



## OPEN *NTHL1* is a recessive cancer susceptibility gene

Anna K. Nurmi<sup>1</sup>, Liisa M. Pelttari<sup>1</sup>, Johanna I. Kiiski<sup>1</sup>, Sofia Khan<sup>1</sup>, Mika Nurmikolu<sup>1</sup>, Maija Suvanto<sup>1</sup>, Niina Aho<sup>1</sup>, Tiina Tasmuth<sup>2</sup>, Eija Kalso<sup>2</sup>, Johanna Schleutker<sup>3</sup>, Anne Kallioniemi<sup>4</sup>, Päivi Heikkilä<sup>5</sup>, FinnGen\*, Kristiina Aittomäki<sup>6</sup>, Carl Blomqvist<sup>7</sup> & Heli Nevanlinna<sup>1</sup>✉

In search of novel breast cancer (BC) risk variants, we performed a whole-exome sequencing and variant analysis of 69 Finnish BC patients as well as analysed loss-of-function variants identified in DNA repair genes in the Finns from the Genome Aggregation Database. Additionally, we carried out a validation study of *SERPINA3* c.918-1G>C, recently suggested for BC predisposition. We estimated the frequencies of 41 rare candidate variants in 38 genes by genotyping them in 2482–4101 BC patients and in 1273–3985 controls. We further evaluated all coding variants in the candidate genes in a dataset of 18,786 BC patients and 182,927 controls from FinnGen. None of the variants associated significantly with cancer risk in the primary BC series; however, in the FinnGen data, *NTHL1* c.244C>T p.(Gln82Ter) associated with BC with a high risk for homozygous (OR = 44.7 [95% CI 6.90–290],  $P = 6.7 \times 10^{-5}$ ) and a low risk for heterozygous women (OR = 1.39 [1.18–1.64],  $P = 7.8 \times 10^{-5}$ ). Furthermore, the results suggested a high risk of colorectal, urinary tract, and basal-cell skin cancer for homozygous individuals, supporting *NTHL1* as a recessive multi-tumour susceptibility gene. No significant association with BC risk was detected for *SERPINA3* or any other evaluated gene.

Cancer is a genetic disease in which accumulating pathogenic variants give growth advantage to malignant cells. Eukaryotic cells have specialized pathways for the repair of different mutation types and others that control the cell cycle checkpoints or initiate apoptosis. Defective DNA damage response mechanisms increase genomic instability and may lead to tumour development<sup>1</sup>.

The validated breast cancer (BC) risk genes to date function primarily in DNA double-strand break and interstrand crosslink repair via the homologous recombination and the Fanconi anaemia (FA) pathways and in DNA damage checkpoint signalling<sup>2,3</sup>. The high-penetrance BC risk genes, *BRCA1* and *BRCA2*, encode proteins at the core of the pathways, promoting DNA repair in response to damage signalling<sup>2</sup>. The validated moderate-to-high risk BC predisposition genes, *PALB2*, *CHEK2*, *ATM*, *BARD1*, *RAD51C*, and *RAD51D*, have functions linked to *BRCA1* and *BRCA2*<sup>2,3</sup>. Studies on hereditary BC risk have most often focused on the DNA damage response genes. Other pathways may also be involved in the BC risk predisposition; for example, the syndromic cancer genes and the low-penetrance variants associated with BC risk show a wide range of affected pathways<sup>2-4</sup>.

The high- and moderate-risk variants in the established BC predisposition genes have an autosomal dominant inheritance pattern, even if with incomplete penetrance. Recessive model has also been suggested for increased risk of BC<sup>5</sup>, but to date, no recessive high- or moderate-risk BC susceptibility gene has been validated. Recently, several BC patients with pathogenic biallelic *NTHL1* variants have been described<sup>6-12</sup>, indicating recessive BC predisposition. Pathogenic variants in the *NTHL1* gene have been determined to cause a recessive multi-tumour syndrome, which is characterized especially by adenomatous polyposis and colorectal cancer (CRC), and with accumulating evidence, BC in women<sup>6-13</sup>.

<sup>1</sup>Department of Obstetrics and Gynecology, University of Helsinki and Helsinki University Hospital, Biomedicum Helsinki, P.O. Box 700, 00290 Helsinki, Finland. <sup>2</sup>Department of Anaesthesiology, Intensive Care and Pain Medicine, University of Helsinki and Helsinki University Hospital, Helsinki, Finland. <sup>3</sup>Institute of Biomedicine, University of Turku, and FICAN West Cancer Centre, and Department of Genomics, Laboratory Division, Turku University Hospital, Turku, Finland. <sup>4</sup>Tays Cancer Center, Tampere University Hospital, and BioMediTech Institute and Faculty of Medicine and Health Technology, Tampere University, and Finlab Laboratories, Tampere, Finland. <sup>5</sup>Department of Pathology, University of Helsinki and Helsinki University Hospital, Helsinki, Finland. <sup>6</sup>Department of Clinical Genetics, University of Helsinki and Helsinki University Hospital, Helsinki, Finland. <sup>7</sup>Department of Oncology, University of Helsinki and Helsinki University Hospital, Helsinki, Finland. \*A list of authors and their affiliations appears at the end of the paper. ✉email: heli.nevanlinna@hus.fi

The genes and causal variants contributing to a large proportion of the hereditary BC risk are yet to be discovered<sup>4</sup>. The genetic bottleneck events in the Finnish population have resulted in less overall variation and a higher frequency of loss-of-function (LoF) variants, including recessive disease variants, in the Finns compared to other Europeans<sup>14,15</sup>. This founder effect present in the Finns is advantageous for genetic research as it facilitates the detection of novel disease genes and variants. Only a few recurrent variants account for most of the pathogenic burden in the validated BC risk genes in Finnish BC patients<sup>16</sup>. High-risk *BRCA1/2* variants have been identified in about 21% of Finnish BC families and 1.8% of unselected BC patients<sup>16–18</sup>. The combined frequency of pathogenic variants in the other validated high- and moderate-risk BC susceptibility genes is about 10% in Finnish BC families and 5% in unselected BC patients<sup>16</sup>.

With the aim of identifying novel BC risk variants, we have performed a whole-exome sequencing (WES) and variant analysis of 69 patients from Finnish BC families as well as an analysis of predicted loss-of-function (pLoF) variants in 520 DNA repair genes, detected in approximately 11,000 Finns from the Genome Aggregation Database (gnomAD), and selected candidate risk variants for a case–control study. Additionally, a recent Finnish study reported a putative novel moderate-risk BC susceptibility variant *SERPINA3* c.918-1G>C<sup>19</sup>, warranting further validation. Here, we evaluated *SERPINA3* c.918-1G>C alongside the other candidate variants for BC risk.

## Results

We selected altogether 41 candidate variants in 38 genes, presented in detail in the Supplementary Table S1, for genotyping in BC patients and controls from the Helsinki and Tampere regions in Southern Finland and assessed the variants for cancer risk (Fig. 1). Finally, we retrieved the data for cancer risk association analyses from the FinnGen project and examined the candidate genes and variants in this large series of cancer patients and controls.

### Breast cancer risk association analyses in the Helsinki and Tampere series

We genotyped 19 of the selected candidate variants in 2482 BC patients and 1273 controls, 20 of the variants in 3151 BC patients and 2089 controls, and two of the variants in 4101 BC patients and 3985 controls from the Helsinki and Tampere regions. After the Bonferroni correction for multiple comparisons ( $P < 0.0012$ ), none of the studied variants associated significantly with BC risk in this primary study (Table 1, Supplementary Table S2). We detected two variants, *MAD1L1* NM\_001013836.2:c.1947C>G p.(Tyr649Ter) and *USP45* NM\_001346022.3:c.2190C>A p.(Tyr730Ter), with a higher frequency in the patients than in the controls on a nominally significant level ( $P < 0.05$ ) (Table 1); however, another pLoF in the *USP45* gene, NM\_001346022.3:c.1008del p.(Val337SerfsTer9), was found only slightly more often in the patients than in the controls.

*FANCG* NM\_004629.2:c.1182\_1192delinsC p.(Glu395TrpfsTer5), *NTHL1* NM\_002528.7:c.244C>T p.(Gln82Ter) (also known as NM\_002528.6:c.268C>T p.(Gln90Ter) in reference to the previous transcript version), and *ERCC6L2* NM\_020207.7:c.1424del p.(Ile475ThrfsTer36) (previously denoted as NM\_020207.5:c.1457del p.(Ile486ThrfsTer36)) have been identified to cause recessive hereditary diseases with increased risk of cancer<sup>6,9,20–23</sup>. Here, we detected no significant association between the heterozygous pLoFs and BC risk. Of note, *FANCG* c.1182\_1192delinsC was very rare in our patient series and only detected in 0.2% (6/3147) of the patients and in 0.05% (1/2086) of the controls. Only two patients were homozygous for the *NTHL1* c.244C>T variant, and we were unable to study any recessive BC risk associated with *NTHL1* in our patient series. No study subject was homozygous for *ERCC6L2* c.1424del.

We found *SERPINA3* NM\_001085.5:c.918-1G>C with a similar frequency in the patients and in the controls and detected no association between the variant and BC risk. Previously, the c.918-1G>C carriers were reported to have a medullary breast tumour type more often than noncarriers<sup>19</sup>. Here, no c.918-1G>C carrier had medullary BC: eight patients had ductal, one patient had lobular, and two patients had carcinomas of mixed type.

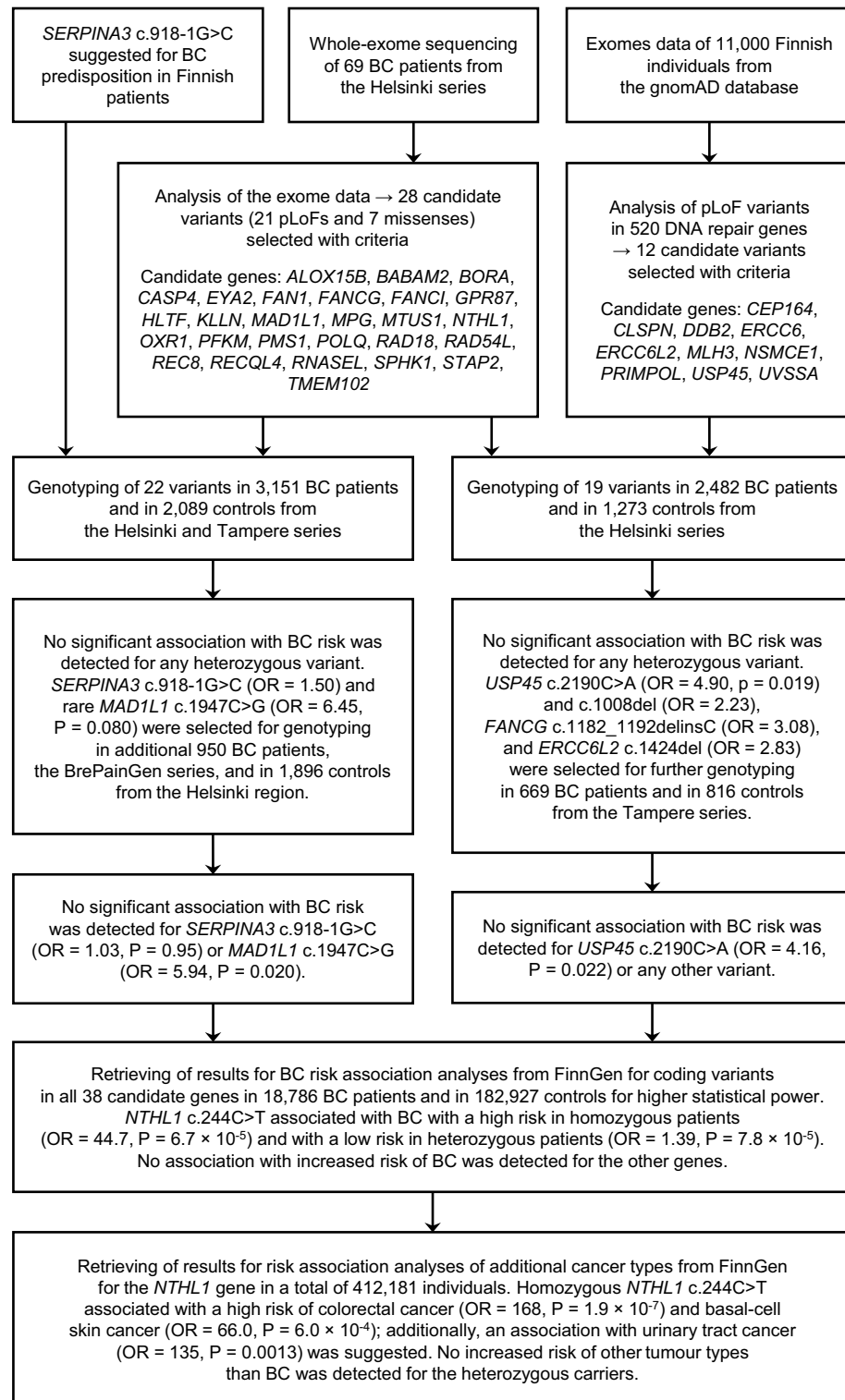
The other studied variants were either detected only in a few patients or the analyses did not suggest an increased risk of BC (Supplementary Table S2).

### Breast cancer risk association analyses from FinnGen

To further evaluate the candidate genes and variants in a dataset with higher statistical power, we retrieved the results for BC risk association analyses from the FinnGen study, data release 10, for all coding variants in the studied genes in 18,786 Finnish BC patients and in 182,927 controls<sup>15,24</sup>. The FinnGen data also provided recessive association analyses for *NTHL1* c.244C>T and *ERCC6L2* c.1424del, which we were unable to perform in the Helsinki and Tampere BC series.

The genotype data suggested a low increased risk of BC for heterozygous *NTHL1* c.244C>T carriers in the additive model (odds ratio (OR) = 1.39 [95% confidence interval (CI) 1.18–1.64],  $P = 7.8 \times 10^{-5}$ ) (Tables 2, 3). Carriers were detected with a similar frequency in the oestrogen receptor (ER)-positive patient group (OR = 1.41 [1.14–1.73],  $P = 0.0012$ ) and in the ER-negative patient group (OR = 1.44 [1.06–1.95],  $P = 0.020$ ) (Table 3). The recessive model suggested a notable risk of BC for homozygous individuals (OR = 44.7 [6.90–290],  $P = 6.7 \times 10^{-5}$ ), both in the ER-positive patient group (OR = 82.1 [10.2–660],  $P = 3.4 \times 10^{-5}$ ) and in the ER-negative patient group (OR = 86.3 [4.89–1523],  $P = 0.0023$ ) (Table 3). Another, a much rarer pLoF in the *NTHL1* gene, c.674dup p.(Ser226ValfsTer39), was found only in heterozygous state (OR = 3.01 [0.67–13.6],  $P = 0.15$ ) (Table 2); therefore, recessive analysis was not available for this variant.

No variant significantly associated with BC risk ( $P < 0.0012$ ) in the other candidate genes (Table 2, Supplementary Table S3). In more detail, no risk association was detected for *MAD1L1* c.1947C>G (OR = 0.87 [0.59–1.27],  $P = 0.47$ ), *SERPINA3* c.918-1G>C (OR = 1.15 [0.86–1.54],  $P = 0.35$ ), or *USP45* c.2190C>A (OR = 0.90 [0.67–1.21],  $P = 0.48$ ). *FANCG* c.1182\_1192delinsC was not included in the FinnGen data, but two other, albeit very rare



**Figure 1.** An overview of the work process and findings of the study.

*FANCG* pLoFs, c.832dup p.(Ala278GlyfsTer11) and c.1076+1G>A, were detected in the study subjects. *ERCC6L2* c.1424del was found with a similar frequency in the patients as in the controls (OR = 1.09 [0.89–1.33], P = 0.42); however, another pLoF in *ERCC6L2*, c.123dup p.(Ile42TyrfsTer5), was more frequent in the patients compared with the controls (OR = 5.08 [1.56–16.5], P = 0.0070). Of the *ERCC6L2* variants, recessive analysis was available only for c.1424del (recessive OR = 20.6 [1.40–303], P = 0.027).

Variant	Group	Total	Carriers	%	OR	95% CI	P value
<i>MAD1L1</i>	Controls	3974	2	0.05			
c.1947C>G	All BC	4083	12	0.3	5.94	1.62–38.2	0.020
p.(Tyr649Ter)	Familial BC	1524	6	0.4	8.13	1.86–55.7	0.011
rs121908981	Unselected BC	3327	8	0.2	4.81	1.20–31.9	0.047
	ER-positive BC	3172	10	0.3	6.46	1.70–42.1	0.016
	ER-negative BC	720	2	0.3	5.61	0.67–46.9	0.085
<i>SERPINA3</i>	Controls	3978	10	0.3			
c.918-1G>C	All BC	4095	11	0.3	1.03	0.43–2.47	0.95
rs199710314	Familial BC	1527	3	0.2	0.75	0.17–2.46	0.66
	Unselected BC	3339	9	0.3	1.07	0.42–2.65	0.89
	ER-positive BC	3184	9	0.3	1.08	0.43–2.69	0.86
	ER-negative BC	720	1	0.1	0.54	0.03–2.83	0.56
<i>ERCC6L2</i>	Controls	2083	12	0.6			
c.1424del	All BC	3142	12	0.4	0.76	0.33–1.75	0.51
p.(Ile475ThrfsTer36)	Familial BC	1369	6	0.4	1.00	0.34–2.68	1.00
rs768081343	Unselected BC	2386	8	0.3	0.64	0.25–1.57	0.34
	ER-positive BC	2386	10	0.4	0.86	0.35–2.03	0.72
	ER-negative BC	573	1	0.2	0.42	0.02–2.20	0.41
<i>FANCG</i>	Controls	2086	1	0.05			
c.1182_1192delinsC	All BC	3147	6	0.2	3.08	0.53–58.2	0.30
p.(Glu395TrpfsTer5)	Familial BC	1368	4	0.3	4.49	0.66–87.9	0.18
rs397507559	Unselected BC	2391	3	0.1	2.21	0.28–44.8	0.49
	ER-positive BC	2389	5	0.2	3.39	0.55–65.0	0.27
	ER-negative BC	576	1	0.2	2.82	0.11–71.4	0.46
<i>NTHL1</i>	Controls	2081	17	0.8			
c.244C>T	All BC	3117	30	1.0	1.35	0.74–2.54	0.33
p.(Gln82Ter)	Familial BC	1357	16	1.2	1.77	0.86–3.64	0.12
rs150766139	Unselected BC	2370	22	0.9	1.25	0.66–2.41	0.49
	ER-positive BC	2366	17	0.7	0.99	0.49–1.99	0.98
	ER-negative BC	573	9	1.6	2.32	0.97–5.21	0.048
<i>USP45</i>	Controls	2088	3	0.1			
c.2190C>A	All BC	3148	21	0.7	4.16	1.42–17.7	0.022
p.(Tyr730Ter)	Familial BC	1368	12	0.9	5.49	1.70–24.6	0.0097
rs118066385	Unselected BC	2392	17	0.7	4.59	1.53–19.7	0.015
	ER-positive BC	2389	17	0.7	4.31	1.43–18.6	0.021
	ER-negative BC	576	2	0.3	2.59	0.33–16.2	0.31
<i>USP45</i>	Controls	2086	6	0.3			
c.1008del	All BC	3148	13	0.4	1.29	0.50–3.73	0.61
p.(Val337SerfsTer9)	Familial BC	1369	8	0.6	1.92	0.65–6.01	0.24
rs554927779	Unselected BC	2392	9	0.4	1.23	0.44–3.69	0.70
	ER-positive BC	2390	11	0.5	1.46	0.55–4.30	0.46
	ER-negative BC	575	2	0.3	1.24	0.18–5.53	0.80

**Table 1.** Variant frequencies in breast cancer patients from the Helsinki and Tampere regions. The familial and the unselected patient groups overlap: 775 patients were included in both groups in the analyses of the *MAD1L1* and *SERPINA3* variants and 614 in the analyses of the other variants. Two of the *NTHL1* c.244C>T carriers included in the analysis were homozygous.

### Breast tumour characteristics of the patients with the *NTHL1* c.244C>T variant

We were able to evaluate the breast tumours of the patients with the *NTHL1* c.244C>T variant further in the Helsinki and Tampere BC series. Two patients from Helsinki were homozygous for the variant. One homozygous patient had been diagnosed with BC at the age of 41 years and with rectal and cecum cancers at the age of 47 years. The breast tumour of this patient was ER-positive and progesterone receptor (PR)-positive ductal carcinoma with grade 3. The other homozygous patient had BC at the age of 47 years and cancer of the sigmoid colon at the age of 51 years. This patient had an ER-positive, PR-positive, and HER2-negative ductal breast carcinoma with grade 2. Neither of the homozygous patients had a family history of BC or OC.

The average age of BC diagnosis among the 28 heterozygous carriers was 58.3 years (range 39–88 years), which was higher than the average age of 56.5 years (range 21–95) for all patients in the Helsinki and Tampere series.

Gene	Variant	Effect allele frequency	OR	95% CI	P value
<i>ERCC6L2</i>	c.123dup p.(Ile42TyrfsTer5)	$8.08 \times 10^{-5}$	5.08	1.56–16.5	0.0070
<i>ERCC6L2</i>	c.1125dup p.(Ile376TyrfsTer7)	$4.30 \times 10^{-5}$	1.07	0.19–6.06	0.94
<i>ERCC6L2</i>	c.1424del p.(Ile475ThrfsTer36) <sup>a</sup>	$3.78 \times 10^{-3}$	1.09	0.89–1.33	0.42
<i>ERCC6L2</i>	c.1930C>T p.(Arg644Ter)	$1.55 \times 10^{-4}$	0.69	0.29–1.65	0.40
<i>FANCG</i>	c.832dup p.(Ala278GlyfsTer11)	$1.10 \times 10^{-4}$	0.84	0.26–2.74	0.78
<i>FANCG</i>	c.1076+1G>A	$4.00 \times 10^{-5}$	1.45	0.23–9.18	0.69
<i>MAD1L1</i>	c.150+1G>T	$4.61 \times 10^{-5}$	0.33	0.04–2.57	0.29
<i>MAD1L1</i>	c.538dup p.(Val180GlyfsTer47)	$5.91 \times 10^{-5}$	0.48	0.11–2.12	0.33
<i>MAD1L1</i>	c.1396C>T p.(Gln466Ter)	$6.15 \times 10^{-5}$	1.09	0.21–5.53	0.92
<i>MAD1L1</i>	c.1505+2T>A	$1.54 \times 10^{-5}$	3.73	0.40–34.8	0.25
<i>MAD1L1</i>	c.1947C>G p.(Tyr649Ter) <sup>a</sup>	$9.47 \times 10^{-4}$	0.87	0.59–1.27	0.47
<i>NTHL1</i>	c.244C>T p.(Gln82Ter) <sup>a</sup>	$4.65 \times 10^{-3}$	1.39	1.18–1.64	$7.8 \times 10^{-5}$
<i>NTHL1</i>	c.674dup p.(Ser226ValfsTer39)	$6.15 \times 10^{-5}$	3.01	0.67–13.6	0.15
<i>SERPINA3</i>	c.511C>T p.(Gln171Ter)	$2.49 \times 10^{-4}$	1.03	0.49–2.14	0.95
<i>SERPINA3</i>	c.918-1G>C <sup>a</sup>	$1.96 \times 10^{-3}$	1.15	0.86–1.54	0.35
<i>USP45</i>	c.7del p.(Val3Ter)	$1.51 \times 10^{-4}$	0.62	0.23–1.66	0.34
<i>USP45</i>	c.658G>T p.(Glu220Ter)	$8.92 \times 10^{-5}$	1.97	0.67–5.77	0.22
<i>USP45</i>	c.845+2T>C	$4.28 \times 10^{-3}$	0.86	0.71–1.03	0.10
<i>USP45</i>	c.1008del p.(Val337SerfsTer9) <sup>a</sup>	$4.43 \times 10^{-4}$	0.78	0.41–1.46	0.43
<i>USP45</i>	c.2190C>A p.(Tyr730Ter) <sup>a</sup>	$1.70 \times 10^{-3}$	0.90	0.67–1.21	0.48

**Table 2.** Breast cancer risk association analyses from FinnGen for heterozygous pLoF variants in the candidate genes. The variants denoted with <sup>a</sup> were genotyped in the Helsinki and Tampere BC series. Reference transcripts: *ERCC6L2* NM\_020207.7, *FANCG* NM\_004629.2, *MAD1L1* NM\_001013836.2, *NTHL1* NM\_002528.7, *SERPINA3* NM\_001085.5, and *USP45* NM\_001346022.3.

Cancer type		Total number of individuals		Recessive model			Additive model		
		Patients	Controls	OR	95% CI	P value	OR	95% CI	P value
Breast	Breast cancer	18,786	182,927	44.7	6.90–290	$6.7 \times 10^{-5}$	1.39	1.18–1.64	$7.8 \times 10^{-5}$
	ER-positive breast cancer	10,404	182,678	82.1	10.2–660	$3.4 \times 10^{-5}$	1.41	1.14–1.73	0.0012
	ER-negative breast cancer	6188	182,678	86.3	4.89–1523	0.0023	1.44	1.06–1.95	0.020
Colon	Colorectal cancer	6847	314,193	168	24.4–1152	$1.9 \times 10^{-7}$	1.14	0.86–1.52	0.35
	Colorectal adenocarcinoma	5610	314,193	204	22.7–1837	$2.1 \times 10^{-6}$	0.99	0.73–1.36	0.96
	Colon cancer	4143	314,193	166	14.8–1856	$3.4 \times 10^{-5}$	1.19	0.83–1.70	0.35
	Colon adenocarcinoma	3212	314,193	224	10.8–4643	$4.7 \times 10^{-4}$	0.99	0.65–1.49	0.95
	Rectal cancer	2490	314,193	447	49.7–4023	$5.2 \times 10^{-8}$	1.04	0.67–1.63	0.85
	Adenocarcinoma, papilloma adenocarcinoma, and mucinous carcinoma of rectum	2545	314,193	472	52.1–4279	$4.4 \times 10^{-8}$	1.13	0.72–1.75	0.60
Urinary tract	Cancer of the urinary organs	2619	314,193	135	6.73–2713	0.0013	0.94	0.60–1.47	0.79
	Cancer of the renal pelvis	138	314,193	146	3.55–5985	0.0086	3.14	0.41–24.0	0.27
	Bladder cancer	2193	314,193	238	8.92–6334	0.0011	1.45	0.87–2.39	0.15
Other	Basal-cell carcinoma of the skin	20,506	314,193	66.0	6.02–723	$6.0 \times 10^{-4}$	1.16	0.97–1.38	0.11
	Prostate cancer	15,199	131,266	365	1.97–67,342	0.027	1.04	0.84–1.29	0.73

**Table 3.** Cancer risk association analyses from FinnGen for the *NTHL1* c.244C>T variant. The controls included only women for BC and only men for prostate cancer. The risk association analyses of other cancer types for heterozygous *NTHL1* c.244C>T carriers are presented in the Supplementary Table S4.

Of the heterozygous carriers, 75.0% (21/28) had ductal, 17.9% (5/28) had lobular, and 7.1% (2/28) had other invasive breast tumour type. Additionally, 65.4% (17/26) of the patients had ER-positive and 34.6% (9/26) had ER-negative BC, including three patients with triple-negative BC, and 78.3% (18/23) of the patients had a breast tumour with a grade 2 or 3. Additional cancer diagnoses were available only for the patients from the Helsinki BC series: of the 18 heterozygous carriers, two patients had bilateral BC, one had BC and uterus cancer, and one had BC and pancreatic cancer. One patient with bilateral BC and one other heterozygous BC patient also carried a pathogenic *CHEK2* c.1100del variant; no other high- or moderate-risk BC predisposition variants had been found in the *NTHL1* c.244C>T carriers from Helsinki.



### Association of *NTHL1* c.244C>T with increased risk of other cancer types than breast cancer

We obtained the data for recessive risk association analyses from FinnGen for all malignant tumour types diagnosed in the individuals homozygous for the *NTHL1* c.244C>T variant. Besides BC, homozygous *NTHL1* c.244C>T significantly associated with a high risk of CRC (OR = 168 [24.4–1152],  $P = 1.9 \times 10^{-7}$ ) and basal-cell skin cancer (OR = 66.0 [6.02–723],  $P = 6.0 \times 10^{-4}$ ) (Table 3). Additionally, the results suggested an increased risk of urinary tract cancers (OR = 135 [6.73–2713],  $P = 0.0013$ ).

Ten individuals with the homozygous *NTHL1* c.244C>T variant were identified in the FinnGen study: nine of them had been diagnosed with one or multiple tumour types as verified by the Finnish Cancer Registry, and one had no cancer diagnosis. The diagnosed malignant tumour types were rectal, colon, breast, bladder, renal pelvis, basal-cell skin, and prostate cancer, and the non-invasive tumour types were rectal, bladder, and meningial tumour. Altogether, the nine patients had 19 tumour diagnoses.

To examine the cancer risks for the heterozygous carriers, we retrieved the results for additive risk association analyses from FinnGen for the available malignant tumour types, which have been diagnosed in the patients with biallelic *NTHL1* variants in the FinnGen data or reported previously<sup>6–11,13</sup>. No increased risk of cancer was suggested for the heterozygous carriers for other cancer types than BC (Table 3, Supplementary Table S4).

### Discussion

We have performed a WES study of BC patients and a gnomAD database analysis of pLoF variants, with the aim of identifying novel BC risk variants. Furthermore, a recent exome-sequencing study of Finnish patients identified *SERPINA3* as a novel candidate gene for moderate-risk BC predisposition<sup>19</sup>. We assessed the cancer risk associated with the candidate variants by evaluating them in series of BC patients and controls from the Helsinki and Tampere regions and from the FinnGen project.

Even though we did not detect a significant association between *NTHL1* c.244C>T p.(Gln82Ter) and BC risk in our patient series, a much larger genotype dataset from FinnGen showed a high increased risk of BC for homozygous (OR = 44.7,  $P = 6.7 \times 10^{-5}$ ) and a low increased risk for heterozygous women (OR = 1.39,  $P = 7.8 \times 10^{-5}$ ). Different cancer studies have reported a high frequency of BC (55%) among women with biallelic pathogenic *NTHL1* variants, as reviewed by Beck et al.<sup>6–13</sup>. The association of *NTHL1* variants with BC predisposition has previously been evaluated in a large international case–control study; however, just one biallelic patient was identified and the BC risk remained unclear also for the heterozygous carriers<sup>25</sup>. In that study, the carrier frequencies and associated BC risk for the c.244C>T variant varied between patient series, but the results for other, rarer heterozygous pLoF and pathogenic missense variants suggested a low increased risk of BC<sup>25</sup>. The c.244C>T variant (previously reported as c.268C>T p.(Gln90Ter)) is the most frequent LoF variant identified in the patients with *NTHL1* tumour syndrome as well as in the *NTHL1* gene in the gnomAD database<sup>13,26</sup>. The variant is enriched in the uniform Finnish population—it was found with a minor allele frequency (MAF) of 0.0044 in the controls from the FinnGen study—which facilitates the detection of increased risk.

Biallelic pathogenic variants in the *NTHL1* gene cause a high-penetrance multi-tumour syndrome, which is especially manifested with colorectal, breast, endometrial, urothelial, and basal-cell skin cancer, as well as meningial tumours<sup>6–13</sup>. Of the previously reported homozygous and compound heterozygous individuals, 49% had CRC, and of the individuals who had undergone a colonoscopy, even 93% had colonic adenomas<sup>13</sup>. The FinnGen results support the previous findings on high risk of CRC for the individuals with biallelic variants<sup>6,9–11,13</sup>. The present study also indicates a high recessive risk of BC; furthermore, high risks of basal-cell skin carcinoma and urinary tract cancer are suggested. Combining the FinnGen and the Helsinki patient series, 11 out of the identified 12 homozygous individuals had a total of 24 tumour diagnoses, further supporting high-penetrance cancer risk. Other cancer types, which have been reported in more than one biallelic case, include hematologic malignancies, squamous cell carcinomas of the head and neck, thyroid, pancreatic, and prostate cancer, and melanoma<sup>6,7,9–11,13</sup>.

Monoallelic *NTHL1* variants are unlikely to cause a substantially increased risk of cancer if any<sup>8,12,25,27</sup>. In the current study, we examined the risks for the heterozygous carriers to malignant tumours, which have been detected in the patients with biallelic *NTHL1* variants<sup>6–13</sup>. We observed no increased risk of any other cancer type than BC; however, for some tumour types, the case groups were small. In addition to BC, the risk associated with monoallelic *NTHL1* variants has previously been investigated in CRC, polyposis, and in a pan-cancer patient population<sup>8,12,27</sup>. In line with our results, no increased risk of other cancer types was detected.

The premature stop codon caused by the *NTHL1* c.244C>T variant has been reported to activate the nonsense-mediated mRNA decay surveillance mechanism<sup>6</sup>, resulting in loss of the *NTHL1* gene product in homozygous individuals. Consistently, reduced *NTHL1* protein expression has been observed in heterozygous carriers<sup>25</sup>. The *NTHL1* protein is a bifunctional DNA glycosylase, which catalyses the initial step of base excision-repair pathway to remove oxidative DNA damage<sup>28–30</sup>. *NTHL1* has glycosylase activity on damaged bases, with a preference for oxidized pyrimidines as the substrate, and apurinic/aprimidinic lyase activity on the DNA phosphate backbone<sup>28,29</sup>. Disruption of the *NTHL1* function may lead to mispairing of damaged bases in replication and accumulation of sequence-specific mutations<sup>30</sup>. Biallelic LoF variants in the *NTHL1* gene have been shown to drive a mutational process causing the COSMIC signature SBS30, which is characterized by somatic C>T transitions at non-CpG sites over different tumour types, including BC<sup>6,9,12,25,31,32</sup>. Although there is some contradiction, the mutational signature 30, somatic loss of a second allele, or promoter methylation have typically not been observed in heterozygous *NTHL1* variant carriers<sup>12,25,27,32,33</sup>—in these individuals, the possible increased risk of cancer has been suggested to be caused by haploinsufficiency<sup>25</sup>.

The current study is a comprehensive cancer risk analysis for *NTHL1* in an extensive case–control material. Previous studies have been unable to estimate the associated risks for the biallelic individuals in a case–control setting. In the FinnGen data, the prevalence of individuals homozygous for the *NTHL1* c.244C>T variant was 1

in every 41,200. This is higher than the estimate of 1 in 114,770 Europeans<sup>30</sup>. Still, due to the rarity of homozygous individuals, the observed effect sizes for the increased recessive risk associated with the c.244C>T variant, here, are uncertain and the CIs are wide, and the *NTHL1* gene warrants further evaluation for more precise risk estimates for different cancer types. Nevertheless, because of the high cancer risk, we suggest that *NTHL1* should be included in cancer gene panels in clinical diagnostics, at least for the most common tumour types reported in the patients with pathogenic biallelic *NTHL1* variants. Additionally, the susceptibility to multiple tumour types should be considered in surveillance and cancer-prevention strategies for the individuals with biallelic variants, and clinical practice guidelines should be developed for the *NTHL1* gene.

*FANCG* c.1182\_1192delinsC p.(Glu395TrpfsTer5) was rare in our patient series, and it was not included in the FinnGen dataset; hence, we were unable to statistically assess any BC risk associated with it. *FANCG* is an established FA risk gene, with p.(Glu395fs) among the first described causative *FANCG* mutations for the syndrome<sup>20,21</sup>. Monoallelic variants in several FA genes are known to predispose to BC<sup>3</sup>. Two other *FANCG* pLoF variants, c.832insG p.(Ala278GlyfsTer11) and c.1076+1G>A, identified in the BC patients in the FinnGen study, have been discovered also in Finnish FA patients<sup>34</sup>. No association with increased risk of BC was detected for these two variants in the FinnGen data; however, both variants were very rare in the study subjects. We did not find heterozygous *ERCC6L2* variants associated with BC risk. The additive ORs were inconsistent between the different *ERCC6L2* variants in the FinnGen data, which may have been influenced by the rarity of the variants. Biallelic LoF variants in the *ERCC6L2* gene, including homozygous c.1424del p.(Ile475ThrfsTer36) (previously known as c.1457del), have been described in patients with inherited bone marrow failure and acute myeloid leukaemia<sup>22,23</sup>. Additionally, a BC patient with biallelic variants has been reported<sup>23</sup>. The homozygous c.1424del variant was detected among the BC patients also in the current study, and the contribution of *ERCC6L2* to BC remains unclear.

We identified *MAD1L1* c.1947C>G p.(Tyr649Ter) and *USP45* c.2190C>A p.(Tyr730Ter) in about four- to fivefold higher frequency in the unselected patient group compared with the controls from the Helsinki and Tampere regions. A recent copy number variant analysis reported a twofold increased frequency of *MAD1L1* gene deletions among patients in a large BC dataset<sup>35</sup>; additionally, p.(Tyr649Ter) has been suggested to have a dominant-negative effect on the *MAD1L1* protein function and impair the mitotic spindle-assembly checkpoint<sup>36</sup>. Other studies have connected *USP45* to hypersensitivity to mitomycin C -induced interstrand crosslinks and as a candidate gene to multiple myeloma<sup>37,38</sup>. Our results did not remain significant after adjusting the P value threshold for multiple comparisons and no association with BC risk was detected for the *MAD1L1* and *USP45* genes in the FinnGen data. We found the *SERPINA3* c.918-1G>C variant with a similar frequency in the BC patients and in the controls both in the Helsinki and Tampere BC series and in the FinnGen data; therefore, in the current study, no association with increased BC risk was detected.

In conclusion, our results indicate that biallelic LoF variants in the *NTHL1* gene cause a high risk of multiple cancer types, including BC. We also suggest *NTHL1* as a low-risk gene for BC predisposition in heterozygous women. However, further studies are required to estimate the effect sizes for the increased risk of different cancer types more precisely. Finally, we propose that *NTHL1* should be included in cancer gene panels in clinical diagnostics and clinical practice guidelines should be developed for cancer screening strategies for individuals with pathogenic biallelic *NTHL1* variants.

## Materials and methods

### Whole-exome sequencing and variant calling

We included 69 BC patients from 44 families in the WES. Of the families, 39 had at least three patients with BC or OC among first- and second-degree relatives and 4 had two affected first-degree relatives. Furthermore, 10 of the families included male BC patients, 19 families had uterine cancer cases, and 8 families were suspected of Li-Fraumeni-like syndrome. None of the exome-sequenced patients had a pathogenic *BRCA1/2* or *TP53* variant. The index patients and their family members were collected among the Helsinki BC series as described below. The WES was carried out using genomic DNA extracted from peripheral blood samples.

The sequencing and variant calling was performed at the McGill University and Génome Québec Innovation Centre, Montreal, Canada. Exome libraries were created with Roche Nimblegen SeqCap EZ Exome + UTR capture kit for 39 of the samples and Roche Nimblegen SeqCap EZ Exome v3 kit for 30 of the samples. Sequencing of the libraries was performed with Illumina HiSeq 2000 platform with 100 bp paired-end reads. The read quality trimming of FASTQ files was executed with FASTX-toolkit ([http://hannonlab.cshl.edu/fastx\\_toolkit/](http://hannonlab.cshl.edu/fastx_toolkit/)). The reads were aligned to the human reference genome GRCh37/hg19 with Burrows-Wheeler Aligner<sup>39</sup>. Insertion and deletion variants (indels) were realigned and duplicates were marked with Picard (<https://broadinstitute.github.io/picard/>). The metrics were computed with Integrative Genomics Viewer<sup>40</sup>, and the variant calling was performed with SAMtools and BCFtools<sup>41,42</sup>.

### Variant selection from the whole-exome sequencing data

The candidate variants were selected for genotyping based on MAF, pathogenicity of the variant, and relevant gene function. We annotated the variants with Annovar<sup>43</sup> and retrieved gene ontology (GO) terms from the AmiGO2 website by the Gene Ontology Consortium<sup>44,45</sup>. We excluded variants with a raw read depth of < 30 and a phred-scaled quality probability of < 10. Common variants with a MAF of > 0.03 were excluded using the Exome Aggregation Consortium (in any population) and the 1000Genomes variant databases<sup>46,47</sup>. This selection stage yielded 22,531 variants, which were predicted to alter the protein sequence. We included pLoF variants, defined as stop-gain, frameshift, and essential splice site variants, involved in DNA repair (GO:0006281), cell cycle (GO:0007049), or apoptotic pathways (GO:0006915), totalling in 178 variants in 160 genes. pLoF variants outside of these pathways were considered based on relevance in tumorigenesis. Missense variants involved in

DNA repair or cell cycle pathways were considered if predicted to be pathogenic by CADD<sup>48</sup> (phred  $\geq 20$ ) and by the majority of the other pathogenicity prediction tools included in the LJB (dbNSFP) database in Annovar<sup>43</sup> (201 variants in 174 genes). Finally, we focused on plausible candidate genes based on gene function, queried from the UniProt and the NCBI Gene databases<sup>49,50</sup>, and selected 28 variants in well-supported transcripts<sup>51</sup> for genotyping, including 21 pLoFs and seven missenses (Supplementary Table S1). All selected variants had a raw read depth of  $\geq 600$  and a phred-scaled quality probability of  $\geq 150$  in the WES data. We further confirmed the indel variants with Sanger sequencing. The variant descriptions were confirmed with Mutalyzer 3 and comply with the current HGVS nomenclature<sup>52,53</sup>.

### Variant selection from the gnomAD database

We downloaded the exomes data of approximately 11,000 Finns from the gnomAD database, release 2.0.1, for about 520 DNA repair genes (GO:0006281, release 2017-07-01)<sup>26,44,45</sup>. We selected only high-confidence pLoF variants with a MAF of 0.0001–0.03 in the Finnish population; furthermore, we excluded the variants with a MAF of  $> 0.03$  in any other population. We excluded the variants in the validated BC risk genes and in the candidate risk genes previously published from the Helsinki BC series<sup>3,54,55</sup>. This selection stage yielded 124 pLoF variants in 92 genes in well-supported transcripts (transcript support level 1 and 2), annotated with transcript flags from the Ensembl database through BioMart<sup>51,56</sup>. We prioritized the candidate variants based on gene function<sup>49,50</sup>, similarly as for variants chosen from the WES data, and selected twelve pLoF variants in ten candidate genes for genotyping (Supplementary Table S1).

### Patient and control series

The case–control series included a total of 4101 BC patients and 3985 population controls from the Helsinki and Tampere regions. All study subjects from Helsinki were women, whereas the Tampere control group also included men. The genomic DNA used in genotyping had been extracted from peripheral blood samples.

### Breast cancer patients

The unselected Helsinki BC series consisted of 1726 patients who had been diagnosed with their first primary invasive BC. The patients were recruited consecutively in the Helsinki University Hospital at the Department of Oncology in 1997–1998 and 2000 ( $n = 847$ ) and at the Department of Surgery in 2001–2004 ( $n = 879$ )<sup>18,57,58</sup> without any selection criteria for family history of BC or age of diagnosis. The familial Helsinki BC series was combined from 380 index patients with a family history of BC or OC from the unselected series and from 756 additional index patients who were recruited at the Department of Oncology and at the Department of Clinical Genetics until 2015<sup>58–60</sup>. Of these 1136 familial patients, 606 had a family history of at least three BC or OC patients among first- or second-degree relatives (including the proband) and 530 had one affected first-degree relative. The familial patients had been tested at least for *BRCA1/2* founder mutations in Finland and the carriers had been excluded from the series. The cancer diagnoses of the patients and their family members were confirmed from hospital records and/or the Finnish Cancer Registry. Altogether, the Helsinki BC series included a total of 2482 patients.

Additional unselected BC patients from the Helsinki region, the BrePainGen series, had been collected in the Helsinki University Hospital at the Breast Surgery Unit in 2006–2010<sup>61</sup>. The series consisted of 950 patients with invasive breast tumour, which had been unilateral and non-metastasised at the time of recruitment; however, no selection for family history of the disease or age of diagnosis had been performed. Of the patients, 161 had at least one first- or second-degree relative diagnosed with BC or OC and were classified as familial.

The unselected Tampere BC series consisted of 669 patients who had been recruited at the Tampere University Hospital consecutively in 1997–1999 and additionally in 1996–2004<sup>18,58</sup>. All patients had been newly diagnosed with invasive BC. Altogether 234 patients had at least one first- or second-degree relative diagnosed with BC or OC and were defined familial.

### Population controls

The geographically matched population controls from the Helsinki region consisted of 1273 anonymous blood donors, collected in 2002–2003, and 1896 additional controls with no cancer diagnosis from the Helsinki Biobank. The population controls from the Tampere region consisted of 816 blood donors.

### Variant genotyping

Twenty-one variants selected from the WES data were genotyped in 3143 BC patients and 2089 controls from the Helsinki and Tampere BC series with the Sequenom MassARRAY. Seven indel variants from the WES data were genotyped outside of the array for technical reasons. Changes of  $\leq 6$  base pairs were genotyped with TaqMan real-time PCR and larger indels with 3% agarose gel electrophoresis in 2482 BC patients and 1273 controls from Helsinki. Positive control samples were included in all analyses and the carriers detected with 3% agarose gel electrophoresis were confirmed with Sanger sequencing. Twelve variants selected from the gnomAD data were genotyped in 2482 BC patients and 1273 controls from Helsinki with the Sequenom MassARRAY.

The genotyping of four variants, which had been analysed in the Helsinki BC series, was continued to the 669 BC patients and 816 controls of the Tampere BC series. We genotyped *ERCC6L2* c.1424del and *USP45* c.2190C>A with TaqMan real-time PCR, *USP45* c.1008del with Sanger sequencing, and *FANCG* c.1182\_1192delinsC with 3% agarose gel electrophoresis. The genotyping of *MAD1L1* c.1947C>G was further continued to additional 950 BC patients from the BrePainGen series and 1896 controls from the Helsinki Biobank with TaqMan real-time PCR. *SERPINA3* c.918-1G>C, selected for genotyping outside of the WES or the gnomAD variant data, was genotyped in all 4101 BC patients and 3985 controls with TaqMan real-time PCR. We confirmed the detected



carriers for the *ERCC6L2* c.1424del, *FANCG* c.1182\_1192delinsC, *MAD1L1* c.1947C>G, *NTHL1* c.244C>T, *SERPINA3* c.918-1G>C, and *USP45* c.2190C>A and c.1008del variants with Sanger sequencing. Further details on genotyping are given in the Supplementary Information Methods.

### Statistical analyses

We performed the statistical analyses using the R environment for statistical computing (version 4.2.2)<sup>62</sup>. We used region-adjusted logistic regression for the combined analyses including patients from Helsinki and Tampere BC series and Fisher's exact test for the Helsinki BC series, with two-sided P values. After the Bonferroni correction for multiple comparisons, P values < 0.0012 were considered statistically significant.

### FinnGen data

To further evaluate the candidate genes, we obtained the data for cancer risk association analyses for a total of 412,181 individuals (230,310 women and 181,871 men) from the FinnGen research project (<https://www.finnngen.fi/en>), which produces genotype data from samples of Finnish biobank participants and combines it with longitudinal data from Finnish health registries<sup>24</sup>. The biobank sample and data accession numbers for the FinnGen data release 10 are presented in the Supplementary Information Materials.

We retrieved the results for BC risk association analyses for all 38 candidate genes with the endpoint C3\_BREAST\_EXALLC, which included 18,786 female BC patients and 182,927 female controls with no cancer diagnosis. We annotated the variants with Annovar<sup>43</sup>; from these results, we included pLoF, missense, and in-frame indel variants with a MAF of  $\leq 0.03$  in the controls. Additionally, we retrieved the data for risk association analyses for all available tumour types, which had been detected in cancer patients with biallelic pathogenic variants in the *NTHL1* gene in the FinnGen study and in previous reports<sup>6–13</sup>. We excluded the endpoints for benign and in situ tumours (ICD-10 D-coded tumours), as the registry entries may be incomplete for them, except for the endpoint C3\_BREAST\_EXALLC, which included both malignant and in situ tumours (ICD-O-3 behaviour codes 3 and 2). We used the analyses in which the controls with any cancer diagnosis had been excluded. All included cancer endpoints are given in the Supplementary Table S5 and the endpoint definitions are available at <https://risteys.finregistry.fi>.

The cancer risk associated with heterozygous variants was detected with the additive model in the FinnGen data; homozygous and compound heterozygotes had been excluded from the analyses as described in<sup>15</sup>. The recessive model compared homozygous individuals against heterozygotes and noncarriers<sup>15</sup>. Of the additive analyses, we included only variants which had been genotyped on array, whereas the recessive analyses for *NTHL1* c.244C>T and *ERCC6L2* c.1424del included also imputed genotypes. The imputation quality scores were 0.9974 for *NTHL1* c.244C>T and 0.9951 for *ERCC6L2* c.1424del. The association analyses in the FinnGen data had been performed with the REGENIE software (version 2.2.4)<sup>63</sup>. The genotyping and production of the FinnGen dataset has been described in<sup>24</sup> and at <https://finngen.gitbook.io/documentation>.

### Ethics declarations

The study was conducted in accordance with the Declaration of Helsinki and with approval by the Ethics Committee of the Helsinki University Hospital (Dnro207/E9/07 and HUS71597/2016). The Tampere study protocol was approved by the Ethics Committee of the Pirkanmaa Hospital District (97247) and the BrePainGen study protocol by the Coordinating Ethics Committee (136/E0/2006) and the Ethics Committee of the Department of Surgery (Dnro 148/E6/05) of the Hospital District of Helsinki and Uusimaa. The ethics statement for FinnGen is given in the Supplementary Information Materials. Informed consent was obtained from all patients.

### Data availability

For the Helsinki and Tampere BC series, the data that support the findings of our study are available on reasonable request from the corresponding author. The data are not publicly available due to privacy or ethical restrictions. Instructions on accessing the FinnGen data are available at [https://www.finnngen.fi/en/access\\_results](https://www.finnngen.fi/en/access_results).

Received: 6 July 2023; Accepted: 14 November 2023

Published online: 30 November 2023

### References

1. Hoeijmakers, J. H. Genome maintenance mechanisms for preventing cancer. *Nature* **411**, 366–374. <https://doi.org/10.1038/35077232> (2001).
2. Nielsen, F. C., van-Overeem-Hansen, T. & Sorensen, C. S. Hereditary breast and ovarian cancer: New genes in confined pathways. *Nat. Rev. Cancer* **16**, 599–612. <https://doi.org/10.1038/nrc.2016.72> (2016).
3. Breast Cancer Association Consortium *et al.* Breast cancer risk genes—association analysis in more than 113,000 Women. *N Engl J Med* **384**, 428–439. <https://doi.org/10.1056/NEJMoa1913948> (2021).
4. Michailidou, K. *et al.* Association analysis identifies 65 new breast cancer risk loci. *Nature* **551**, 92–94. <https://doi.org/10.1038/nature24284> (2017).
5. Antoniou, A. C. *et al.* A comprehensive model for familial breast cancer incorporating BRCA1, BRCA2 and other genes. *Br. J. Cancer* **86**, 76–83. <https://doi.org/10.1038/sj.bjc.6600008> (2002).
6. Weren, R. D. *et al.* A germline homozygous mutation in the base-excision repair gene NTHL1 causes adenomatous polyposis and colorectal cancer. *Nat. Genet.* **47**, 668–671. <https://doi.org/10.1038/ng.3287> (2015).
7. Rivera, B., Castellsague, E., Bah, I., van Kempen, L. C. & Foulkes, W. D. Biallelic NTHL1 mutations in a woman with multiple primary tumors. *N. Engl. J. Med.* **373**, 1985–1986. <https://doi.org/10.1056/NEJMc1506878> (2015).
8. Belhadj, S. *et al.* Delineating the phenotypic spectrum of the NTHL1-associated polyposis. *Clin. Gastroenterol. Hepatol.* **15**, 461–462. <https://doi.org/10.1016/j.cgh.2016.09.153> (2017).

9. Grolleman, J. E. *et al.* Mutational signature analysis reveals NTHL1 deficiency to cause a multi-tumor phenotype. *Cancer Cell* **35**, 256–266. <https://doi.org/10.1016/j.ccell.2018.12.011> (2019).
10. Boulouard, F. *et al.* Further delineation of the NTHL1 associated syndrome: A report from the French Oncogenetic Consortium. *Clin. Genet.* **99**, 662–672. <https://doi.org/10.1111/cge.13925> (2021).
11. Weatherill, C. B. *et al.* Six case reports of NTHL1-associated tumor syndrome further support it as a multi-tumor predisposition syndrome. *Clin. Genet.* **103**, 231–235. <https://doi.org/10.1111/cge.14242> (2023).
12. Salo-Mullen, E. E. *et al.* Prevalence and Characterization of Biallelic and Monoallelic NTHL1 and MSH3 Variant Carriers from a Pan-Cancer Patient Population. *JCO Precis. Oncol.* **5**, 455. <https://doi.org/10.1200/PO.20.00443> (2021).
13. Beck, S. H. *et al.* Intestinal and Extraintestinal Neoplasms in Patients with NTHL1 Tumor Syndrome: A Systematic Review. *Fam. Cancer* **21**, 453–462. <https://doi.org/10.1007/s10689-022-00291-3> (2022).
14. Lim, E. T. *et al.* Distribution and Medical Impact of Loss-of-Function Variants in the Finnish Founder Population. *PLoS Genet.* **10**, e1004494. <https://doi.org/10.1371/journal.pgen.1004494> (2014).
15. Heyne, H. O. *et al.* Mono- and Biallelic Variant Effects on Disease at Biobank Scale. *Nature* **613**, 519–525. <https://doi.org/10.1038/s41586-022-05420-7> (2023).
16. Nurmi, A. K. *et al.* Pathogenic Variant Spectrum in Breast Cancer Risk Genes in Finnish Patients. *Cancers (Basel)* **14**, 6158. <https://doi.org/10.3390/cancers14246158> (2022).
17. Vehmanen, P. *et al.* Low Proportion of BRCA1 and BRCA2 Mutations in Finnish Breast Cancer Families: Evidence for Additional Susceptibility Genes. *Hum. Mol. Genet.* **6**, 2309–2315. <https://doi.org/10.1093/hmg/6.13.2309> (1997).
18. Syrjäkoski, K. *et al.* Population-based Study of BRCA1 and BRCA2 Mutations in 1035 Unselected Finnish Breast Cancer Patients. *J. Natl. Cancer Inst.* **92**, 1529–1531. <https://doi.org/10.1093/jnci/92.18.1529> (2000).
19. Koivuluoma, S. *et al.* Exome Sequencing Identifies a Recurrent Variant in SERPINA3 Associating with Hereditary Susceptibility to Breast Cancer. *Eur. J. Cancer* **143**, 46–51. <https://doi.org/10.1016/j.ejca.2020.10.033> (2021).
20. de Winter, J. P. *et al.* The Fanconi Anemia Group G Gene FANCG is Identical with XRCC9. *Nat. Genet.* **20**, 281–283. <https://doi.org/10.1038/3093> (1998).
21. Demuth, I. *et al.* Spectrum of Mutations in the Fanconi Anemia Group G Gene, FANCG/XRCC9. *Eur. J. Hum. Genet.* **8**, 861–868. <https://doi.org/10.1038/sj.ejhg.5200552> (2000).
22. Douglas, S. P. M. *et al.* ERCC6L2 Defines a Novel Entity within Inherited Acute Myeloid Leukemia. *Blood* **133**, 2724–2728. <https://doi.org/10.1182/blood-2019-01-896233> (2019).
23. Hakkarainen, M. *et al.* The Clinical Picture of ERCC6L2 Disease: From Bone Marrow Failure to Acute Leukemia. *Blood* **141**, 2853–2866. <https://doi.org/10.1182/blood.2022019425> (2023).
24. Kurki, M. I. *et al.* FinnGen Provides Genetic Insights from a Well-Phenotyped Isolated Population. *Nature* **613**, 508–518. <https://doi.org/10.1038/s41586-022-05473-8> (2023).
25. Li, N. *et al.* Evaluation of the Association of Heterozygous Germline Variants in NTHL1 with Breast Cancer Predisposition: An International Multi-center Study of 47,180 Subjects. *NPJ Breast Cancer* **7**, 52. <https://doi.org/10.1038/s41523-021-00255-3> (2021).
26. Karczewski, K. J. *et al.* The Mutational Constraint Spectrum Quantified from Variation in 141,456 Humans. *Nature* **581**, 434–443. <https://doi.org/10.1038/s41586-020-2308-7> (2020).
27. Elsayed, F. A. *et al.* Monoallelic NTHL1 Loss-of-Function Variants and Risk of Polyposis and Colorectal Cancer. *Gastroenterology* **159**, 2241–2243. <https://doi.org/10.1053/j.gastro.2020.08.042> (2020).
28. Aspinwall, R. *et al.* Cloning and Characterization of a Functional Human Homolog of *Escherichia coli* Endonuclease III. *Proc. Natl. Acad. Sci. U. S. A.* **94**, 109–114. <https://doi.org/10.1073/pnas.94.1.109> (1997).
29. Dizdaroglu, M., Karahalil, B., Senturker, S., Buckley, T. J. & Roldan-Arjona, T. Excision of Products of Oxidative DNA Base Damage by Human NTH1 Protein. *Biochemistry* **38**, 243–246. <https://doi.org/10.1021/bi9819071> (1999).
30. Weren, R. D. *et al.* NTHL1 and MUTYH Polyposis Syndromes: Two Sides of the Same Coin?. *J. Pathol.* **244**, 135–142. <https://doi.org/10.1002/path.5002> (2018).
31. Nik-Zainal, S. *et al.* Landscape of Somatic Mutations in 560 Breast Cancer Whole-Genome Sequences. *Nature* **534**, 47–54. <https://doi.org/10.1038/nature17676> (2016).
32. Drost, J. *et al.* Use of CRISPR-Modified Human Stem Cell Organoids to Study the Origin of Mutational Signatures in Cancer. *Science* **358**, 234–238. <https://doi.org/10.1126/science.aao3130> (2017).
33. Belhadj, S. *et al.* NTHL1 Biallelic Mutations Seldom Cause Colorectal Cancer, Serrated Polyposis or a Multi-Tumor Phenotype, in Absence of Colorectal Adenomas. *Sci. Rep.* **9**, 9020. <https://doi.org/10.1038/s41598-019-45281-1> (2019).
34. Mantere, T. *et al.* Finnish Fanconi Anemia Mutations and Hereditary Predisposition to Breast and Prostate Cancer. *Clin. Genet.* **88**, 68–73. <https://doi.org/10.1111/cge.12447> (2015).
35. Dennis, J. *et al.* Rare Germline Copy Number Variants (CNVs) and Breast Cancer Risk. *Commun. Biol.* **5**, 65. <https://doi.org/10.1038/s42003-021-02990-6> (2022).
36. Tsukasaki, K. *et al.* Mutations in the Mitotic Check Point Gene, MAD1L1, in Human Cancers. *Oncogene* **20**, 3301–3305. <https://doi.org/10.1038/sj.onc.1204421> (2001).
37. Perez-Oliva, A. B. *et al.* USP45 Deubiquitylase Controls ERCC1-XPF Endonuclease-Mediated DNA Damage Responses. *EMBO J.* **34**, 326–343. <https://doi.org/10.15252/embj.201489184> (2015).
38. Waller, R. G. *et al.* Novel Pedigree Analysis Implicates DNA Repair and Chromatin Remodeling in Multiple Myeloma Risk. *PLoS Genet.* **14**, e1007111. <https://doi.org/10.1371/journal.pgen.1007111> (2018).
39. Li, H. & Durbin, R. Fast and Accurate Short Read Alignment with Burrows-Wheeler Transform. *Bioinformatics* **25**, 1754–1760. <https://doi.org/10.1093/bioinformatics/btp324> (2009).
40. Thorvaldsdottir, H., Robinson, J. T. & Mesirov, J. P. Integrative Genomics Viewer (IGV): High-Performance Genomics Data Visualization and Exploration. *Brief. Bioinform.* **14**, 178–192. <https://doi.org/10.1093/bib/bbs017> (2013).
41. Li, H. *et al.* The Sequence Alignment/Map Format and SAMtools. *Bioinformatics* **25**, 2078–2079. <https://doi.org/10.1093/bioinformatics/btp352> (2009).
42. Danecsek, P. *et al.* Twelve Years of SAMtools and BCFtools. *Gigascience* **10**, 2. <https://doi.org/10.1093/gigascience/giab008> (2021).
43. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: Functional Annotation of Genetic Variants from High-Throughput Sequencing Data. *Nucleic Acids Res.* **38**, e164. <https://doi.org/10.1093/nar/gkq603> (2010).
44. Ashburner, M. *et al.* Gene Ontology: Tool for the Unification of Biology. The Gene Ontology Consortium. *Nat. Genet.* **25**, 25–29. <https://doi.org/10.1038/75556> (2000).
45. Carbon, S. *et al.* AmiGO: Online Access to Ontology and Annotation Data. *Bioinformatics* **25**, 288–289. <https://doi.org/10.1093/bioinformatics/btn615> (2009).
46. Lek, M. *et al.* Analysis of Protein-Coding Genetic Variation in 60,706 Humans. *Nature* **536**, 285–291. <https://doi.org/10.1038/nature19057> (2016).
47. 1000 Genomes Project Consortium *et al.* A Global Reference for Human Genetic Variation. *Nature* **526**, 68–74. <https://doi.org/10.1038/nature15393> (2015).
48. Kircher, M. *et al.* A General Framework for Estimating the Relative Pathogenicity of Human Genetic Variants. *Nat. Genet.* **46**, 310–315. <https://doi.org/10.1038/ng.2892> (2014).
49. UniProt Consortium. UniProt: The Universal Protein Knowledgebase in 2023. *Nucleic Acids Res.* **51**, D523–D531. <https://doi.org/10.1093/nar/gkac1052> (2023).

50. Brown, G. R. *et al.* Gene: A gene-centered information resource at NCBI. *Nucleic Acids Res.* **43**, D36–42. <https://doi.org/10.1093/nar/gku1055> (2015).
51. Cunningham, F. *et al.* Ensembl 2022. *Nucleic Acids Res.* **50**, D988–D995. <https://doi.org/10.1093/nar/gkab1049> (2022).
52. Lefter, M. *et al.* Mutalyzer 2: Next generation HGVS nomenclature checker. *Bioinformatics* **37**, 2811–2817. <https://doi.org/10.1093/bioinformatics/btab051> (2021).
53. Higgins, J. *et al.* Verifying nomenclature of DNA variants in submitted manuscripts: Guidance for journals. *Hum. Mutat.* **42**, 3–7. <https://doi.org/10.1002/humu.24144> (2021).
54. Kiiski, J. I. *et al.* Exome sequencing identifies FANCM as a susceptibility gene for triple-negative breast cancer. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 15172–15177. <https://doi.org/10.1073/pnas.1407909111> (2014).
55. Mantere, T. *et al.* Case-control analysis of truncating mutations in DNA damage response genes connects TEX15 and FANCD2 with hereditary breast cancer susceptibility. *Sci. Rep.* **7**, 681. <https://doi.org/10.1038/s41598-017-00766-9> (2017).
56. Kinsella, R. J. *et al.* Ensembl BioMart: A hub for data retrieval across taxonomic space. *Database (Oxf.)* **2011**, 030. <https://doi.org/10.1093/database/bar030> (2011).
57. Kilpivaara, O. *et al.* Correlation of CHEK2 protein expression and c.1100delC mutation status with tumor characteristics among unselected breast cancer patients. *Int. J. Cancer* **113**, 575–580. <https://doi.org/10.1002/ijc.20638> (2005).
58. Fagerholm, R. *et al.* NAD(P)H:quinone oxidoreductase 1 NQO1\*2 genotype (P187S) is a strong prognostic and predictive factor in breast cancer. *Nat. Genet.* **40**, 844–853. <https://doi.org/10.1038/ng.155> (2008).
59. Vahteristo, P. *et al.* A CHEK2 genetic variant contributing to a substantial fraction of familial breast cancer. *Am. J. Hum. Genet.* **71**, 432–438. <https://doi.org/10.1086/341943> (2002).
60. Eerola, H., Blomqvist, C., Pukkala, E., Pyrhonen, S. & Nevanlinna, H. Familial breast cancer in southern Finland: How prevalent are breast cancer families and can we trust the family history reported by patients?. *Eur. J. Cancer* **36**, 1143–1148. [https://doi.org/10.1016/s0959-8049\(00\)00093-9](https://doi.org/10.1016/s0959-8049(00)00093-9) (2000).
61. Kaunisto, M. A. *et al.* Pain in 1,000 women treated for breast cancer: A prospective study of pain sensitivity and postoperative pain. *Anesthesiology* **119**, 1410–1421. <https://doi.org/10.1097/ALN.000000000000012> (2013).
62. R Core Team. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. <https://www.r-project.org/> (2022).
63. Mbatchou, J. *et al.* Computationally efficient whole-genome regression for quantitative and binary traits. *Nat. Genet.* **53**, 1097–1103. <https://doi.org/10.1038/s41588-021-00870-7> (2021).

## Acknowledgements

We thank the study participants for their valuable participation. We thank research nurses Outi Malkavaara and Irja Erkkilä for their help with the patient data, M.Sc. Erja Nynäs with bioinformatics and with the FinnGen data, and M.Sc. Lotta Mielikäinen with the genotyping, Dr. Taru A. Muranen for valuable advice, and Prof. Katri Pylkäs for providing a positive control sample for genotyping. We gratefully acknowledge the participants and investigators of the FinnGen study and thank Prof. Aarno Palotie and Dr. Risto Kajanne for their kind help. The FinnGen consortium members are presented in the Supplementary Table S6. The following biobanks are acknowledged for delivering biobank samples to FinnGen: Auria Biobank (<https://www.auria.fi/biopankki>), THL Biobank (<https://www.thl.fi/biobank>), Helsinki Biobank (<https://www.helsinginbiopankki.fi>), Biobank Borealis of Northern Finland (<https://www.ppsfp.fi/Tutkimus-ja-opetus/Biopankki/Pages/Biobank-Borealis-briefly-in-English.aspx>), Finnish Clinical Biobank Tampere ([https://www.tays.fi/en-US/Research\\_and\\_development/Finnish\\_Clinical\\_Biobank\\_Tampere](https://www.tays.fi/en-US/Research_and_development/Finnish_Clinical_Biobank_Tampere)), Biobank of Eastern Finland (<https://www.ita-suomenbiopankki.fi/en>), Central Finland Biobank (<https://www.ksshp.fi/fi-FI/Potilaalle/Biopankki>), Finnish Red Cross Blood Service Biobank (<https://www.veripalvelu.fi/verenluovutus/biopankkitoiminta>), Terveystalo Biobank (<https://www.terveystalo.com/fi/Yritystietoa/Terveystalo-Biopankki/Biopankki/>), and Arctic Biobank (<https://www.oulu.fi/en/university/faculties-and-units/faculty-medicine/northern-finland-birth-cohorts-and-arctic-biobank>). All Finnish Biobanks are members of BBMRI.fi infrastructure (<https://www.bbMRI.fi>). Finnish Biobank Cooperative -FINBB (<https://finbb.fi/>) is the coordinator of BBMRI-ERIC operations in Finland. The Finnish biobank data can be accessed through the Fingenious® services (<https://site.fingenious.fi/en/>) managed by FINBB. We thank the Helsinki Biobank also for providing control samples for genotyping in the current study, the staff at the McGill University and Génome Québec Innovation Centre for exome-sequencing services, the staff at the Genotyping laboratory of Institute for Molecular Medicine Finland (FIMM) Technology Centre, University of Helsinki, for variant genotyping services, and the Genome Aggregation Database (gnomAD) and the groups that provided exome and genome variant data to the database. A full list of contributing groups for the gnomAD database can be found at <https://gnomad.broadinstitute.org/about>.

## Author contributions

H.N., A.K.N., and L.M.P. designed the study. S.K. annotated and pre-processed the WES data. M.N. pre-processed the gnomAD data. A.K.N., J.I.K., L.M.P., and H.N. performed the variant selection. A.K.N., M.S., N.A., L.M.P., and J.I.K. carried out the variant genotyping. A.K.N. did the statistical analyses and processed the FinnGen data. K.A., C.B., P.H., J.S., A.K., E.K., and T.T. contributed samples and patient information. A.K.N. and H.N. wrote the manuscript. All authors read and approved the final manuscript.

## Funding

This work was supported by the Helsinki University Hospital Research Fund, the Sigrid Jusélius Foundation, and the Cancer Foundation Finland. The FinnGen project is funded by two grants from Business Finland (HUS 4685/31/2016 and UH 4386/31/2016) and the following industry partners: AbbVie Inc., AstraZeneca UK Ltd, Biogen MA Inc., Bristol Myers Squibb (and Celgene Corporation & Celgene International II Sàrl), Genentech Inc., Merck Sharp & Dohme LCC, Pfizer Inc., GlaxoSmithKline Intellectual Property Development Ltd., Sanofi US Services Inc., Maze Therapeutics Inc., Janssen Biotech Inc., Novartis AG, and Boehringer Ingelheim International GmbH.

## Competing interests

L.M.P. and M.S. are currently employed by Blueprint Genetics. The other authors declare no competing interests.

## Additional information

**Supplementary Information** The online version contains supplementary material available at <https://doi.org/10.1038/s41598-023-47441-w>.

**Correspondence** and requests for materials should be addressed to H.N.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023

## FinnGen

Aarno Palotie<sup>8,9,10</sup>, Mark Daly<sup>8,9,10</sup>, Bridget Riley-Gillis<sup>11</sup>, Howard Jacob<sup>11</sup>, Dirk Paul<sup>12</sup>, Slavé Petrovski<sup>12</sup>, Heiko Runz<sup>13</sup>, Sally John<sup>13</sup>, George Okafo<sup>14</sup>, Nathan Lawless<sup>14</sup>, Heli Salminen-Mankonen<sup>14</sup>, Robert Plenge<sup>15</sup>, Joseph Maraville<sup>15</sup>, Mark McCarthy<sup>16</sup>, Margaret G. Ehm<sup>17</sup>, Kirsi Auro<sup>18</sup>, Simonne Longrich<sup>19</sup>, Anders Mälarstig<sup>20</sup>, Katherine Klinger<sup>21</sup>, Clement Chatelain<sup>21</sup>, Matthias Gossel<sup>21</sup>, Karol Estrada<sup>22</sup>, Robert Graham<sup>22</sup>, Robert Yang<sup>23</sup>, Chris O'Donnell<sup>24</sup>, Tomi P. Mäkelä<sup>25</sup>, Jaakko Kaprio<sup>8</sup>, Petri Virolainen<sup>26</sup>, Antti Hakanen<sup>26</sup>, Terhi Kilpi<sup>27</sup>, Markus Perola<sup>27</sup>, Jukka Partanen<sup>28</sup>, Anne Pitkäranta<sup>29</sup>, Taneli Raivio<sup>29</sup>, Jani Tikkanen<sup>30</sup>, Raisa Serpi<sup>30</sup>, Tarja Laitinen<sup>31</sup>, Veli-Matti Kosma<sup>32</sup>, Jari Laukkanen<sup>33</sup>, Marco Hautalahti<sup>34</sup>, Outi Tuovila<sup>35</sup>, Raimo Pakkanen<sup>35</sup>, Jeffrey Waring<sup>11</sup>, Fedik Rahimov<sup>11</sup>, Ioanna Tachmazidou<sup>12</sup>, Chia-Yen Chen<sup>13</sup>, Zhihao Ding<sup>14</sup>, Marc Jung<sup>14</sup>, Shameek Biswas<sup>15</sup>, Rion Pendergrass<sup>16</sup>, David Pulford<sup>36</sup>, Neha Raghavan<sup>19</sup>, Adriana Huertas-Vazquez<sup>19</sup>, Jae-Hoon Sul<sup>19</sup>, Xinli Hu<sup>20</sup>, Åsa Hedman<sup>20</sup>, Manuel Rivas<sup>22,37</sup>, Dawn Waterworth<sup>38</sup>, Nicole Renaud<sup>24</sup>, Ma'en Obeidat<sup>24</sup>, Samuli Ripatti<sup>8</sup>, Johanna Schleutker<sup>26</sup>, Mikko Arvas<sup>28</sup>, Olli Carpén<sup>29</sup>, Reetta Hinttala<sup>30</sup>, Johannes Kettunen<sup>30</sup>, Arto Mannermaa<sup>32</sup>, Katriina Aalto-Setälä<sup>39</sup>, Mika Kähönen<sup>31</sup>, Johanna Mäkelä<sup>34</sup>, Reetta Kälviäinen<sup>40</sup>, Valtteri Julkunen<sup>40</sup>, Hilikka Soininen<sup>40</sup>, Anne Remes<sup>41</sup>, Mikko Hiltunen<sup>42</sup>, Jukka Peltola<sup>43</sup>, Minna Raivio<sup>44</sup>, Pentti Tienari<sup>44</sup>, Juha Rinne<sup>45</sup>, Roosa Kallionpää<sup>45</sup>, Juulia Partanen<sup>8</sup>, Ali Abbasi<sup>11</sup>, Adam Ziemann<sup>11</sup>, Nizar Smaoui<sup>11</sup>, Anne Lehtonen<sup>11</sup>, Susan Eaton<sup>13</sup>, Sanni Lahdenperä<sup>13</sup>, Natalie Bowers<sup>16</sup>, Edmond Teng<sup>16</sup>, Fanli Xu<sup>46</sup>, Laura Addis<sup>46</sup>, John Eicher<sup>46</sup>, Qingqin S. Li<sup>47</sup>, Karen He<sup>38</sup>, Ekaterina Khramtsova<sup>38</sup>, Martti Färkkilä<sup>44</sup>, Jukka Koskela<sup>44</sup>, Sampsa Pikkarainen<sup>44</sup>, Airi Jussila<sup>43</sup>, Katri Kaukinen<sup>43</sup>, Timo Blomster<sup>41</sup>, Mikko Kiviniemi<sup>40</sup>, Markku Voutilainen<sup>45</sup>, Tim Lu<sup>16</sup>, Linda McCarthy<sup>46</sup>, Amy Hart<sup>38</sup>, Meijian Guan<sup>38</sup>, Jason Miller<sup>19</sup>, Kirsi Kalpala<sup>20</sup>, Melissa Miller<sup>20</sup>, Kari Eklund<sup>44</sup>, Antti Palomäki<sup>45</sup>, Pia Isomäki<sup>43</sup>, Laura Pirilä<sup>45</sup>, Oili Kaipainen-Seppänen<sup>40</sup>, Johanna Huhtakangas<sup>41</sup>, Nina Mars<sup>8</sup>, Apinya Lertratanakul<sup>11</sup>, Coralie Viollet<sup>12</sup>, Marla Hochfeld<sup>15</sup>, Jorge Esparza Gordillo<sup>46</sup>, Fabiana Farias<sup>19</sup>, Nan Bing<sup>20</sup>, Margit Pelkonen<sup>40</sup>, Paula Kauppi<sup>44</sup>, Hannu Kankaanranta<sup>48,49,50</sup>, Terttu Harju<sup>41</sup>, Riitta Lahesmaa<sup>45</sup>, Hubert Chen<sup>16</sup>, Joanna Betts<sup>46</sup>, Rajashree Mishra<sup>46</sup>, Majd Mouded<sup>51</sup>, Debby Ngo<sup>51</sup>, Teemu Niiranen<sup>52</sup>, Felix Vaura<sup>52</sup>, Veikko Salomaa<sup>52</sup>, Kaj Metsärinne<sup>45</sup>, Jenni Aittokallio<sup>45</sup>, Jussi Hernesniemi<sup>43</sup>, Daniel Gordin<sup>44</sup>, Juha Sinisalo<sup>44</sup>, Marja-Riitta Taskinen<sup>44</sup>, Tiinamaija Tuomi<sup>44</sup>, Timo Hiltunen<sup>44</sup>, Amanda Elliott<sup>8,9,10</sup>, Mary Pat Reeve<sup>8</sup>, Sanni Ruotsalainen<sup>8</sup>, Audrey Chu<sup>46</sup>, Dermot Reilly<sup>53</sup>, Mike Mendelson<sup>54</sup>, Jaakko Parkkinen<sup>20</sup>, Tuomo Meretoja<sup>44</sup>, Heikki Joensuu<sup>44</sup>, Johanna Mattson<sup>44</sup>, Eveliina Salminen<sup>44</sup>, Annika Auranen<sup>43</sup>, Peeter Karihtala<sup>41</sup>, Päivi Auvinen<sup>40</sup>, Klaus Elenius<sup>45</sup>, Esa Pitkänen<sup>8</sup>, Relja Popovic<sup>11</sup>, Margarete Fabre<sup>12</sup>, Jennifer Schutzman<sup>16</sup>, Diptee Kulkarni<sup>46</sup>, Alessandro Porello<sup>38</sup>, Andrey Loboda<sup>19</sup>, Heli Lehtonen<sup>20</sup>, Stefan McDonough<sup>20</sup>, Sauli Vuoti<sup>55</sup>, Kai Kaarniranta<sup>40,56</sup>, Joni A. Turunen<sup>57,58</sup>, Terhi Ollila<sup>44</sup>, Hannu Uusitalo<sup>43</sup>, Juha Karjalainen<sup>8</sup>,



Mengzhen Liu<sup>11</sup>, Stephanie Loomis<sup>13</sup>, Erich Strauss<sup>16</sup>, Hao Chen<sup>16</sup>, Kaisa Tasanen<sup>41</sup>, Laura Huilaja<sup>41</sup>, Katariina Hannula-Jouppi<sup>44</sup>, Tea Salmi<sup>43</sup>, Sirkku Peltonen<sup>45</sup>, Leena Koulu<sup>45</sup>, David Choy<sup>16</sup>, Ying Wu<sup>20</sup>, Pirkko Pussinen<sup>44</sup>, Aino Salminen<sup>44</sup>, Tuula Salo<sup>44</sup>, David Rice<sup>44</sup>, Pekka Nieminen<sup>44</sup>, Ulla Palotie<sup>44</sup>, Maria Siponen<sup>40</sup>, Liisa Suominen<sup>40</sup>, Päivi Mäntylä<sup>40</sup>, Ulvi Gursoy<sup>45</sup>, Vuokko Anttonen<sup>41</sup>, Kirsi Sipilä<sup>59,60</sup>, Hannele Laivuori<sup>8</sup>, Venla Kurra<sup>43</sup>, Laura Kotaniemi-Talonen<sup>43</sup>, Oskari Heikinheimo<sup>44</sup>, Ilkka Kalliala<sup>44</sup>, Lauri Aaltonen<sup>44</sup>, Varpu Jokimaa<sup>45</sup>, Marja Väärasmäki<sup>41</sup>, Outi Uimari<sup>41</sup>, Laure Morin-Papunen<sup>41</sup>, Maarit Niinimäki<sup>41</sup>, Terhi Piltonen<sup>41</sup>, Katja Kivinen<sup>8</sup>, Elisabeth Widen<sup>8</sup>, Taru Tukiainen<sup>8</sup>, Niko Välimäki<sup>61</sup>, Eija Laakkonen<sup>62</sup>, Jaakko Tyrmi<sup>48,63</sup>, Heidi Silven<sup>63</sup>, Eeva Sliz<sup>63</sup>, Riikka Arffman<sup>63</sup>, Susanna Savukoski<sup>63</sup>, Triin Laisk<sup>64</sup>, Natalia Pujol<sup>64</sup>, Janet Kumar<sup>17</sup>, Iiris Hovatta<sup>61</sup>, Erkki Isometsä<sup>44</sup>, Hanna Ollila<sup>8</sup>, Jaana Suvisaari<sup>52</sup>, Thomas Damm Als<sup>65</sup>, Antti Mäkitie<sup>66</sup>, Argyro Bizaki-Vallaskangas<sup>43</sup>, Sanna Toppila-Salmi<sup>61</sup>, Tytti Willberg<sup>45</sup>, Elmo Saarentaus<sup>8</sup>, Antti Aarnisalo<sup>44</sup>, Elisa Rahikkala<sup>41</sup>, Kristiina Aittomäki<sup>67</sup>, Fredrik Åberg<sup>68</sup>, Mitja Kurki<sup>8,9</sup>, Aki Havulinna<sup>8,52</sup>, Juha Mehtonen<sup>8</sup>, Priit Palta<sup>8</sup>, Shabbeer Hassan<sup>8</sup>, Pietro Della Briotta Parolo<sup>8</sup>, Wei Zhou<sup>9</sup>, Mutaamba Maasha<sup>9</sup>, Susanna Lemmelä<sup>8</sup>, Aoxing Liu<sup>8</sup>, Arto Lehisto<sup>8</sup>, Andrea Ganna<sup>8</sup>, Vincent Llorens<sup>8</sup>, Henrike Heyne<sup>8</sup>, Joel Rämö<sup>8</sup>, Rodos Rodosthenous<sup>8</sup>, Satu Strausz<sup>8</sup>, Tuula Palotie<sup>44,61</sup>, Kimmo Palin<sup>61</sup>, Javier Gracia-Tabuenca<sup>48</sup>, Harri Siirtola<sup>48</sup>, Tuomo Kiiskinen<sup>8</sup>, Jiwoo Lee<sup>8,9</sup>, Kristin Tsuo<sup>8,9</sup>, Kati Kristiansson<sup>27</sup>, Kati Hyvärinen<sup>69</sup>, Jarmo Ritari<sup>69</sup>, Katri Pylkäs<sup>63</sup>, Minna Karjalainen<sup>63</sup>, Tuomo Mantere<sup>30</sup>, Eeva Kangasniemi<sup>31</sup>, Sami Heikkinen<sup>42</sup>, Nina Pitkänen<sup>26</sup>, Samuel Lessard<sup>21</sup>, Clément Chatelain<sup>21</sup>, Lila Kallio<sup>26</sup>, Tiina Wahlfors<sup>27</sup>, Eero Punkka<sup>29</sup>, Sanna Siltanen<sup>31</sup>, Teijo Kuopio<sup>33</sup>, Anu Jalanko<sup>8</sup>, Hwei-Yi Shen<sup>8</sup>, Risto Kajanne<sup>8</sup>, Mervi Aavikko<sup>8</sup>, Helen Cooper<sup>8</sup>, Denise Öller<sup>8</sup>, Rasko Leinonen<sup>8,70</sup>, Henna Palin<sup>31</sup>, Malla-Maria Linna<sup>29</sup>, Masahiro Kanai<sup>9</sup>, Zhili Zheng<sup>9</sup>, L. Elisa Lahtela<sup>8</sup>, Mari Kaunisto<sup>8</sup>, Elina Kilpeläinen<sup>8</sup>, Timo P. Sipilä<sup>8</sup>, Oluwaseun Alexander Dada<sup>8</sup>, Awaisa Ghazal<sup>8</sup>, Anastasia Kytölä<sup>8</sup>, Rigbe Weldatsadik<sup>8</sup>, Kati Donner<sup>8</sup>, Anu Loukola<sup>29</sup>, Päivi Laiho<sup>27</sup>, Tuuli Sistonen<sup>27</sup>, Essi Kaiharju<sup>27</sup>, Markku Laukkanen<sup>27</sup>, Elina Järvensivu<sup>27</sup>, Sini Lähteenmäki<sup>27</sup>, Lotta Männikkö<sup>27</sup>, Regis Wong<sup>27</sup>, Auli Toivola<sup>27</sup>, Minna Brunfeldt<sup>27</sup>, Hannele Mattsson<sup>27</sup>, Sami Koskelainen<sup>27</sup>, Tero Hiekkalinna<sup>27</sup>, Teemu Paajanen<sup>27</sup>, Kalle Pärn<sup>8</sup>, Mart Kals<sup>8</sup>, Shuang Luo<sup>8</sup>, Shanmukha Sampath Padmanabhuni<sup>8</sup>, Marianna Niemi<sup>48</sup>, Mika Helminen<sup>48</sup>, Tiina Luukkaala<sup>48</sup>, Iida Vähätalo<sup>48</sup>, Jyrki Tammerluoto<sup>8</sup>, Sarah Smith<sup>34</sup>, Tom Southerington<sup>34</sup> & Petri Lehto<sup>34</sup>

<sup>8</sup>Institute for Molecular Medicine Finland (FIMM), HiLIFE, University of Helsinki, Helsinki, Finland. <sup>9</sup>Broad Institute of MIT and Harvard, Cambridge, MA, USA. <sup>10</sup>Massachusetts General Hospital, Boston, MA, USA. <sup>11</sup>AbbVie, Chicago, IL, USA. <sup>12</sup>AstraZeneca, Cambridge, UK. <sup>13</sup>Biogen, Cambridge, MA, USA. <sup>14</sup>Boehringer Ingelheim, Ingelheim am Rhein, Germany. <sup>15</sup>Bristol Myers Squibb, New York, NY, USA. <sup>16</sup>Genentech, San Francisco, CA, USA. <sup>17</sup>GlaxoSmithKline, Collegeville, PA, USA. <sup>18</sup>GlaxoSmithKline, Espoo, Finland. <sup>19</sup>Merck, Kenilworth, NJ, USA. <sup>20</sup>Pfizer, New York, NY, USA. <sup>21</sup>Translational Sciences, Sanofi R&D, Framingham, MA, USA. <sup>22</sup>Maze Therapeutics, San Francisco, CA, USA. <sup>23</sup>Janssen Biotech, Beerse, Belgium. <sup>24</sup>Novartis Institutes for BioMedical Research, Cambridge, MA, USA. <sup>25</sup>HiLIFE, University of Helsinki, Helsinki, Finland. <sup>26</sup>Auria Biobank, University of Turku and Hospital District of Southwest Finland, Turku, Finland. <sup>27</sup>THL Biobank, Finnish Institute for Health and Welfare (THL), Helsinki, Finland. <sup>28</sup>Finnish Red Cross Blood Service and Finnish Hematology Registry and Clinical Biobank, Helsinki, Finland. <sup>29</sup>Helsinki Biobank, University of Helsinki and Hospital District of Helsinki and Uusimaa, Helsinki, Finland. <sup>30</sup>Northern Finland Biobank Borealis, University of Oulu and Northern Ostrobothnia Hospital District, Oulu, Finland. <sup>31</sup>Finnish Clinical Biobank Tampere, Tampere University and Pirkanmaa Hospital District, Tampere, Finland. <sup>32</sup>Biobank of Eastern Finland, University of Eastern Finland and Northern Savo Hospital District, Kuopio, Finland. <sup>33</sup>Central Finland Biobank, University of Jyväskylä and Central Finland Health Care District, Jyväskylä, Finland. <sup>34</sup>FINBB, Finnish Biobank Cooperative, Helsinki, Finland. <sup>35</sup>Business Finland, Helsinki, Finland. <sup>36</sup>GlaxoSmithKline, Stevenage, UK. <sup>37</sup>University of Stanford, Stanford, CA, USA. <sup>38</sup>Janssen Research & Development, LLC, Spring House, PA, USA. <sup>39</sup>Faculty of Medicine and Health Technology, Tampere University, Tampere, Finland. <sup>40</sup>Northern Savo Hospital District, Kuopio, Finland. <sup>41</sup>Northern Ostrobothnia Hospital District, Oulu, Finland. <sup>42</sup>University of Eastern Finland, Kuopio, Finland. <sup>43</sup>Pirkanmaa Hospital District, Tampere, Finland. <sup>44</sup>Hospital District of Helsinki and Uusimaa, Helsinki, Finland. <sup>45</sup>Hospital District of Southwest Finland, Turku, Finland. <sup>46</sup>GlaxoSmithKline, Brentford, UK. <sup>47</sup>Janssen Research & Development, LLC, Titusville, NJ, USA. <sup>48</sup>Tampere University, Tampere, Finland. <sup>49</sup>Seinäjäki Central Hospital, Seinäjoki, Finland. <sup>50</sup>University of Gothenburg, Gothenburg, Sweden. <sup>51</sup>Novartis, Basel, Switzerland. <sup>52</sup>Finnish Institute for Health and Welfare (THL), Helsinki, Finland. <sup>53</sup>Janssen Research & Development, LLC, Boston, MA, USA. <sup>54</sup>Novartis, Boston, MA, USA. <sup>55</sup>Janssen-Cilag Oy, Espoo, Finland. <sup>56</sup>Department of Molecular Genetics, University of Lodz, Lodz, Poland. <sup>57</sup>Helsinki University Hospital and University of Helsinki, Helsinki, Finland. <sup>58</sup>Eye Genetics Group, Folkhälsan Research Center, Helsinki, Finland. <sup>59</sup>Research Unit of Oral Health Sciences, Faculty of Medicine, University of Oulu, Oulu, Finland. <sup>60</sup>Medical Research Center Oulu, Oulu University Hospital and University of Oulu, Oulu, Finland. <sup>61</sup>University of Helsinki,



Helsinki, Finland. <sup>62</sup>University of Jyväskylä, Jyväskylä, Finland. <sup>63</sup>University of Oulu, Oulu, Finland. <sup>64</sup>Estonian Biobank, Tartu, Estonia. <sup>65</sup>Aarhus University, Aarhus, Denmark. <sup>66</sup>Department of Otorhinolaryngology-Head and Neck Surgery, University of Helsinki and Helsinki University Hospital, Helsinki, Finland. <sup>67</sup>Department of Medical Genetics, Helsinki University Central Hospital, Helsinki, Finland. <sup>68</sup>Transplantation and Liver Surgery Clinic, Helsinki University Hospital and University of Helsinki, Helsinki, Finland. <sup>69</sup>Finnish Red Cross Blood Service, Helsinki, Finland. <sup>70</sup>European Molecular Biology Laboratory, European Bioinformatics Institute, Cambridge, UK.