of Family Affairs, Senior Citizens, Women and Youth (Berlin, Germany).

The GENICA was funded by the Federal Ministry of Education and Research (BMBF) Germany grants 01KW9975/5, 01KW9976/8, 01KW9977/0 and 01KW0114, the Robert Bosch Foundation, Stuttgart, Deutsches Krebsforschungszentrum (DKFZ), Heidelberg, Institute for Prevention and Occupational Medicine of the German Social Accident Insurance, Institute of the Ruhr University Bochum (IPA), Germany, as well as the Department of Internal Medicine, Evangelische Kliniken Bonn gGmbH, Johanniter Krankenhaus, Bonn, Germany.

Financial support for the KARBAC study was provided through the regional agreement on medical training and clinical research (ALF) between Stockholm County Council and Karolinska Institutet, as well as the Swedish Cancer Society.

The MARIE study was supported by the Deutsche Krebshilfe e.V. [70-2892-BR I], the Hamburg Cancer Society and the German Cancer Research Center

MAY was supported by R01-CA122443, P50-CA136393, the Fred C. and Katherine B. Andersen Foundation, and the Mayo Foundation.

MCBCS recognizes funding from the Breast Cancer Research Foundation (BCRF), the David F. and Margaret T. Grohne Family Foundation, and the Ting Tsung and Wei Fong Chao Foundation

MCCS recruitment was funded by VicHealth and Cancer Council Victoria, and its follow-up has been continuously supported by infrastructure provided by Cancer Council Victoria.

UKBGS was funded by Breakthrough Breast Cancer and the Institute of Cancer Research, which acknowledges NHS funding to the NIHR Biomedical Research Centre