

# Genome-wide analyses using UK Biobank data provide insights into the genetic architecture of osteoarthritis

Eleni Zengini<sup>1,2,19</sup>, Konstantinos Hatzikotoulas<sup>3,19</sup>, Ioanna Tachmazidou<sup>3,4,19</sup>, Julia Steinberg<sup>3,5</sup>, Fernando P. Hartwig<sup>6,7</sup>, Lorraine Southam<sup>3,8</sup>, Sophie Hackinger<sup>3</sup>, Cindy G. Boer<sup>9</sup>, Unnur Styrkarsdottir<sup>10</sup>, Arthur Gilly<sup>3</sup>, Daniel Suveges<sup>3</sup>, Britt Killian<sup>3</sup>, Thorvaldur Ingvarsson<sup>11,12,13</sup>, Helgi Jonsson<sup>12,14</sup>, George C. Babis<sup>15</sup>, Andrew McCaskie<sup>16</sup>, Andre G. Uitterlinden<sup>9</sup>, Joyce B. J. van Meurs<sup>9</sup>, Unnur Thorsteinsdottir<sup>10,12</sup>, Kari Stefansson<sup>10,12</sup>, George Davey Smith<sup>7,17,18</sup>, Jeremy M. Wilkinson<sup>1</sup> and Eleftheria Zeggini<sup>1</sup><sup>3\*</sup>

**Osteoarthritis is a common complex disease imposing a large public-health burden. Here, we performed a genome-wide association study for osteoarthritis, using data across 16.5 million variants from the UK Biobank resource. After performing replication and meta-analysis in up to 30,727 cases and 297,191 controls, we identified nine new osteoarthritis loci, in all of which the most likely causal variant was noncoding. For three loci, we detected association with biologically relevant radiographic endophenotypes, and in five signals we identified genes that were differentially expressed in degraded compared with intact articular cartilage from patients with osteoarthritis. We established causal effects on osteoarthritis for higher body mass index but not for triglyceride levels or genetic predisposition to type 2 diabetes.**

Osteoarthritis is the most prevalent musculoskeletal disease and the most common form of arthritis<sup>1</sup>. The hallmarks of osteoarthritis are degeneration of articular cartilage, remodeling of the underlying bone and synovitis<sup>2</sup>. A leading cause of disability worldwide, osteoarthritis affects 40% of individuals over the age of 70 and is associated with an elevated risk of comorbidity and death<sup>3</sup>. The rising health economic burden of osteoarthritis is commensurate with rising longevity and obesity rates, and there is currently no curative therapy. The heritability of osteoarthritis is ~50%, and previous genetic studies have identified 21 loci in total, traversing hip, knee and hand osteoarthritis with limited overlap<sup>3</sup>. Here, we conducted a large osteoarthritis genome-wide association study (GWAS), using genotype data across 16.5 million variants from UK Biobank. We defined osteoarthritis on the basis of both self-reported status and linkage to Hospital Episode Statistics data, as well as the joint specificity of the disease (knee and/or hip) (Supplementary Fig. 1).

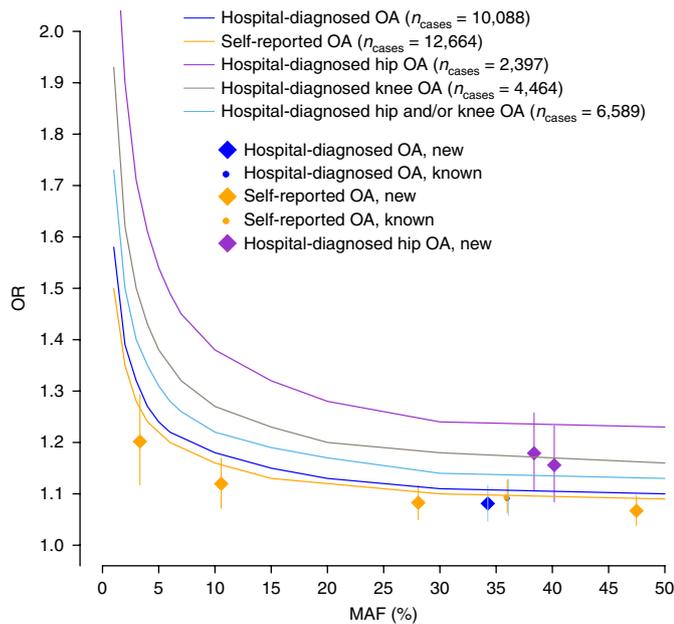
## Results

**Disease definition and power to detect genetic associations.** We compared and contrasted the hospital-diagnosed ( $n = 10,083$  cases)

and self-reported ( $n = 12,658$  cases) osteoarthritis GWAS drawn from the same UK Biobank dataset (with selection of approximately four times more nonosteoarthritis controls than cases to preserve power for common alleles while avoiding case-control imbalance that might cause association tests to misbehave for low-frequency variants<sup>4</sup>) (Supplementary Tables 1–3, Supplementary Figs. 2–4 and Methods). We found power advantages with the self-reported dataset, thus indicating that the higher sample size overcame the limitations associated with phenotype uncertainty. When evaluating the accuracy of disease definition, we found that self-reported osteoarthritis had a modest positive predictive value (PPV; 30%) and sensitivity (37%), but high negative predictive value (95%) and specificity, correctly identifying 93% of individuals who did not have osteoarthritis (Supplementary Table 4). In terms of power to detect genetic associations, the self-reported-osteoarthritis dataset had clear advantages commensurate with its larger sample size (Fig. 1). For example, for a representative complex-disease-associated variant with a minor allele frequency (MAF) of 30% and an allelic odds ratio (OR) of 1.10, the self-reported and hospital-diagnosed osteoarthritis analyses had 80% and 56% power, respectively,

<sup>1</sup>Department of Oncology and Metabolism, University of Sheffield, Sheffield, UK. <sup>2</sup>5th Psychiatric Department, Dromokaiteio Psychiatric Hospital, Athens, Greece. <sup>3</sup>Human Genetics, Wellcome Trust Sanger Institute, Hinxton, UK. <sup>4</sup>GSK, R&D Target Sciences, Medicines Research Centre, Stevenage, UK.

<sup>5</sup>Cancer Research Division, Cancer Council NSW, Sydney, New South Wales, Australia. <sup>6</sup>Postgraduate Program in Epidemiology, Federal University of Pelotas, Pelotas, Brazil. <sup>7</sup>Medical Research Council Integrative Epidemiology Unit, University of Bristol, Bristol, UK. <sup>8</sup>Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford, UK. <sup>9</sup>Departments of Internal Medicine and Epidemiology, Erasmus MC, Rotterdam, the Netherlands. <sup>10</sup>deCODE genetics/Amgen, Reykjavik, Iceland. <sup>11</sup>Department of Orthopaedic Surgery, Akureyri Hospital, Akureyri, Iceland. <sup>12</sup>Faculty of Medicine, University of Iceland, Reykjavik, Iceland. <sup>13</sup>Institution of Health Science, University of Akureyri, Akureyri, Iceland. <sup>14</sup>Department of Medicine, Landspítali, National University Hospital of Iceland, Reykjavik, Iceland. <sup>15</sup>2nd Department of Orthopaedic Surgery, Konstantopouleio General Hospital, National and Kapodistrian University of Athens, Athens, Greece. <sup>16</sup>Division of Trauma & Orthopaedic Surgery, Department of Surgery, University of Cambridge, Addenbrooke's Hospital, Cambridge, UK. <sup>17</sup>Population Health Sciences, University of Bristol, Bristol, UK. <sup>18</sup>National Institute for Health Research, Bristol Biomedical Research Centre, University Hospitals Bristol, NHS Foundation Trust and University of Bristol, Bristol, UK. <sup>19</sup>These authors contributed equally: Eleni Zengini, Konstantinos Hatzikotoulas and Ioanna Tachmazidou. \*e-mail: [eleftheria@sanger.ac.uk](mailto:eleftheria@sanger.ac.uk)



**Fig. 1 | Power to detect association in the discovery stage.** OR and 95% CI values are shown as a function of MAF. Diamonds, newly reported variants; circles, known variants. The curves indicate 80% power at the genome-wide-significance threshold of  $P < 5.0 \times 10^{-8}$  for the number of cases and controls of each trait at the discovery stage (likelihood ratio test). OA, osteoarthritis.

to detect an effect at genome-wide significance (i.e.,  $P < 5.0 \times 10^{-8}$ ; Supplementary Table 5).

We found nominally significant evidence of concordance between the direction of effect at previously reported osteoarthritis loci and the discovery analyses for hospital-diagnosed-osteoarthritis definitions (Supplementary Tables 6 and 7, and Supplementary Note), thus indicating that a narrower definition of disease may provide better effect-size estimates despite being limited by power to identify robust statistical evidence of association.

**Heritability estimates across osteoarthritis definitions.** We found that common-frequency variants explained 12% of osteoarthritis heritability when using self-reported status and explained 16% of osteoarthritis heritability when using hospital records (19% of hip-osteoarthritis and 15% of knee-osteoarthritis heritability) (Supplementary Table 8). The heritability estimates from self-reported and hospital records were not significantly different (Supplementary Table 9). The concordance between self-reported and hospital-diagnosed osteoarthritis was further substantiated by the high genetic-correlation estimate of the two disease definitions (87%,  $P = 3.14 \times 10^{-53}$ ) (Supplementary Table 10). We found strong genome-wide correlation between hip osteoarthritis and knee osteoarthritis (88%,  $P = 1.96 \times 10^{-6}$ ), even though the previously reported osteoarthritis loci are predominantly not shared between the two osteoarthritis joint sites. From this new observation of a substantial shared genetic etiology, we sought replication of association signals across joint sites.

**Identification of novel osteoarthritis loci.** We used 173 variants with  $P < 1.0 \times 10^{-5}$  and  $\text{MAF} > 0.01$  for replication in an Icelandic cohort of up to 18,069 cases and 246,293 controls (Supplementary Fig. 1, Supplementary Tables 11–15 and Methods). Given the number of variants, the replication significance threshold was  $P < 2.9 \times 10^{-4}$ . After meta-analysis in up to 30,727 cases and 297,191 controls, we identified six genome-wide-significant associations at

novel loci and three further replicating signals just below the corrected genome-wide-significance threshold (Table 1 and Fig. 2).

We identified association between rs2521349 and hip osteoarthritis (OR 1.13 (95% confidence interval (CI) 1.09–1.17),  $P = 9.95 \times 10^{-10}$ , effect-allele frequency (EAF) 0.37). rs2521349 resides in an intron of *MAP2K6* on chromosome 17. *MAP2K6* encodes an essential component of the p38 MAP kinase-mediated signal-transduction pathway, which is involved in various cellular processes in bone, muscle, fat-tissue homeostasis and differentiation<sup>5</sup>. The MAPK signaling pathway is closely associated with osteoblast differentiation<sup>6</sup>, chondrocyte apoptosis and necrosis<sup>7</sup>, and has been reported to be differentially expressed in osteoarthritis synovial-tissue samples<sup>6–12</sup>. In animal-model studies, p38 MAP kinase activity has been found to be important in maintaining cartilage health, and it has been proposed as a potential osteoarthritis diagnosis and treatment target<sup>10,13,14</sup>.

rs11780978 on chromosome 8 is also associated with hip osteoarthritis with a similar effect size (OR 1.13 (95% CI 1.08–1.17),  $P = 1.98 \times 10^{-9}$ , EAF 0.39). This variant is located in the intronic region of *PLEC* (plectin gene). We found rs11780978 to be nominally associated with the radiographically derived endophenotype of minimal joint-space width ( $\beta -0.0291$ , s.e.m. 0.0129,  $P = 0.024$ ) (Table 2 and Methods). The direction of the effect was consistent with the established clinical association between joint-space narrowing and osteoarthritis. *PLEC* encodes plectin, a structural protein that interlinks components of the cytoskeleton. Functional studies in mice have shown an effect on skeletal-muscle tissue correlated with low body weight, small size and slow postnatal growth<sup>15</sup>.

rs2820436, an intergenic variant located 24 kb upstream of the long-noncoding-RNA gene *RP11-392O17.1* and 142 kb downstream of *ZC3H11B* (zinc-finger CCCH-type containing 11B pseudogene), is associated with osteoarthritis across any joint site (OR 0.93 (95% CI 0.91–0.95),  $P = 2.01 \times 10^{-9}$ , EAF 0.65). It also resides within a region with multiple metabolic- and anthropometric-trait-associated variants, with which it was found to correlate ( $r^2$  0.18–0.88).

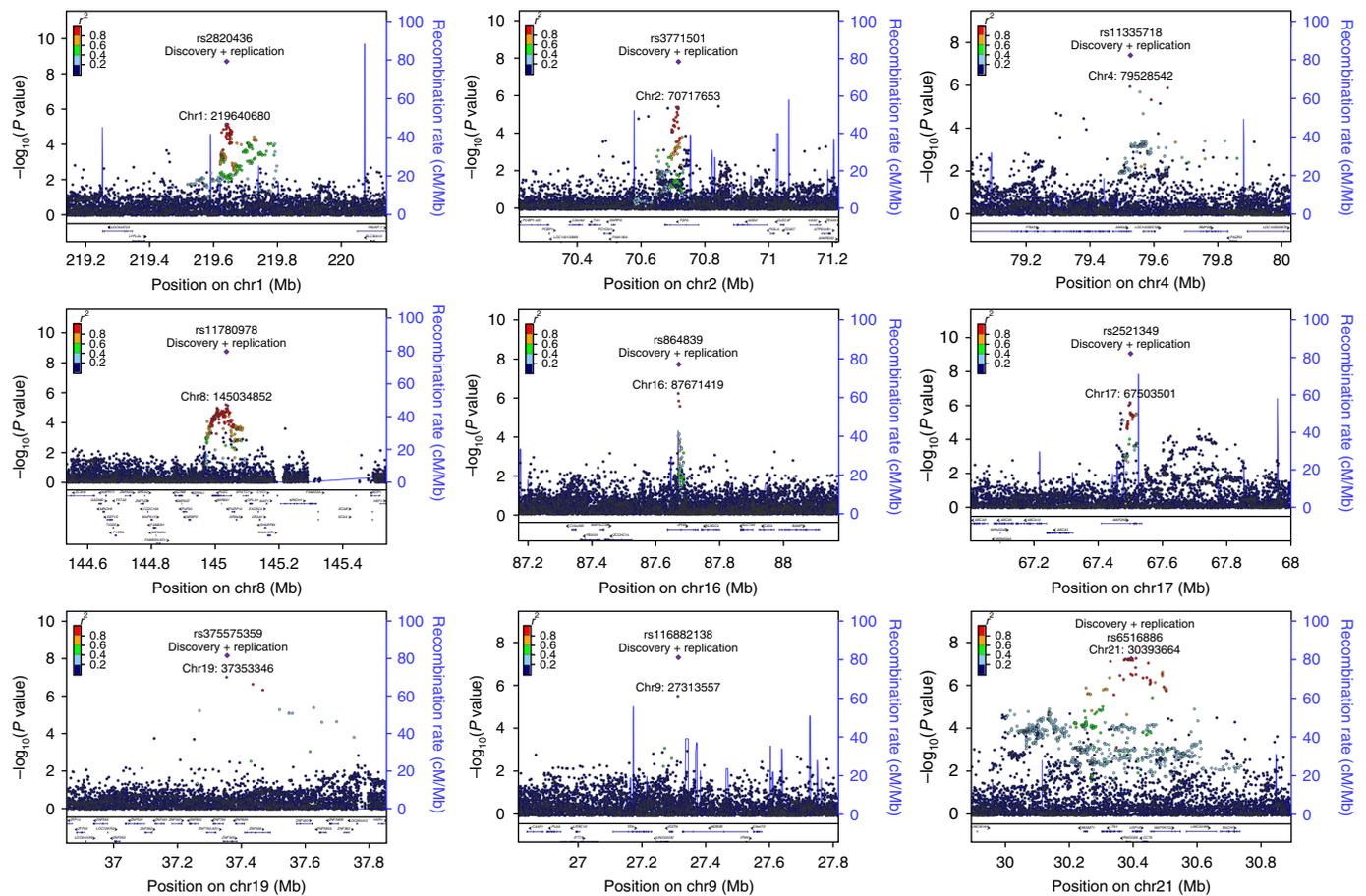
rs375575359 resides in an intron of *ZNF345* (zinc-finger-protein 345 gene) on chromosome 19. It was prioritized on the basis of osteoarthritis at any joint site and was more strongly associated with knee osteoarthritis in the replication dataset (OR 1.21 (95% CI 1.14–1.30),  $P = 7.54 \times 10^{-9}$ , EAF 0.04). Similarly, rs11335718 on chromosome 4 was associated with osteoarthritis in the discovery stage and with knee osteoarthritis in the replication stage (OR 1.11 (95% CI 1.07–1.16),  $P = 4.26 \times 10^{-8}$ , EAF 0.10). We note that Bonferroni correction for the effective number of traits tested caused rs11335718 to no longer reach genome-wide significance, with a meta-analysis  $P = 4.26 \times 10^{-8}$ . rs11335718 is an intronic variant in *ANXA3*, the annexin A3 gene. Through meta-analysis of the any-site-osteoarthritis phenotype across the discovery and replication datasets, we determined  $P = 2.6 \times 10^{-5}$  and  $P = 1.32 \times 10^{-7}$  for rs375575359 and rs11335718, respectively (Supplementary Table 11). A recent mouse-model study supports the involvement of expression of a similar-motif zinc-finger-protein (ZFP36L1) with osteoblastic differentiation<sup>16</sup>.

rs3771501 (OR 0.94 (95% CI 0.92–0.96),  $P = 1.66 \times 10^{-8}$ , EAF 0.53) is associated with osteoarthritis at any site and resides in an intron of *TGFA* (transforming growth factor alpha gene). *TGFA* encodes an epidermal-growth-factor-receptor ligand and is an important integrator of cellular signaling and function. We detected association of rs3771501 with minimal joint-space width ( $\beta -0.0699$ , s.e.m. 0.0127,  $P = 3.45 \times 10^{-8}$ ) (Table 2 and Methods); i.e., the osteoarthritis-risk-increasing allele was also associated with lower joint-cartilage thickness in humans. A perfectly correlated variant in this gene has previously been associated with cartilage thickness and with hip osteoarthritis; moreover, this variant has been found to be differentially expressed in osteoarthritis cartilage lesions compared with nonlesioned cartilage<sup>17</sup>. Functional

**Table 1 | Association summary statistics for the nine signals**

rsID	EA	Discovery phenotype	Discovery coverage EAF	Dis-covery OR	Dis-covery OR, lower 95% CI	Discovery OR, upper 95% CI	Discovery P value <sup>c</sup>	Discovery n (cases/controls)	Discovery imputation-accuracy score <sup>a</sup>	Replication phenotype	Repli-cation OR	Repli-cation OR, lower 95% CI	Repli-cation OR, upper 95% CI	Replication P value <sup>c</sup>	Replication n cases/controls	Replication imputation-accuracy score <sup>a</sup>	Overall OR	Overall OR, lower 95% CI	Overall OR, upper 95% CI	Overall P value <sup>c</sup>	Heterogeneity P value <sup>b</sup>	Overall n (cases/controls)
rs2820436	C	Hospital-diagnosed osteoarthritis	0.66	0.92	0.9	0.96	$6.45 \times 10^{-6}$	10083/40,425	Directly typed	Osteoarthritis at any site	0.64	0.91	0.97	$8.71 \times 10^{-5}$	18,069/246,293	0.99972	0.93	0.91	0.96	$2.01 \times 10^{-3}$	0.5739	28,152/286,718
rs3771501	G	Self-reported osteoarthritis	0.53	0.94	0.91	0.96	$3.81 \times 10^{-6}$	12,658/50,898	0.991707	Osteoarthritis at any site	0.54	0.92	0.98	0.001069	18,069/246,293	0.999808	0.94	0.92	0.96	$1.66 \times 10^{-8}$	0.4825	30,727/297,191
rs11335718	A	Self-reported osteoarthritis	0.11	1.12	1.07	1.17	$1.12 \times 10^{-6}$	12,658/50,898	0.968932	Knee osteoarthritis	0.11	1.1	1.2	0.014675	4,672/172,791	0.998899	1.11	1.07	1.16	$4.26 \times 10^{-8}$	0.792	17,330/223,689
rs11335718	A	Self-reported osteoarthritis	0.11	1.12	1.07	1.17	$1.12 \times 10^{-6}$	12,658/50,898	0.968932	Osteoarthritis at any site	0.11	1.06	1.11	0.013023	18,069/246,293	0.998899	1.09	1.06	1.13	$1.32 \times 10^{-7}$	0.1254	30,727/297,191
rs11780978	A	Hospital-diagnosed hip osteoarthritis	0.4	1.16	1.08	1.23	$6.24 \times 10^{-6}$	2,396/9,593	0.983752	Hip osteoarthritis	0.39	1.11	1.16	$4.55 \times 10^{-5}$	5,714/199,421	0.999673	1.13	1.08	1.17	$1.98 \times 10^{-3}$	0.2424	8,110/209,014
rs116882138	A	Hospital-diagnosed hip and/or knee osteoarthritis	0.02	1.4	1.22	1.6	$2.96 \times 10^{-6}$	6,586/26,384	Directly typed	Knee osteoarthritis	0.02	1.27	1.5	0.006552	4,672/172,791	0.998087	1.34	1.21	1.49	$5.09 \times 10^{-8}$	0.3988	11,258/199,175
rs116882138	A	Hospital-diagnosed hip and/or knee osteoarthritis	0.02	1.4	1.22	1.6	$2.96 \times 10^{-6}$	6,586/26,384	Directly typed	Hip and/or knee osteoarthritis	0.02	1.13	1.29	0.069018	9,429/199,421	0.998087	1.25	1.14	1.38	$2.93 \times 10^{-4}$	0.03456	16,015/222,805
rs2521349	A	Hospital-diagnosed hip osteoarthritis	0.38	1.18	1.11	1.26	$6.85 \times 10^{-7}$	2,396/9,593	0.996479	Hip osteoarthritis	0.37	1.1	1.16	0.000103	5,714/199,421	0.999925	1.13	1.09	1.18	$9.95 \times 10^{-10}$	0.1151	8,110/209,014
rs864839	T	Self-reported osteoarthritis	0.72	1.08	1.05	1.12	$6.21 \times 10^{-7}$	12,658/50,898	0.97115	Hip osteoarthritis	0.7	1.07	1.13	0.008275	5,714/199,421	0.997756	1.08	1.05	1.11	$2.01 \times 10^{-8}$	0.7886	18,372/250,319
rs864839	T	Self-reported osteoarthritis	0.72	1.08	1.05	1.12	$6.21 \times 10^{-7}$	12,658/50,898	0.97115	Osteoarthritis at any site	0.7	1.02	1.06	0.18218	18,069/246,293	0.997756	1.05	1.03	1.08	$7.02 \times 10^{-4}$	0.0121	30,727/297,191
rs375575359	C	Self-reported osteoarthritis	0.03	1.2	1.12	1.29	$9.96 \times 10^{-8}$	12,658/50,898	0.823533	Knee osteoarthritis	0.05	1.15	1.29	0.025177	4,672/172,791	0.992996	1.21	1.14	1.3	$7.54 \times 10^{-9}$	0.25	17,330/223,689
rs375575359	C	Self-reported osteoarthritis	0.03	1.2	1.12	1.29	$9.96 \times 10^{-8}$	12,658/50,898	0.823533	Osteoarthritis at any site	0.05	1.03	1.1	0.47234	18,069/246,293	0.992996	1.12	1.06	1.18	$2.6 \times 10^{-5}$	0.000403	30,727/297,191
rs6516886	T	Hospital-diagnosed hip and/or knee osteoarthritis	0.75	1.13	1.08	1.19	$5.36 \times 10^{-8}$	6,586/26,384	0.998499	Hip osteoarthritis	0.76	1.06	1.12	0.055135	5,714/199,421	0.999981	1.1	1.06	1.14	$5.84 \times 10^{-8}$	0.06276	12,300/225,805
rs6516886	T	Hospital-diagnosed hip and/or knee osteoarthritis	0.75	1.13	1.08	1.19	$5.36 \times 10^{-8}$	6,586/26,384	0.998499	Hip and/or knee osteoarthritis	0.76	1.05	1.11	0.043467	9,429/199,421	0.999981	1.09	1.06	1.13	$1.42 \times 10^{-7}$	0.01914	16,015/222,805

<sup>a</sup>Imputation accuracy was assessed with IMPUTE-infoscore. <sup>b</sup>Heterogeneity P values were derived from Cochran's Q test. <sup>c</sup>Two-sided P value, likelihood ratio test. EA, effect allele; n, sample size.



**Fig. 2 | Regional association plots for the nine novel osteoarthritis loci.** The y axis represents the negative logarithm (base 10) of the variant  $P$  value (likelihood ratio test), and the x axis represents the position on the chromosome (chr), with the names and location of genes and nearest genes shown at the bottom. The variant with the lowest  $P$  value in the region after combined discovery and replication is marked by a purple diamond. The same variant is marked by a purple dot showing the discovery  $P$  value. The colors of the other variants indicate their  $r^2$  with the lead variant.

studies have shown that *TGFA* regulates the conversion of cartilage to bone during the process of endochondral bone growth, and that it is a dysregulated cytokine present in degrading cartilage in osteoarthritis and a strong stimulator of cartilage degradation upregulated by articular chondrocytes in experimentally induced and human osteoarthritis<sup>18–21</sup>. The function of *TGFA* has also been associated with craniofacial development, palate closure and small body size<sup>22</sup>.

rs864839 resides in the intronic region of *JPH3* (junctophilin 3 gene) on chromosome 16 and was discovered in the any-joint-site osteoarthritis analysis. It was more strongly associated with hip osteoarthritis in the replication dataset (OR 1.08 (95% CI 1.05–1.11),  $P=2.1 \times 10^{-8}$ , EAF 0.71). Through meta-analysis of the any-site-osteoarthritis phenotype across the discovery and replication datasets, we determined  $P=7.02 \times 10^{-6}$  (Supplementary Table 11). *JPH3* is involved in the formation of junctional membrane structure, and it regulates neuronal calcium flux and has been reported to be expressed in pancreatic beta cells and in the regulation of insulin secretion.

rs116882138 was most strongly associated with hip and/or knee osteoarthritis in the discovery dataset and with knee osteoarthritis in the replication dataset (OR 1.34 (95% CI 1.21–1.49),  $P=5.09 \times 10^{-8}$ , EAF 0.02). It is an intergenic variant located 11 kb downstream of *MOB3B* (kinase activator 3B gene) and 16 kb upstream of *EQTN* (equatorin sperm-acrosome-associated gene) on chromosome 9. We found rs116882138 to be nominally associated with acetabular dysplasia, as determined by the center-edge angle ( $\beta -1.1388$ , s.e.m. 0.5276,  $P=0.031$ ) (Table 2 and Methods).

Finally, rs6516886 was prioritized on the basis of the hip and/or knee osteoarthritis-discovery analysis and was more strongly associated in the hip-osteoarthritis replication dataset (OR 1.10 (95% CI 1.06–1.14),  $P=5.84 \times 10^{-8}$ , EAF 0.75). rs6516886 is situated 1 kb upstream of *RWDD2B* (RWD-domain-containing 2B gene) on chromosome 21. *LTN1* (listerin E3 ubiquitin protein ligase 1 gene), which is located 28 kb from the variant, has been reported to affect musculoskeletal development in a mouse model<sup>23</sup>.

**Functional analysis.** Using molecular phenotyping through quantitative proteomics and RNA sequencing, we tested whether coding genes within 1 Mb of the novel osteoarthritis-associated variants were differentially expressed at 1% false discovery rate (FDR) in chondrocytes extracted from intact compared with degraded cartilage from patients with osteoarthritis undergoing total-joint-replacement surgery (Table 3 and Methods).

*PCYOX1*, located 209 kb downstream of rs3771501, showed significant evidence of differential expression (1.21-fold higher post-normalization in degraded cartilage at the RNA level,  $q=0.0047$ ; and 1.17-fold lower abundance at the protein level,  $q=0.0042$ ). This discrepancy may indicate potential clinical relevance, because the gene product is a candidate biomarker for osteoarthritis progression. Prenylcysteine oxidase 1, the protein product of this gene, is a secreted protein that catalyzes the degradation of prenylated proteins<sup>24</sup> and has been identified in urinary exosomes<sup>25</sup>. Further investigation into the chondrocyte and peripheral secretome is warranted to assess the potential of this molecule as a biomarker for

**Table 2 | Association of the nine osteoarthritis loci with radiographically derived osteoarthritis endophenotypes**

rsID	Minimal joint-space width <sup>b</sup>					Center-edge angle <sup>c</sup>					Alpha angle (cam deformity) <sup>d,e</sup>				
	EA	EAF	$\beta$	s.e.	P value <sup>f</sup>	EA	EAF	$\beta$	s.e.	P value <sup>f</sup>	EA	EAF	$\beta$	s.e.	P value <sup>f</sup>
rs2820436	A	0.317	-0.0146	0.0135	0.2817	A	0.317	-0.0104	0.1301	0.9363	A	0.318	0.0165	0.0675	0.8073
rs3771501	A	0.484	-0.0699	0.0127	$3.454 \times 10^{-8}$	A	0.474	0.1943	0.1199	0.1051	A	0.4779	-0.0144	0.0626	0.8176
rs11335718	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
rs11780978	A	0.389	-0.0291	0.0129	0.02419	A	0.386	0.078	0.1239	0.5291	A	0.3866	0.0035	0.0644	0.9564
rs2521349	A	0.398	0.0229	0.0128	0.07404	A	0.391	0.0998	0.1236	0.4192	A	0.3921	-0.0262	0.0644	0.6835
rs864839	N/A	N/A	N/A	N/A	N/A	T	0.702	-0.0206	0.1325	0.8766	T	0.7026	-0.0081	0.0691	0.907
rs375575359	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
rs116882138	N/A	N/A	N/A	N/A	N/A	A	0.0137	-1.1388	0.5276	0.0309	A	0.0135	0.1814	0.2607	0.4865
rs6516886 <sup>a</sup>	T	0.272	-0.0222	0.0143	0.1206	A	0.265	-0.1491	0.1373	0.2773	A	0.263	0.0544	0.0713	0.4458

<sup>a</sup>For minimal joint-space width, proxy variant rs2150403 ( $r^2 = 0.99$  with rs6516886) was used. <sup>b</sup>Sample size = 13,013. <sup>c</sup>Sample size = 6,880. <sup>d</sup>Sample size, cases = 639. <sup>e</sup>Sample size, controls = 4,339. <sup>f</sup>Two-sided P value following inverse variance-based meta-analysis. EA, effect allele; s.e., standard error; N/A, not available.

osteoarthritis progression. *PCYOX1* has been reported to be over-expressed in human dental-pulp-derived osteoblasts compared with osteosarcoma cells<sup>26</sup>. *FAM136A*, located 188 kb upstream of the same variant (rs3771501), showed 1.13-fold-lower transcriptional levels in chondrocytes from degraded articular cartilage ( $q = 0.0066$ ).

*BACH1* and *MAP3K7CL*, located in the vicinity of rs6516886, showed evidence of differential transcription (1.26-fold higher,  $q = 0.0019$ , and 1.37-fold higher,  $q = 0.0021$ , respectively, in degraded tissue). *BACH1* is a transcriptional repressor of heme oxygenase-1. Studies in *Bach1*-deficient mice have independently suggested inactivation of *Bach1* as a novel target for the prevention and treatment of meniscal degeneration<sup>27</sup> and of osteoarthritis<sup>28</sup>.

Finally, *PLAA* and *ZNF382*, located proximal to rs116882138 and rs375575359, respectively, showed higher transcription levels in degraded compared with intact cartilage (1.15-fold,  $q = 0.0027$ , and 1.31-fold,  $q = 0.0031$ , respectively). *BOP1*, located 451 kb downstream of rs11780978, showed 1.17-fold lower levels of transcription in degraded tissue ( $q = 0.003$ ).

We examined evidence for expression quantitative trait loci (eQTLs) in the Genotype-Tissue Expression GTEx tissues and found that none of the eQTLs identified at 5% FDR overlapped with the genes identified as differentially expressed between osteoarthritis intact and degraded cartilage (Supplementary Note and Supplementary Table 16).

**Fine mapping indicates noncoding variants at all loci.** For five of the new loci, the sum of probabilities of causality of all variants in the fine-mapped region was  $\geq 0.95$  ( $>0.99$  for two signals) and was  $>0.93$  for two further loci (Supplementary Table 17 and Methods). Most variants within each credible set had marginal posterior probabilities, whereas only a small number of variants had a posterior probability of association (PPA)  $> 0.1$ ; these accounted for 25–92% of PPA across the different regions. The credible set of four signals was narrowed down to three variants, one signal to two variants, and one signal to one variant, with a probability of causality  $> 0.1$ . For all nine regions, the variant identified as most likely to be causal was noncoding (Supplementary Table 18, Supplementary Note and Supplementary Fig. 5).

**Gene-based analyses.** Gene-set analysis identified *UQCCL1* and *GDF5*, located close to each other on chromosome 20, as key genes with consistent evidence of significant association with osteoarthritis across phenotype definitions (Supplementary Table 19 and Supplementary Note). *UQCCL1* and *GDF5* were significantly associated with four and three of the five osteoarthritis definitions, respectively. *GDF5* encodes growth differentiation factor 5,

a member of the TGF $\beta$  superfamily, and accruing evidence indicates that it plays a central role in skeletal health and development<sup>29–32</sup>. Pathway analyses identified significant associations between self-reported osteoarthritis and anatomical-structure morphogenesis ( $P = 4.76 \times 10^{-5}$ ) or ion-channel transport ( $P = 8.98 \times 10^{-5}$ ); hospital-diagnosed hip osteoarthritis and activation of MAPK activity ( $P = 1.61 \times 10^{-5}$ ); hospital-diagnosed knee osteoarthritis and histidine metabolism ( $P = 1.02 \times 10^{-5}$ ); and hospital-diagnosed hip and/or knee osteoarthritis and recruitment of mitotic centrosome proteins and complexes ( $P = 8.88 \times 10^{-5}$ ) (Supplementary Table 20 and Supplementary Fig. 6).

**Genetic links between osteoarthritis and other traits.** Established clinical risk factors for osteoarthritis include old age, female sex, obesity, occupational exposure to high levels of joint loading activity, previous injury, smoking status and family history of osteoarthritis. We estimated the genome-wide genetic correlation between osteoarthritis and 219 other traits and diseases and identified 35 phenotypes with significant (5% FDR) genetic correlation with osteoarthritis across definitions, with large overlap between the identified phenotypes (Supplementary Fig. 7, Fig. 3, Supplementary Table 21 and Methods).

The phenotypes with significant genetic correlations ( $r_g$ ) fell into the following broad categories: obesity, body mass index (BMI) and related anthropometric traits ( $r_g > 0$ ); type 2 diabetes ( $r_g > 0$ ); educational achievement ( $r_g < 0$ ); neuroticism, depressive symptoms ( $r_g > 0$ ) and sleep duration ( $r_g < 0$ ); mother's, father's or parents' age at death ( $r_g < 0$ ); reproductive phenotypes, including age at first birth ( $r_g < 0$ ) and number of children born ( $r_g > 0$ ); smoking, including age of smoking initiation ( $r_g < 0$ ) and having ever smoked ( $r_g > 0$ ), and lung cancer ( $r_g > 0$ ) (Fig. 3, Supplementary Table 21). The four phenotypes with significant genetic correlation in all analyses were number of years of schooling, waist circumference, hip circumference and BMI.

We found a nominally significant positive genetic correlation with rheumatoid arthritis, which did not pass multiple-testing correction for self-reported and hospital-diagnosed osteoarthritis ( $r_g = 0.14$ – $0.19$ , FDR 10–12%). Among musculoskeletal phenotypes, lumbar-spine bone mineral density showed a positive genetic correlation with hospital-diagnosed hip and/or knee osteoarthritis ( $r_g = 0.2$ , FDR = 3%) but did not reach significance in other analyses.

**Disentangling causality.** We undertook Mendelian randomization (MR) analyses<sup>33</sup> to strengthen causal inference regarding modifiable exposures that might influence osteoarthritis risk (Supplementary Tables 22–25 and Methods). Each  $\text{kg/m}^2$  increment in body mass

**Table 3 | Genes in the osteoarthritis-associated signals with significantly different gene expression and/or protein abundance in intact versus degraded articular cartilage**

Index variant	Gene	Position (chromosome: start-end)	Distance from index variant (kb)	Proteomics logFC	Proteomics FDR <i>q</i> value	RNA-seq logFC	RNA-seq FDR <i>q</i> value
rs3771501	<i>PCYOX1</i>	2: 70484518–70508323	209.3	–0.27	0.0042	0.27	0.0047
rs3771501	<i>FAM136A</i>	2: 70523107–70529222	188.4	N/A	N/A	–0.20	0.0066
rs6516886	<i>BACH1</i>	21: 30566392–31003071	172.7	N/A	N/A	0.32	0.0019
rs6516886	<i>MAP3K7CL</i>	21: 30449792–30548210	56.1	N/A	N/A	0.41	0.0021
rs11780978	<i>BOP1</i>	8: 145486055–145515082	451.2	N/A	N/A	–0.26	0.0030
rs116882138	<i>PLAA</i>	9: 26904081–26947461	366	–0.07	0.601	0.20	0.0027
rs375575359	<i>ZNF382</i>	19: 37095719–37119499	233.8	N/A	N/A	0.39	0.0031

logFC, log<sub>2</sub> fold change based on normalized values (increase indicates higher value in degraded cartilage); FDR, Benjamini–Hochberg FDR; N/A, proteomics data not available.

index was predicted to increase the risk of self-reported osteoarthritis by 1.11 (95% CI 1.07–1.15,  $P=8.3 \times 10^{-7}$ ). This result was consistent across MR methods (OR 1.52–1.66) and disease definition (OR 1.66–2.01). Consistent results were also observed for other obesity-related measures, such as waist circumference (OR 1.03 per cm increment; 95% CI 1.02–1.05,  $P=5 \times 10^{-4}$ ) and hip circumference (OR 1.03 per cm increment; 95% CI 1.01–1.06,  $P=0.021$ ). The OR values for type 2 diabetes liability and triglycerides were consistently small across estimators and osteoarthritis definitions; given that the analyses involving those traits were well powered (Supplementary Table 26), these results are compatible with either a weak or no causal effect. The results for years of schooling were not consistent across estimators, and there was evidence of directional horizontal pleiotropy, thus hampering any causal interpretation (Fig. 4). For lumbar-spine bone mineral density, there was evidence of a causal effect with OR per s.d. increment of 1.28 (95% CI 1.11–1.47,  $P=0.002$ ) for hip and/or knee osteoarthritis. This effect appeared to be site specific, with OR of 1.29 (95% CI 1.06–1.57,  $P=0.014$ ) for knee osteoarthritis, whereas the OR for hip osteoarthritis ranged from 0.71 to 1.57. There was also some evidence of a site-specific causal effect of height on knee osteoarthritis (OR 1.13 per s.d. increment; 95% CI 1.02–1.25,  $P=0.023$ ), which was consistent across estimators. One-sample MR analyses corroborated these findings, and obesity-related phenotypes presented strong statistical evidence after multiple-testing correction (Supplementary Table 27). These analyses did not detect reliable effects of smoking or reproductive traits on osteoarthritis (Supplementary Tables 28 and 29).

## Discussion

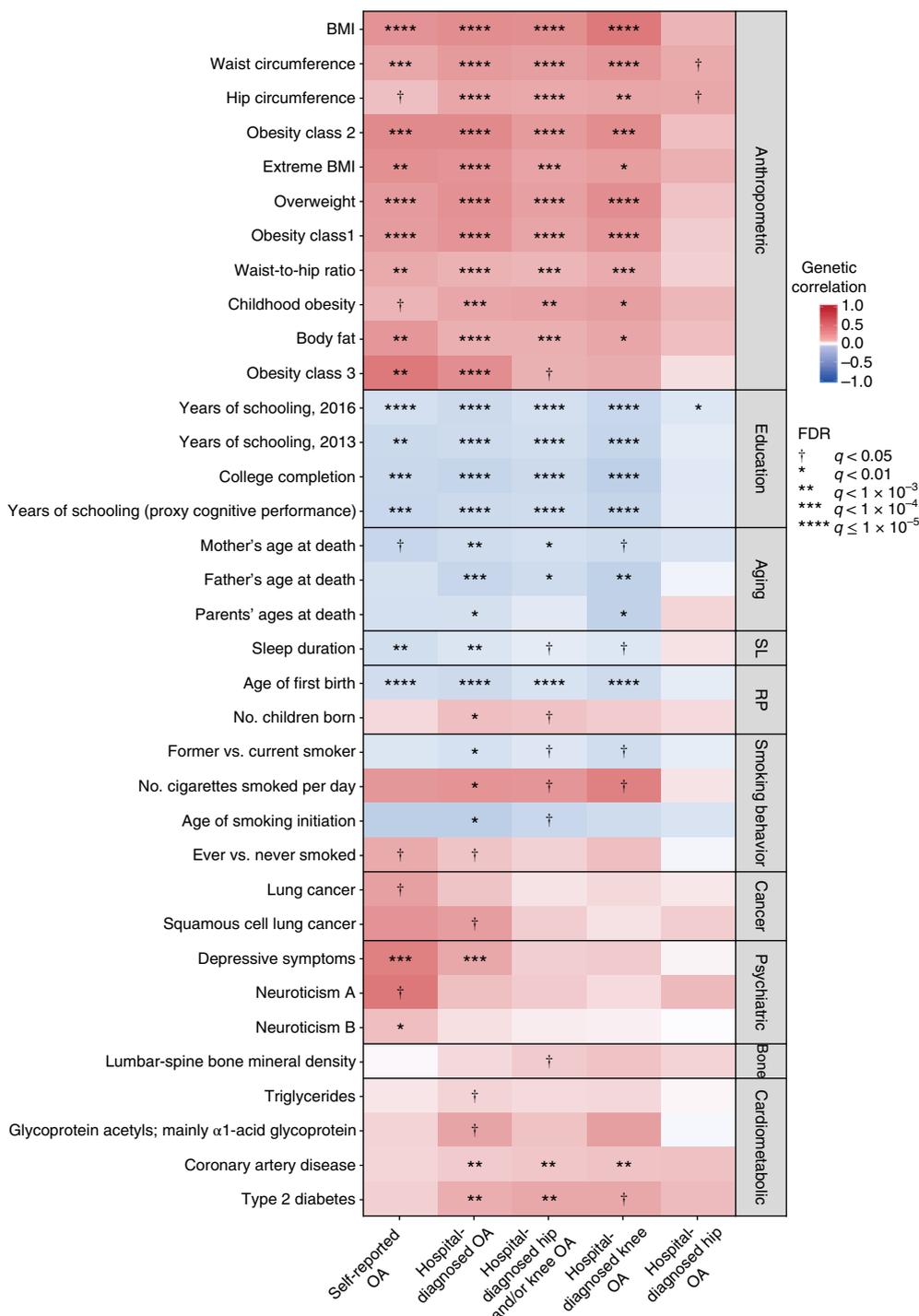
To improve understanding of the genetic etiology of osteoarthritis, we conducted a study combining genotype data in up to 327,918 individuals. We identified six novel, robustly replicating loci associated with osteoarthritis, three of which fell just under the corrected genome-wide-significance threshold. These loci provide a substantial increase in the number of known osteoarthritis loci. Together, all established osteoarthritis loci accounted for 26.3% of trait variance (Supplementary Fig. 8). The key attributes of this study were the large sample size and the homogeneity of the UK Biobank dataset, coupled with independent replication, independent association with clinically relevant radiographic endophenotypes and functional genomics follow-up in primary osteoarthritis tissue. We further capitalized on the wealth of available genome-wide summary statistics across complex traits to identify genetic correlations between osteoarthritis and multiple molecular, physiological and behavioral phenotypes, and we performed formal MR analyses to assess causality and disentangle complex cross-trait epidemiological relationships.

Most novel signals were at common frequency variants and conferred small to modest effects, in line with a highly polygenic model underpinning osteoarthritis risk. We identified one low-frequency variant associated with osteoarthritis (MAF 0.02) with a modest effect size (combined OR 1.34). Even though our study was well powered to detect such variants, we found no evidence of a role of low-frequency variation of large effect in osteoarthritis susceptibility (Supplementary Table 5). The power of this study was very limited for low-frequency variants with OR < 1.50 and for rare variants. We estimate a requirement of up to 40,000 osteoarthritis cases and 160,000 controls to recapitulate the effects identified in this study at genome-wide significance, on the basis of the sample-size-weighted effect-allele frequencies and replication-cohort odds-ratio estimates (Table 1 and Supplementary Table 30).

We integrated functional information with statistical evidence for association to fine map the locations of likely causal variants and genes. All the predicted most likely causal variants resided in noncoding sequence: six were intronic, and three were intergenic. We were able to refine the association signal to a single variant in one instance, and to variants residing within a single gene in three instances, although the mechanisms of action may be mediated through other genes in the vicinity.

We empirically found self-reported osteoarthritis definition to be a powerful tool for genetic association studies, as evidenced, for example, by the genome-wide significance reached for the established *GDF5* osteoarthritis locus in only the self-reported disease-status analyses. Published epidemiological studies investigating osteoarthritis via self-reporting<sup>34,35</sup> and validation of self-reported status against primary-care records has yielded similar conclusions<sup>34</sup>. We also found very high genetic correlation between self-reported and hospital-diagnosed osteoarthritis, as well as similar variant-based heritability estimates, thus corroborating the validity of self-reported osteoarthritis status for genetic studies. However, we also note that the hospital-diagnosed-osteoarthritis analyses had higher heritability and yielded stronger evidence of effect-direction concordance at established loci, thus indicating that larger sample sizes would afford the power required to convincingly detect the established loci. Hospital-diagnosed-osteoarthritis data may potentially capture a different patient demographic than self-reported data (Supplementary Note). From the results of this study, we deduce that there is no gold standard for osteoarthritis definition in genetics studies, and we identified advantages in using both methods of defining disease to broadly maximize discovery power.

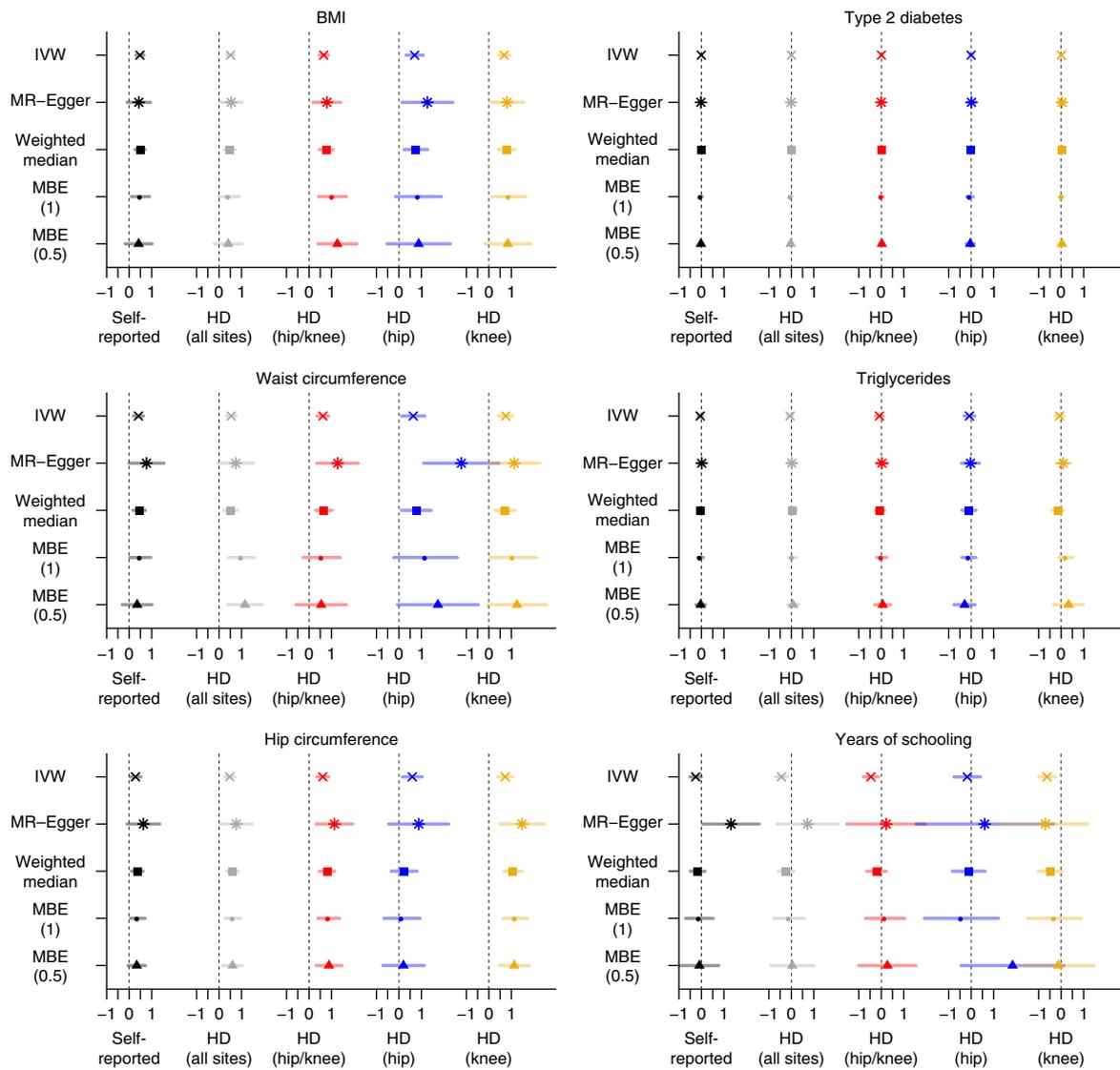
We identified strong genome-wide correlation between hip and knee osteoarthritis, thus indicating a substantial shared genetic etiology that has been hitherto overlooked. We therefore sought to replicate signals across these highly correlated phenotypes and



**Fig. 3 | Heat map of genetic correlations between osteoarthritis phenotypes in UK Biobank and 35 traits grouped in ten categories from GWAS consortia.** Symbols and hues depict the two-tailed Benjamini–Hochberg FDR  $q$  values and strength of the genetic correlation (darker shade denotes stronger correlation), respectively. Red and blue indicate positive and negative correlations, respectively. RP, reproductive; SL, sleep; OA, osteoarthritis.

to identify multiple instances of signals detected in the larger discovery analysis of osteoarthritis and independently replicated in joint-specific definitions of disease. Indeed, when examining the replication phenotypes, we found no instances of confirmed replication in which the replication phenotype was not captured within the accompanying discovery-phenotype definition. Further analysis in larger sample sets with precise phenotyping should help distinguish signal specificity.

Two of the newly identified signals, indexed by rs11780978 and rs2820436, resided in regions with established metabolic- and anthropometric-trait associations. Osteoarthritis is epidemiologically associated with high BMI, and the association is stronger for knee osteoarthritis. In line with this finding, we observed higher genetic correlation between BMI and knee osteoarthritis ( $r_g = 0.52$ ,  $P = 2.2 \times 10^{-11}$ ) compared with hip osteoarthritis ( $r_g = 0.28$ ,  $P = 4 \times 10^{-4}$ ). BMI is also known to be genetically correlated with



**Fig. 4 | Two-sample MR estimates and 95% CI values of the effects of obesity-related measures, triglyceride levels, years of schooling and type 2 diabetes liability on different definitions of osteoarthritis.** All values are shown in s.d. units except for type 2 diabetes liability, which is shown in  $\ln(\text{OR})$  units. HD, hospital diagnosed; IVW, inverse-variance weighting; MBE, mode-based estimate; MBE (1), tuning parameter  $\varphi = 1$ ; MBE (0.5), tuning parameter  $\varphi = 0.5$ .

education phenotypes, depressive symptoms, and reproductive and other phenotypes; hence, some of the genetic correlations for osteoarthritis observed here may be mediated through BMI. However, for the education and personality/psychiatric phenotypes, the strength of the genetic correlations observed here for osteoarthritis was substantially higher than that observed for BMI (for example, hospital-diagnosed osteoarthritis and years of schooling had  $rg = -0.45$ ,  $P = 5 \times 10^{-27}$ , whereas BMI and years of schooling had  $rg = -0.27$ ,  $P = 9 \times 10^{-32}$ ; hospital-diagnosed osteoarthritis and depressive symptoms had  $rg = 0.49$ ,  $P = 6 \times 10^{-7}$ , whereas BMI and depressive symptoms had  $rg = 0.10$ ,  $P = 0.023$ ). Epidemiologically, lower educational levels are known to be particularly associated with risk of knee osteoarthritis, even with adjustment for BMI<sup>36</sup>.

MR provided further insight into the nature of the genetic correlations that we observed. In the case of BMI and other obesity-related measures, there was evidence of a causal effect of those phenotypes on osteoarthritis. This result corroborated findings from conventional observational studies<sup>37</sup>, which are prone to important limitations (such as reverse causation and residual

confounding) regarding causal inference<sup>38</sup>. For all other exposure phenotypes, there was no convincing evidence of a causal effect on osteoarthritis risk, thus suggesting that the genetic correlations detected by linkage disequilibrium (LD)-score regression may be mostly due to horizontal pleiotropy, although for some phenotypes the MR analyses were underpowered (Supplementary Table 26). In the case of triglycerides and liability to type 2 diabetes, the MR analyses had sufficient power to rule out nonsmall causal effects, thus suggesting that these phenotypes have at most weak effects on osteoarthritis risk.

Importantly, structural changes in joints usually precede the onset of osteoarthritis symptoms. Articular cartilage is an avascular, aneural tissue. It provides tensile strength, compressive resilience and a low-friction articulating surface. Chondrocytes are the only cell type in cartilage. The mode of function of noncoding DNA is linked to context-dependent regulation of gene expression, and identification of the causal variants and the genes that they affect requires experimental analysis of genome regulation in the proper cell type. Our functional analysis of genes in osteoarthritis-associated

regions and pathways identified differentially expressed molecules in chondrocytes extracted from degraded compared with intact articular cartilage. Cartilage degeneration is a key hallmark of osteoarthritis pathogenesis, and regulation of these genes may be implicated in disease development and progression.

Osteoarthritis is a leading cause of disability worldwide, and it imposes a substantial public-health and health economic burden. Here, we gleaned novel insights into the genetic etiology of osteoarthritis and implicated genes with translational potential<sup>10,13,14,27,28</sup>. The cohorts contributing to this study were composed of European-descent populations. In the future, large-scale whole-genome-sequencing studies of well-phenotyped individuals across diverse populations should capture the full allele frequency and variation type spectrum, and afford further insights into the causes of this debilitating disease.

**URLs.** Quanto, <http://biostats.usc.edu/Quanto.html>; genotyping and quality control of UK Biobank, [http://biobank.ctsu.ox.ac.uk/crystal/docs/genotyping\\_qc.pdf](http://biobank.ctsu.ox.ac.uk/crystal/docs/genotyping_qc.pdf); genotype imputation of UK Biobank, [http://www.ukbiobank.ac.uk/wp-content/uploads/2014/04/imputation\\_documentation\\_May2015.pdf](http://www.ukbiobank.ac.uk/wp-content/uploads/2014/04/imputation_documentation_May2015.pdf); GRCh38 cDNA assembly release 87, [http://ftp.ensembl.org/pub/release-87/fasta/homo\\_sapiens/cdna/](http://ftp.ensembl.org/pub/release-87/fasta/homo_sapiens/cdna/).

## Methods

Methods, including statements of data availability and any associated accession codes and references, are available at <https://doi.org/10.1038/s41588-018-0079-y>.

Received: 9 June 2017; Accepted: 29 January 2018;

Published online: 20 March 2018

## References

- Loeser, R. F., Goldring, S. R., Scanzello, C. R. & Goldring, M. B. Osteoarthritis: a disease of the joint as an organ. *Arthritis Rheum.* **64**, 1697–1707 (2012).
- Felson, D. T. et al. Osteoarthritis: new insights. Part 1: the disease and its risk factors. *Ann. Intern. Med.* **133**, 635–646 (2000).
- Cibrián Uhalte, E., Wilkinson, J. M., Southam, L. & Zeggini, E. Pathways to understanding the genomic aetiology of osteoarthritis. *Hum. Mol. Genet.* **26**, R193–R201 (2017).
- Ma, C., Blackwell, T., Boehnke, M. & Scott, L. J., GoT2D investigators. Recommended joint and meta-analysis strategies for case-control association testing of single low-count variants. *Genet. Epidemiol.* **37**, 539–550 (2013).
- Broome, D. T. & Datta, N. S. Mitogen-activated protein kinase phosphatase-1: function and regulation in bone and related tissues. *Connect. Tissue Res.* **57**, 175–189 (2016).
- Rodríguez-Carballo, E., Gámez, B. & Ventura, F. p38 MAPK signaling in osteoblast differentiation. *Front. Cell Dev. Biol.* **4**, 40 (2016).
- Wei, L., Sun, X. J., Wang, Z. & Chen, Q. CD95-induced osteoarthritic chondrocyte apoptosis and necrosis: dependency on p38 mitogen-activated protein kinase. *Arthritis Res. Ther.* **8**, R37 (2006).
- Wang, Q. et al. Bioinformatics analysis of gene expression profiles of osteoarthritis. *Acta Histochem.* **117**, 40–46 (2015).
- Prasad, I. et al. Osteoarthritic cartilage chondrocytes alter subchondral bone osteoblast differentiation via MAPK signalling pathway involving ERK1/2. *Bone* **46**, 226–235 (2010).
- Prasad, I. et al. Inhibition of p38 pathway leads to OA-like changes in a rat animal model. *Rheumatology (Oxford)* **51**, 813–823 (2012).
- Prasad, I. et al. ERK-1/2 and p38 in the regulation of hypertrophic changes of normal articular cartilage chondrocytes induced by osteoarthritic subchondral osteoblasts. *Arthritis Rheum.* **62**, 1349–1360 (2010).
- Zhang, Y., Pizzute, T. & Pei, M. A review of crosstalk between MAPK and Wnt signals and its impact on cartilage regeneration. *Cell Tissue Res.* **358**, 633–649 (2014).
- Namdari, S., Wei, L., Moore, D. & Chen, Q. Reduced limb length and worsened osteoarthritis in adult mice after genetic inhibition of p38 MAPK kinase activity in cartilage. *Arthritis Rheum.* **58**, 3520–3529 (2008).
- Zhang, R., Murakami, S., Coustry, F., Wang, Y. & de Crombrughe, B. Constitutive activation of MKK6 in chondrocytes of transgenic mice inhibits proliferation and delays endochondral bone formation. *Proc. Natl Acad. Sci. USA* **103**, 365–370 (2006).
- Castañón, M. J., Walko, G., Winter, L. & Wiche, G. Plectin-intermediate filament partnership in skin, skeletal muscle, and peripheral nerve. *Histochem. Cell Biol.* **140**, 33–53 (2013).
- Tseng, K. Y., Chen, Y. H. & Lin, S. Zinc finger protein ZFP36L1 promotes osteoblastic differentiation but represses adipogenic differentiation of mouse multipotent cells. *Oncotarget* **8**, 20588–20601 (2017).
- Castaño-Betancourt, M. C. et al. Novel genetic variants for cartilage thickness and hip osteoarthritis. *PLoS Genet.* **12**, e1006260 (2016).
- Usmani, S. E. et al. Transforming growth factor alpha controls the transition from hypertrophic cartilage to bone during endochondral bone growth. *Bone* **51**, 131–141 (2012).
- Appleton, C. T., Usmani, S. E., Bernier, S. M., Aigner, T. & Beier, F. Transforming growth factor alpha suppression of articular chondrocyte phenotype and Sox9 expression in a rat model of osteoarthritis. *Arthritis Rheum.* **56**, 3693–3705 (2007).
- Appleton, C. T., Usmani, S. E., Mort, J. S. & Beier, F. Rho/ROCK and MEK/ERK activation by transforming growth factor-alpha induces articular cartilage degradation. *Lab. Invest.* **90**, 20–30 (2010).
- Usmani, S. E. et al. Context-specific protection of TGF $\alpha$  null mice from osteoarthritis. *Sci. Rep.* **6**, 30434 (2016).
- Miettinen, P. J. et al. Epidermal growth factor receptor function is necessary for normal craniofacial development and palate closure. *Nat. Genet.* **22**, 69–73 (1999).
- Chu, J. et al. A mouse forward genetics screen identifies LISTERIN as an E3 ubiquitin ligase involved in neurodegeneration. *Proc. Natl Acad. Sci. USA* **106**, 2097–2103 (2009).
- Wang, M. & Casey, P. J. Protein prenylation: unique fats make their mark on biology. *Nat. Rev. Mol. Cell Biol.* **17**, 110–122 (2016).
- Gonzales, P. A. et al. Large-scale proteomics and phosphoproteomics of urinary exosomes. *J. Am. Soc. Nephrol.* **20**, 363–379 (2009).
- Palmieri, A. et al. Comparison between osteoblasts derived from human dental pulp stem cells and osteosarcoma cell lines. *Cell Biol. Int.* **32**, 733–738 (2008).
- Ochiai, S. et al. Oxidative stress reaction in the meniscus of Bach 1 deficient mice: potential prevention of meniscal degeneration. *J. Orthop. Res.* **26**, 894–898 (2008).
- Takada, T. et al. Bach1 deficiency reduces severity of osteoarthritis through upregulation of heme oxygenase-1. *Arthritis Res. Ther.* **17**, 285 (2015).
- Capellini, T. D. et al. Ancient selection for derived alleles at a GDF5 enhancer influencing human growth and osteoarthritis risk. *Nat. Genet.* **49**, 1202–1210 (2017).
- Daams, M., Luyten, F. P. & Lories, R. J. GDF5 deficiency in mice is associated with instability-driven joint damage, gait and subchondral bone changes. *Ann. Rheum. Dis.* **70**, 208–213 (2011).
- Miyamoto, Y. et al. A functional polymorphism in the 5' UTR of GDF5 is associated with susceptibility to osteoarthritis. *Nat. Genet.* **39**, 529–533 (2007).
- Southam, L. et al. An SNP in the 5'-UTR of GDF5 is associated with osteoarthritis susceptibility in Europeans and with in vivo differences in allelic expression in articular cartilage. *Hum. Mol. Genet.* **16**, 2226–2232 (2007).
- Davey Smith, G. & Ebrahim, S. 'Mendelian randomization': can genetic epidemiology contribute to understanding environmental determinants of disease? *Int. J. Epidemiol.* **32**, 1–22 (2003).
- Prieto-Alhambra, D. et al. An increased rate of falling leads to a rise in fracture risk in postmenopausal women with self-reported osteoarthritis: a prospective multinational cohort study (GLOW). *Ann. Rheum. Dis.* **72**, 911–917 (2013).
- Baldwin, J. N. et al. Self-reported knee pain and disability among healthy individuals: reference data and factors associated with the knee injury and osteoarthritis outcome score (KOOS) and KOOS-Child. *Osteoarthritis Cartilage* **25**, 1282–1290 (2017).
- Callahan, L. F. et al. Associations of educational attainment, occupation and community poverty with knee osteoarthritis in the Johnston County (North Carolina) osteoarthritis project. *Arthritis Res. Ther.* **13**, R169 (2011).
- Hussain, S. M. et al. How are obesity and body composition related to patellar cartilage? a systematic review. *J. Rheumatol.* **44**, 1071–1082 (2017).
- Fewell, Z., Davey Smith, G. & Sterne, J. A. The impact of residual and unmeasured confounding in epidemiologic studies: a simulation study. *Am. J. Epidemiol.* **166**, 646–655 (2007).

## Acknowledgements

This research was conducted by using the UK Biobank Resource under application no. 9979. This work was funded by the Wellcome Trust (WT098051). We are grateful to R. Brooks, J. Choudhary and T. Roumeliotis for their contribution to the functional genomics data collection. The Human Research Tissue Bank is supported by the NIHR Cambridge Biomedical Research Centre. The Rotterdam Study is funded by Erasmus Medical Center and Erasmus University, Organization for Health Research and Development (ZonMw); the Netherlands Genomics Initiative (NGI)/Netherlands Organisation of Scientific Research (NWO); the Netherlands Consortium for Healthy

Aging (NCHA), Research Institute for Diseases in the Elderly (RIDE); the Ministry of Education, Culture and Science; the Ministry for Health, Welfare and Sports; the European Commission (DG XII); and the Municipality of Rotterdam.

### Author contributions

Association analyses: E. Zengini, K.H., I.T., L.S., J.S., S.H. and A.G. Mendelian randomization: F.P.H. and G.D.S. Functional genomics sample collection: A.McC., J.M.W. and E. Zeggini. Functional genomics analyses: J.S. and L.S. Endophenotype analyses: C.G.B., A.G.U. and J.B.J.v.M. Replication analyses: U.S., T.I., H.J., U.T. and K.S. Bioinformatics: A.G., D.S. and B.K. Student supervision: K.H., G.C.B., G.D.S., J.M.W. and E. Zeggini. Manuscript writing: E. Zengini, K.H., I.T., J.S., F.P.H., L.S., C.G.B., U.S., D.S., J.B.J.v.M., G.D.S., J.M.W. and E. Zeggini.

### Competing interests

U.T., U.S. and K.S. are employees of deCODE genetics/Amgen.

### Additional information

**Supplementary information** is available for this paper at <https://doi.org/10.1038/s41588-018-0079-y>.

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

**Correspondence and requests for materials** should be addressed to E.Z.

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

## Methods

**Accuracy of self-reported data.** We evaluated the classification accuracy of self-reported disease status by estimating the sensitivity, specificity, PPV and negative predictive values (NPV) in the self-reported and hospital-diagnosed disease-definition datasets. We performed a sensitivity test to evaluate the true-positive rate by calculating the proportion of individuals diagnosed with osteoarthritis that were correctly identified as such in the self-reported analysis, then performed a specificity test to evaluate the true-negative rate by calculating the proportion of individuals not diagnosed with osteoarthritis that were correctly identified as such in the control set. The number of individuals overlapping between the self-reported ( $n_{SR} = 12,658$ ) and hospital-diagnosed ( $n_{HD} = 10,083$ ) datasets was  $n_{OVER} = 3,748$ . The total number of individuals was  $n_{TOT} = 138,997$ . Sensitivity =  $\frac{n_{OVER}}{n_{SR}}$ ; specificity =  $\frac{n_{TOT} - (n_{HD} + n_{SR} - n_{OVER})}{n_{TOT} - n_{HD}}$ ; PPV =  $\frac{n_{OVER}}{n_{SR}}$ ; NPV =  $\frac{n_{TOT} - (n_{HD} + n_{SR} - n_{OVER})}{n_{TOT} - n_{SR}}$ .

**Discovery GWAS.** UK Biobank's scientific protocol and operational procedures were reviewed and approved by the North West Research Ethics Committee (REC reference no. 06/MRE08/65). The first UK Biobank release of genotype data included ~150,000 volunteers between 40 and 69 years old from the UK, genotyped at approximately 820,967 SNPs. 50,000 samples were genotyped with the UKBiLEVE array, and the remaining samples were genotyped with the UK Biobank Axiom array (Affymetrix; URLs). The UK Biobank Axiom is an update of UKBiLEVE, and the two arrays share 95% content. In total, after sample and SNP quality control (QC), which was carried out centrally, 152,763 individuals and 806,466 directly typed SNPs remained. Phasing, imputation and derivation of principal components were also carried out centrally. Briefly, the combined UK10K/1000 Genomes Project haplotype reference panel was used to impute untyped variants through the IMPUTE3 program (URLs). After imputation, the number of variants reached 73,355,667 in 152,249 individuals. We performed additional QC checks and excluded samples with call rate  $\leq 97\%$ . We checked samples for sex discrepancies, excess heterozygosity, relatedness and ancestry, and removed possibly contaminated and withdrawn samples. After QC, the number of individuals was 138,997. We excluded 528 SNPs that had been centrally flagged as being subject to exclusion due to failure in one or more additional quality metrics.

To define osteoarthritis cases, we used the self-reported status questionnaire and the Hospital Episode Statistics data (Supplementary Table 3 and Supplementary Note). We conducted five osteoarthritis-discovery GWAS and one sensitivity analysis. The case strata were as follows: self-reported osteoarthritis at any site,  $n = 12,658$ ; sensitivity analysis (a random subset of the self-reported cohort equal to the sample size of the hospital-diagnosed cohort),  $n = 10,083$ ; hospital-diagnosed osteoarthritis at any site, on the basis of ICD10 and/or ICD9 hospital-record codes,  $n = 10,083$ ; hospital-diagnosed hip osteoarthritis,  $n = 2,396$ ; hospital-diagnosed knee osteoarthritis,  $n = 4,462$ ; and hospital-diagnosed hip and/or knee osteoarthritis,  $n = 6,586$ . We applied exclusion criteria to minimize misclassification in the control datasets to the extent possible (using approximately four times the number of cases for each definition) (Supplementary Table 2 and Supplementary Fig. 1). We restricted the number of controls used and did not use the full set of available genotyped control samples from UK Biobank to guard against association test statistics behaving anticonservatively in the presence of stark case-control imbalance for alleles with minor allele count (MAC)  $< 400$  (ref. <sup>3</sup>) (analogous to MAF of  $\sim 0.02$  in the self-reported and hospital-diagnosed osteoarthritis datasets). For the control set, we excluded all participants diagnosed with any musculoskeletal disorder or having relevant symptoms or signs, such as pain and arthritis, and we selected older participants to ensure that we minimized the number of controls that might be diagnosed with osteoarthritis in the future, while keeping the number of males and females balanced (Supplementary Table 1).

At the SNP level, we further filtered for Hardy-Weinberg-equilibrium  $P \leq 10^{-6}$ , MAF  $\leq 0.001$  and info score  $< 0.4$  (Supplementary Fig. 1). We tested for association by using the frequentist likelihood ratio test (LRT) and method ml in SNPTEST v2.5.2 (ref. <sup>39</sup>) with adjustment for the first ten principal components to control for population structure. Power calculations were carried out in Quanto v1.2.4 (URLs).

**Replication.** Two hundred independent and novel variants with  $P < 1.0 \times 10^{-5}$  in the discovery analyses were used for in silico replication in an independent cohort from Iceland (deCODE) through fixed-effects inverse-variance-weighted meta-analysis in METAL<sup>40</sup>. One hundred seventy three variants were present in the replication cohort. The remaining 27 variants had ambiguous alleles (i.e., those incompatible because of alignment issues) and were not included in further analyses. The significance threshold for association in the replication study was hence  $0.05/173 = 2.9 \times 10^{-4}$ . The deCODE dataset comprised four osteoarthritis phenotypes: any osteoarthritis site (18,069 cases and 246,293 controls), hip osteoarthritis (5,714 cases and 199,421 controls), knee osteoarthritis (4,672 cases and 172,791 controls) and hip and/or knee osteoarthritis (9,429 cases and 199,421 controls). We performed meta-analyses (across osteoarthritis definitions), using summary statistics from the UK Biobank osteoarthritis analyses and deCODE. We used  $P \leq 2.8 \times 10^{-8}$  as the threshold corrected for the effective number of traits to report genome-wide significance.

**Replication cohort.** The information on hip, knee and vertebral osteoarthritis was obtained from Landspítali University Hospital electronic health records, from Akureyri Hospital electronic health records and from a national Icelandic hip or knee arthroplasty registry<sup>41</sup>. Individuals with secondary osteoarthritis (for example, Perthes disease and hip dysplasia), post-trauma osteoarthritis (for example, ACL rupture) or a diagnosis of rheumatoid arthritis were excluded from these lists. Only those diagnosed with osteoarthritis after the age of 40 were included. Subjects with hand osteoarthritis were drawn from a database of 9,000 patients with hand osteoarthritis that was initiated in 1972 (ref. <sup>42</sup>). The study was approved by the Data Protection Authority of Iceland and the National Bioethics Committee of Iceland. Informed consent was obtained from all participants.

**Association with osteoarthritis-related endophenotypes.** The nine replicating genetic loci were examined for association in radiographic osteoarthritis endophenotypes. This examination was done for three phenotypes: minimal joint-space width (mJSW) and two measures of hip-shape deformities known to be strong predictors of osteoarthritis: acetabular dysplasia (measured with center-edge (CE) angle) and cam deformity (measured with alpha angle). For mJSW, association statistics for the variants were looked up in a previously published GWAS, which performed joint analysis of data from Rotterdam Study I (RS-I), Rotterdam Study II (RS-II), TwinsUK, SOF and MrOS by using standardized age-, sex- and population stratification (four principal components)-adjusted residuals from linear regression<sup>17</sup>. For the two hip-shape phenotypes, CE angle and alpha angle were measured as previously described. CE angle was analyzed as a continuous phenotype. We conducted GWAS on a total of 6,880 individuals from RS-I, RS-II, Rotterdam Study III (RS-III) and CHECK<sup>43</sup> datasets, using standardized age- and sex-adjusted residuals from linear regression. For cam deformity, individuals with an alpha angle  $> 60^\circ$  were defined as cases ( $n = 639$ ), and all others were defined as controls (4,339). The GWAS was done on individuals from RS-I, RS-II and CHECK, by using age, sex and principal components to adjust for population stratification as covariates. The results for the separate cohorts were combined in a meta-analysis using inverse-variance weighting with METAL<sup>40</sup>. Genomic-control correction was applied to the standard errors and  $P$  values before meta-analysis.

**Functional genomics. Patients and samples.** We collected cartilage samples from 38 patients undergoing total-joint-replacement surgery: 12 patients with knee osteoarthritis (cohort 1; 2 women, 10 men, age 50–88 years); 17 patients with knee osteoarthritis (cohort 2; 12 women, 5 men, age 54–82 years); and 9 patients with hip osteoarthritis (cohort 3; 6 women, 3 men, age 44–84 years). We collected matched intact and degraded cartilage samples from each patient (Supplementary Note). Cartilage was separated from bone, and chondrocytes were extracted from each sample. From each isolated chondrocyte sample, we extracted RNA and protein. All patients provided full written informed consent before participation. All sample collection and RNA- and protein-extraction steps were as described in detail in ref. <sup>44</sup>. This work was approved by Oxford NHS REC C (10/H0606/20), and samples were collected under Human Tissue Authority license 12182, Sheffield Musculoskeletal Biobank, University of Sheffield, UK. Samples were also collected under National Research Ethics approval reference 11/EE/0011, Cambridge Biomedical Research Centre Human Research Tissue Bank, Cambridge University Hospitals, UK (additional information in Supplementary Note).

**Proteomics.** Proteomics analysis was performed on intact and degraded cartilage samples from 24 individuals (15 from cohort 2; 9 from cohort 3). LC-MS analysis was performed on a Dionex Ultimate 3000 UHPLC system coupled with an Orbitrap Fusion Tribrid mass spectrometer. To account for protein loading, abundance values were normalized by the sum of all protein abundances in a given sample, then  $\log_2$  transformed and quantile normalized. We restricted the analysis to 3,917 proteins that were quantified in all samples. We tested proteins for differential abundance by using limma<sup>45</sup> in R, on the basis of a within-individual paired sample design. Significance was defined at 1% Benjamini-Hochberg FDR to correct for multiple testing. Of the 3,732 proteins with unique mapping of gene name and Ensembl ID, we used 245 proteins with significantly different abundance between intact and degraded cartilage at 1% FDR.

**RNA sequencing.** We performed a gene expression analysis on samples from all 38 patients. Multiplexed libraries were sequenced on the Illumina HiSeq 2000 platform (75-bp paired-end read length), thus yielding bam files for cohort 1 and cram files for cohorts 2 and 3. The cram files were converted to bam files in samtools 1.3.1 (ref. <sup>46</sup>) and then to fastq files in biobambam 0.0.191 (ref. <sup>47</sup>), after exclusion of reads that failed QC. We obtained transcript-level quantification by using salmon 0.8.2 (ref. <sup>48</sup>) (with `--gcBias` and `--seqBias` flags to account for potential biases) and the GRCh38 cDNA assembly release 87 downloaded from Ensembl (URLs). We converted transcript-level to gene-level count estimates, with estimates for 39,037 genes, on the basis of Ensembl gene IDs. After quality control, we retained expression estimates for 15,994 genes with counts per million of 1 or higher in at least ten samples. Limma-voom<sup>49</sup> was used to remove heteroscedasticity from the estimated expression data. We tested genes for differential expression in limma<sup>45</sup> in R (with `lmFit` and `eBayes`), on the basis of a within-individual paired sample design. Significance was defined at 1% Benjamini-Hochberg FDR to correct for multiple testing. Of the 14,408 genes with unique

mapping of gene name and Ensembl ID, we used 1,705 genes with significantly different abundance between intact and degraded cartilage at 1% FDR.

**Fine mapping.** We constructed regions for fine mapping, by taking a window of at least 0.1 cM to either side of each index variant. The region was extended to the furthest variant with  $r^2 > 0.1$  with the index variant within a 1-Mb window. LD calculations for extending the region were based on whole-genome-sequenced EUR samples from the combined reference panel from UK10K<sup>50</sup> and the 1000 Genomes Project<sup>51</sup>. For each region we implemented the Bayesian fine-mapping method CAVIARBF<sup>52</sup>, which uses association summary statistics and correlations among variants to calculate Bayes' factors and the posterior probability of each variant being causal. We assumed a single causal variant in each region and calculated 95% credible intervals, which contained the minimum set of variants jointly having at least 95% probability of including the causal variant. We also applied the extended CAVIARBF method, which uses functional annotation scores to upweigh variants according to their predicted functional scores. To this end, we downloaded precalculated CADD<sup>53</sup> and Eigen<sup>54</sup> scores from their equivalent websites. We observed better separation of severe-consequence genic variants with the CADD score and better separation of regulatory variants with the Eigen score, and we therefore created a combined score in which splice-acceptor, splice-donor, stop-loss, stop-gain, missense and splice-region variants were assigned their CADD-Phred scores, and the rest were assigned their Eigen-Phred scores.

**Functional enrichment analysis.** We used genome-wide summary statistics to test for enrichment of functional annotations. We used GARFIELD<sup>55</sup> with customized functional annotations, making use of the functional genomics data that we generated in primary articular chondrocytes by using RNA sequencing and quantitative proteomics. We defined differentially expressed genes separately at the RNA (transcriptional) level and at the protein level when comparing intact and degraded cartilage (1% FDR). We extended each differentially regulated gene by 5 kb on each side. Using GARFIELD's approach, we calculated the effective number of independent annotations to be 1.995, which led to an adjusted-*P*-value significance level of 0.025. We tested for enrichment by using variants with  $P < 1.0 \times 10^{-5}$ , and no analysis surpassed the corrected significance threshold.

**LD regression.** We used LDHub<sup>56</sup> (accessed 23–27 January 2017) to estimate the genome-wide genetic correlation between each of the osteoarthritis definitions and 219 other human traits and diseases. In each analysis, we extracted variants with rsIDs (range 11999363–15561966) and uploaded the corresponding association summary statistics to LDHub; the analysis yielded 896,076–1,172,130 variants overlapping with LDHub. We corrected for multiple testing by defining significance at 5% Benjamini–Hochberg FDR for each of the five osteoarthritis analyses.

**Mendelian randomization analysis.** We used MR to assess the potential causal role of the phenotypes identified in the LD-score regression analysis on osteoarthritis. We also included birth weight and height (Supplementary Table 22). In all analyses, the primary outcome variable was self-reported osteoarthritis. We used data from hospital records (which were available for a much smaller number of individuals) as sensitivity analyses and to identify potential site-specific effects.

**Data sources.** Genetic instruments were identified from publicly available summary GWAS results through the TwoSampleMR R package, which allows for extraction of the data available in the MR-Base database<sup>57</sup>. Only results that combined both sexes were extracted. Preference was given to studies restricted to European populations to minimize the risk of bias due to population stratification; however, for several traits, those results were either not available or corresponded to much smaller studies (Supplementary Table 22). However, this aspect was unlikely to substantially bias the results, because all studies used correction methods, and even multiethnic studies were composed of mostly European populations. The exception was for number of children born and age of the individual at birth of the first child: because the GWAS of reproductive traits by Barban and colleagues<sup>58</sup> was not available in MR-Base, we extracted summary association results for the variants that achieved genome-wide significance directly from the paper and used coefficients from each sex in sex-specific analyses. The search was performed on 19 June, 2017. For each trait, all genetic instruments achieved the conventional levels of genome-wide significance (i.e.,  $P < 5.0 \times 10^{-8}$ ) and were mutually independent (i.e.,  $r^2 < 0.001$  between all pairs of instruments).

**Two-sample MR.** For the exposure phenotypes with at least one genetic instrument available, we used two-sample MR analysis to evaluate their causal effects on osteoarthritis risk. The exceptions were smoking and reproductive traits, which were performed with one-sample MR only, because of the need to perform the analysis within specific subgroups. All summary association results used for two-sample MR are shown in Supplementary Table 23, and Supplementary Table 24 provides an overall description of each set of genetic instruments. We applied the following methods:

**Ratio method.** For exposure phenotypes with only one genetic instrument available, MR was performed with the ratio method, which consists of dividing the

instrument-outcome regression coefficient by the instrument-exposure regression coefficient. The standard error of the ratio estimate can be calculated by dividing the instrument-outcome standard error by the instrument-exposure regression coefficient. Confidence intervals and *P* values were calculated with the normal approximation.

**Inverse-variance weighting (IVW).** This method allows for combination of the ratio estimates from multiple instruments into a single pooled estimate. We used a multiplicative random-effects version of the method, which incorporates between-instrument heterogeneity in the confidence intervals.

**MR-Egger regression.** This method yields consistent causal-effect estimates even if all instruments are invalid, provided that horizontal pleiotropic effects are uncorrelated with instrument strength (i.e., the instrument strength independent of direct effects INSIDE) assumption holds).

**Weighted median.** This method allows for consistent causal-effect estimation even if the INSIDE assumption is violated, provided that up to (but not including) 50% of the weights in the analysis come from invalid instruments.

**Mode-based estimate (MBE).** The weighted MBE relies on the zero modal pleiotropy assumption (ZEMPA), which postulates that the largest subgroup (or the subgroup carrying the largest amount of weight in the analysis) of instruments that estimate the same causal-effect estimate is composed of valid instruments. This procedure allows for consistent causal-effect estimation even if most instruments are invalid. The stringency of the method can be regulated by the  $\phi$  parameter. We tested two values of  $\phi$ :  $\phi = 1$  (i.e., the default) and  $\phi = 0.5$  (half the default, or twice as stringent).

For exposure phenotypes with more than one but fewer than ten genetic instruments, only the IVW method was applied, because the remaining methods are typically less powered and require a relatively large number of genetic instruments to provide reliable results. The degree of weak instrument bias (which corresponds to regression-dilution bias in two-sample MR) for the IVW and MR-Egger methods was quantified with the  $\frac{F_{XG}-1}{F_{XG}}$  and  $I_{GX}^2$  statistics, respectively.

Both range from 0% to 100%, and  $100 \left(1 - \frac{F_{XG}-1}{F_{XG}}\right)\%$  and  $100(1 - I_{GX}^2)\%$  can be

interpreted as the amount of dilution in the corresponding causal-effect estimates. Given that only genome-wide-significant variants were selected as instruments, the  $\frac{F_{XG}-1}{F_{XG}}$  statistic was necessarily high (at least ~95%). However, the  $I_{GX}^2$  statistic depends on both instrument strength and heterogeneity between instrument-exposure associations, thus suggesting that regression dilution bias in MR-Egger can be substantial even if instruments are individually strong. Indeed, for some traits, the  $I_{GX}^2$  statistic was very low (Supplementary Table 24). Therefore, all MR-Egger regression analyses were corrected for regression dilution with a simulation extrapolation (SIMEX) approach.

**Horizontal pleiotropy tests.** We additionally assessed the robustness of our findings to potential violations of the assumption of no horizontal pleiotropy by applying two tests of horizontal pleiotropy. One test was the MR-Egger intercept, which can be interpreted as the average instrument-outcome coefficient when the instrument-exposure coefficient is zero. If there is no horizontal pleiotropy, the intercept should be zero. Therefore, the intercept provides an indication of overall unbalanced horizontal pleiotropy. The second test was Cochran's *Q* test of heterogeneity, which relies on the assumption that all valid genetic instruments estimate the same causal effect.

**Power calculations.** We performed power calculations to estimate the power of our two-sample MR analysis to detect odds ratios of 1.2, 1.5 and 2.0 (Supplementary Note).

**One-sample MR.** UK Biobank data were used to perform one-sample MR with the same genetic instruments as in the two-sample MR (Supplementary Note).

**Life Sciences Reporting Summary.** Further information on experimental design is available in the Life Sciences Reporting Summary.

**Data availability.** All RNA-sequencing data have been deposited in the European Genome/Phenome Archive (cohort 1, EGAD00001001331; cohort 2, EGAD00001003355; cohort 3, EGAD00001003354).

## References

- 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).
- Willer, C. J., Li, Y. & Abecasis, G. R. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* **26**, 2190–2191 (2010).
- Franklin, J., Ingvarsson, T., Englund, M. & Lohmander, S. Association between occupation and knee and hip replacement due to osteoarthritis: a case-control study. *Arthritis Res. Ther.* **12**, R102 (2010).

42. Styrkarsdottir, U. et al. Whole-genome sequencing identifies rare genotypes in COMP and CHADL associated with high risk of hip osteoarthritis. *Nat. Genet.* **49**, 801–805 (2017).
43. Wesseling, J. et al. CHECK (Cohort Hip and Cohort Knee): similarities and differences with the Osteoarthritis Initiative. *Ann. Rheum. Dis.* **68**, 1413–1419 (2009).
44. Steinberg, J. et al. Integrative epigenomics, transcriptomics and proteomics of patient chondrocytes reveal genes and pathways involved in osteoarthritis. *Sci. Rep.* **7**, 8935 (2017).
45. Ritchie, M. E. et al. limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.* **43**, e47 (2015).
46. Li, H. et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
47. Tischler, G. & Leonard, S. biobambam: tools for read pair collation based algorithms on BAM files. *Source Code Biol. Med.* **9**, 13 (2014).
48. Patro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon provides fast and bias-aware quantification of transcript expression. *Nat. Methods* **14**, 417–419 (2017).
49. Law, C. W., Chen, Y., Shi, W. & Smyth, G. K. voom: precision weights unlock linear model analysis tools for RNA-seq read counts. *Genome Biol.* **15**, R29 (2014).
50. UK10K Consortium. et al. The UK10K project identifies rare variants in health and disease. *Nature* **526**, 82–90 (2015).
51. 1000 Genomes Project Consortium. et al. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
52. Chen, W. et al. Fine mapping causal variants with an approximate Bayesian method using marginal test statistics. *Genetics* **200**, 719–736 (2015).
53. Kircher, M. et al. A general framework for estimating the relative pathogenicity of human genetic variants. *Nat. Genet.* **46**, 310–315 (2014).
54. Ionita-Laza, I., McCallum, K., Xu, B. & Buxbaum, J. D. A spectral approach integrating functional genomic annotations for coding and noncoding variants. *Nat. Genet.* **48**, 214–220 (2016).
55. Iotchkova, V. et al. Discovery and refinement of genetic loci associated with cardiometabolic risk using dense imputation maps. *Nat. Genet.* **48**, 1303–1312 (2016).
56. Zheng, J. et al. LD Hub: a centralized database and web interface to perform LD score regression that maximizes the potential of summary level GWAS data for SNP heritability and genetic correlation analysis. *Bioinformatics* **33**, 272–279 (2017).
57. Hemani, G. et al. MR-Base: a platform for systematic causal inference across the phenome using billions of genetic associations. Preprint at <https://www.biorxiv.org/content/early/2016/12/16/078972> (2016).
58. Barban, N. et al. Genome-wide analysis identifies 12 loci influencing human reproductive behavior. *Nat. Genet.* **48**, 1462–1472 (2016).

## Life Sciences Reporting Summary

Nature Research wishes to improve the reproducibility of the work we publish. This form is published with all life science papers and is intended to promote consistency and transparency in reporting. All life sciences submissions use this form; while some list items might not apply to an individual manuscript, all fields must be completed for clarity.

For further information on the points included in this form, see [Reporting Life Sciences Research](#). For further information on Nature Research policies, including our [data availability policy](#), see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

### ▶ Experimental design

#### 1. Sample size

Describe how sample size was determined.

We used the data provided by the interim release of the UK Biobank resource, which is a large population-based study of over 152,000 genotyped subjects, to conduct the largest osteoarthritis (OA) genome wide association study (GWAS) to date. We defined 5 OA case strata based on both self-reported status and through linkage to Hospital Episode Statistics data, and on joint-specificity of disease (knee and/or hip). We used the full set of available OA cases with directly genotyped and imputed data. Non-OA controls were drawn from the same UK Biobank dataset. We applied exclusion criteria to minimise misclassification in the control dataset which was selected to be ~4x the number of cases to preserve power for common alleles while avoiding case:control imbalance causing association tests to misbehave for low frequency variants. For each of the 5 OA definitions, we performed power calculations and found that we have >80% power to detect variants at genome-wide significance ( $P < 5 \times 10^{-8}$ ) with modest effect size (odds ratio of 1.15 to 1.50) for common variants (minor allele frequency (MAF) 0.5 to 0.2). The power of this study is very limited for low frequency variants with  $OR < 1.50$ , and for rare variants.

#### 2. Data exclusions

Describe any data exclusions.

Sample and variant quality control (QC) was carried out centrally. We performed additional QC checks with exclusion criteria as follows: i) call rate  $\leq 97\%$  ii) gender discrepancies, iii) excess heterozygosity, iv) duplicates and/or high relatedness, v) ethnicity outliers, vi) possibly contaminated and vii) withdrawn samples. We excluded 528 SNPs that had been centrally flagged as subject to exclusion due to failure in one or more additional quality metrics. At the SNP level, we further filtered for Hardy Weinberg equilibrium (HWE)  $P \leq 10^{-6}$ ,  $MAF \leq 0.001$  and info score  $< 0.4$ . All the above are typical QC steps for quality assurance in GWAS data. We also restricted the number of controls used and did not utilise the full set of available genotyped control samples from UK Biobank in order to guard against association test statistics behaving anti-conservatively in the presence of stark case: control imbalance for alleles with minor allele count (MAC)  $< 400$ .

#### 3. Replication

Describe whether the experimental findings were reliably reproduced.

Our findings were reliably reproduced. We took 173 variants with  $P < 1.0 \times 10^{-5}$  and  $MAF > 0.01$  forward to replication in an independent cohort from Iceland (deCODE) of up to 18,069 cases and 246,293 controls. Following meta-analysis in up to 30,727 cases and 297,191 controls, we report seven genome-wide significant associations at novel loci, and two further new replicating signals just below the genome-wide significance threshold ( $P < 6.0 \times 10^{-8}$ ).

#### 4. Randomization

Describe how samples/organisms/participants were allocated into experimental groups.

Randomization of experimental groups were not required to this study. Participants were allocated into experimental groups according to their OA

## 5. Blinding

Describe whether the investigators were blinded to group allocation during data collection and/or analysis.

Blinding was not applicable to this study.

Note: all studies involving animals and/or human research participants must disclose whether blinding and randomization were used.

## 6. Statistical parameters

For all figures and tables that use statistical methods, confirm that the following items are present in relevant figure legends (or the Methods section if additional space is needed).

- |                          |  |
|--------------------------|--|
| n/a                      | Confirmed  |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> The <u>exact</u> sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement (animals, litters, cultures, etc.)                                    |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A description of how samples were collected, noting whether measurements were taken from distinct samples or whether the same sample was measured repeatedly.  |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A statement indicating how many times each experiment was replicated   |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> The statistical test(s) used and whether they are one- or two-sided (note: only common tests should be described solely by name; more complex techniques should be described in the Methods section) |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A description of any assumptions or corrections, such as an adjustment for multiple comparisons  |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> The test results (e.g. $p$ values) given as exact values whenever possible and with confidence intervals noted   |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A summary of the descriptive statistics, including central tendency (e.g. median, mean) and variation (e.g. standard deviation, interquartile range)   |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> Clearly defined error bars   |

See the web collection on [statistics for biologists](#) for further resources and guidance.

## ► Software

Policy information about [availability of computer code](#)

## 7. Software

Describe the software used to analyze the data in this study.

We used SNPTTEST v2.5.245 for association analysis. Power calculations were carried out using Quanto v1.2.4. We performed a fixed effects inverse-variance weighted meta-analysis in METAL. The results of the discovery cohorts of from RS-I, RS-II and CHECK were quality checked using EASYQC. To investigate the narrow sense heritability and the genetic correlation between the five osteoarthritis disease definitions, we ran the LDscore method. We tested proteins for differential abundance using limma in R. The cram rna-seq files were converted to bam files using samtools 1.3.157 and then to fastq files using biobambam 0.0.191. We obtained transcript-level quantification using salmon 0.8.259 and the GRCh38 cDNA assembly release 87 downloaded from Ensembl. We used tximport to convert transcript-level to gene-level count estimates. Limma-voom was used to remove heteroscedasticity from the estimated expression data. We tested genes for differential expression using limma in R. For fine-mapping we used CAVIARBF. We used GARFIELD for functional enrichment analysis. Gene-based and gene-set analyses were performed using MAGMA v1.06. We used DAPPLE for visualization of the pathways and protein-protein interaction (PPIs) relationships among the genes in each gene-set by integrating data from the InWeb database. LDHub was used to estimate the genome-wide genetic correlation between each of the OA definitions and 219 other human traits and diseases. In Mendelian randomization (MR) analysis, genetic instruments were identified from publicly-available summary GWAS results through the TwoSampleMR R package, which allows extracting the data available in the MR-Base database.

For all studies, we encourage code deposition in a community repository (e.g. GitHub). Authors must make computer code available to editors and reviewers upon request. The *Nature Methods* [guidance for providing algorithms and software for publication](#) may be useful for any submission.

## ► Materials and reagents

Policy information about [availability of materials](#)

### 8. Materials availability

Indicate whether there are restrictions on availability of unique materials or if these materials are only available for distribution by a for-profit company.

No unique materials used.

### 9. Antibodies

Describe the antibodies used and how they were validated for use in the system under study (i.e. assay and species).

No antibodies were used.

### 10. Eukaryotic cell lines

a. State the source of each eukaryotic cell line used.

No eukaryotic cell lines were used.

b. Describe the method of cell line authentication used.

No eukaryotic cell lines were used.

c. Report whether the cell lines were tested for mycoplasma contamination.

No eukaryotic cell lines were used.

d. If any of the cell lines used in the paper are listed in the database of commonly misidentified cell lines maintained by [ICLAC](#), provide a scientific rationale for their use.

No commonly misidentified cell lines were used.

## ► Animals and human research participants

Policy information about [studies involving animals](#); when reporting animal research, follow the [ARRIVE guidelines](#)

### 11. Description of research animals

Provide details on animals and/or animal-derived materials used in the study.

Animals were not used in this study.

Policy information about [studies involving human research participants](#)

### 12. Description of human research participants

Describe the covariate-relevant population characteristics of the human research participants.

This information can be found in the online methods:  
 1) Discovery GWAS  
 2) Replication  
 3) Association with OA-related endophenotypes  
 4) Functional genomics  
 Consent was obtained for each individual as stated in each of the above sections.