# **Decomposing the genetic background of chronic back pain**

| 2  | Elizaveta E. Elgaeva <sup>1,2*</sup> , Irina V. Zorkoltseva <sup>1</sup> , Arina V. Nostaeva <sup>2</sup> , Dmitrii A. Verzun <sup>1,2</sup> , Evgeny |
|----|---|
| 3  | S. Tiys <sup>1</sup> , Anna N. Timoshchuk <sup>3,4</sup> , Anatoliy V. Kirichenko <sup>1</sup> , Gulnara R. Svishcheva <sup>1,5</sup> , Maxim B.      |
| 4  | Freidin <sup>6</sup> , Frances M. K. Williams <sup>7</sup> , Pradeep Suri <sup>8,9</sup> , Yurii S. Aulchenko <sup>1,10</sup> , Tatiana I.            |
| 5  | Axenovich <sup>1</sup> , Yakov A. Tsepilov <sup>1*</sup>  |
| 6  | 1 - Institute of Cytology and Genetics, Siberian Branch of Russian Academy of Sciences,   |
| 7  | Novosibirsk, Russia   |
| 8  | 2 – Novosibirsk State University, Novosibirsk, Russia   |
| 9  | 3 - MSU Institute for Artificial Intelligence, Lomonosov Moscow State University, Moscow,   |
| 10 | Russia  |
| 11 | 4 – Moscow institute of physics and technology, Moscow, Russia  |
| 12 | 5 – Vavilov Institute of General Genetics, RAS, Moscow, Russia  |
| 13 | 6 - Department of Biology, School of Biological and Behavioural Sciences, Queen Mary  |
| 14 | University of London, London, UK  |
| 15 | 7 - Department of Twin Research and Genetic Epidemiology, School of Life Course Sciences,   |
| 16 | King's College London, London, UK   |
| 17 | 8 – Department of Rehabilitation Medicine, University of Washington, Seattle, USA   |
| 18 | 9 – VA Puget Sound Health Care System, Seattle, USA   |
| 19 | 10 – PolyOmica, The Netherlands   |
| 20 | *correspondence to: tsepilov@bionet.nsc.ru, elizabeth.elgaeva@gmail.com   |
|    |   |

- 21 Key words: genome-wide association study, GWAS, shared genetic background,
- 22 subphenotyping, polygenic risk score, PRS, low back pain

## 23 Abstract

Chronic back pain (CBP) is a disabling condition with a lifetime prevalence of 40% and a substantial socioeconomic burden. Because of the high heterogeneity of CBP, subphenotyping may be necessary to improve prediction and support personalized treatment for those with CBP. The lack of distinct cellular and molecular markers for CBP complicates the task of subphenotyping.

29 To investigate CBP subphenotypes, we decomposed the genetic background of CBP into a shared genetic background common to other chronic pain conditions (back, neck, hip, knee, stomach, and 30 31 head pain) and unshared genetic background related only to CBP. We showed that the shared and unshared genetic backgrounds of CBP differ in their biological functions: the first one is likely to 32 control processes mainly in nervous, immune and musculoskeletal systems underlying chronic 33 34 pain development regardless its site, while the second may contribute more to local processes in spine leading to chronic pain precisely in the back. We identified 18 genes with shared impact 35 across different chronic pain conditions and two genes that were specific for CBP. These findings 36 may contribute to future development of targets and new biomarkers for chronic pain management. 37 38 Next, among people with CBP, we demonstrated that polygenic risk scores accounting for the 39 shared and unshared genetic backgrounds of CBP may underpin different subphenotypes of CBP cases. These subphenotypes are characterized by varying genetic predisposition to a wide array of 40 medical conditions and interventions such as diabetes mellitus, myocardial infarction, diagnostic 41 42 endoscopic procedures, and surgery involving muscles, bones, and joints. The proposed genetic 43 decomposition framework holds promise for investigating the genetic underpinnings of other heterogeneous diseases. 44

## 45 Author Summary

Chronic back pain (CBP) is a prevalent disabling health problem with heterogeneous clinical 46 presentation and natural history. This may contribute to generic pain treatment approaches not 47 sufficiently effective when prescribed for patients with certain characteristics. Development of 48 49 more personalized treatment is needed, and may benefit from a deep understanding of CBP biology, such as genomics. It is known that chronic pain is under the control of both environmental 50 51 and genetic background. Here we applied bioinformatic methods to study the genetic background 52 of CBP decomposed into two parts: a shared one common to six distinct chronic pain types, and unshared, which is specific to CBP. This approach allowed us to identify more genes potentially 53 involved in CBP development. Among them 18 belong to the shared genetic background 54 55 contributing to development of chronic pain in general, and two are specific for CBP. We 56 demonstrate that these two parts of the genetic background of CBP are associated with distinct biological pathways and underlie predisposition to different medical states and procedures, 57 involving diabetes, myocardial infarction, and musculoskeletal surgery. Decomposition of CBP 58 genetic background into shared and unshared may provide a better understanding of mechanisms 59 60 of CBP and facilitate development of personalized pain treatment.

## 62 Introduction

Back pain is a prevalent clinical syndrome which affects about 40% of the population [1]. It has 63 tremendous social and economic consequences: according to the Global Burden of Disease Study 64 2016, back pain has been a major cause of disability worldwide for 30 years [2]. Back pain is not 65 66 only highly prevalent, but it is also difficult to treat [3]. One potential cause of this problem is the high heterogeneity of the condition. Patients with a specific back pain "subphenotype" – a set of 67 common features, distinguishing them from other patients with back pain [4] – may respond to 68 69 treatment in a different way than a patient of another back pain subphenotype, decreasing the effectiveness of treatment approaches when not tailored to subphenotype. Generic treatment 70 approaches may contribute to back pain patient care costs, which reach 1/5 of total health care 71 costs in a separate country [5]. Subphenotyping may help to more accurately select treatment for 72 73 patients with back pain and thus decrease these costs.

In 10% of cases acute back pain ceases to be just a symptom and becomes a chronic condition [6]. While the initial cause of acute pain may resolve, an alternative pathophysiological process takes over, leading to anatomical changes and affecting human behavior and mental state [7–9]. Chronic back pain (CBP) has been shown to be a complex trait with heritability estimated between 30 and 68% [10–12]. Genome-wide association studies (GWAS) of chronic pain including but not limited to CBP have revealed about three dozen associated loci [12–17], but only a few of them have been replicated in independent samples.

Studies have demonstrated the presence of a shared genetic background of chronic pain across different pain sites [16,18,19]. This shared genetic background is thought to condition the generation, transduction and processing of pain stimuli in general [16]. At the same time, it is

84 reasonable to expect the existence of unshared genetic background specific to chronic back pain (CBP) and related to local pathological processes in the back and spine (e. g. through SOX5 gene 85 [15]), but this has been little studied. Genetic variants related to unshared genetic background of 86 CBP may suggest promising drug targets and biomarkers for diagnosis and treatment of CBP 87 precisely, while the shared genes could be of interest to understand general pain biology and hence 88 89 contribute in future development of drugs to treat pain irrespective of its site. Therefore, the decomposition of genetic background into shared and unshared may be helpful for subphenotyping 90 and personalizing treatment. 91

Many prior attempts to identify subphenotypes of back pain have been made. Existing approaches 92 93 to do this classify patients based on observable clinical characteristics such as pain-related, social, physiological and anatomic features [20] and some of these studies show the utility of this 94 95 approach for revealing clusters of patients with specific pain trajectories and differentially reacting to treatment [21,22]. A novel and alternative approach is to subphenotype CBP using genotypic 96 information, which is increasingly available in clinical care and commercial use [23]. Polygenic 97 risk scores (PRS) use GWAS data to estimate the personal genetic liability to disease based on an 98 individual's genotype. In medicine, PRSs have been used to predict disease, however, this 99 100 information can also be utilized for subphenotyping. While existing chronic back pain PRS models 101 show modest prediction ability [24], division of CBP genetic background into shared and unshared, followed by PRS calculation for each of them, may result in higher statistical power and provide 102 103 useful information for patients' subphenotyping. To estimate the potential of genetic background 104 decomposition for subphenotyping, a comprehensive examination of the features of shared and 105 unshared genetic background is required.

106 Previously, we have studied the shared genetic background of chronic musculoskeletal pain at different sites (back, neck, hip, knee) by applying principal component analysis of these traits [16]. 107 We interpreted the first genetically independent phenotype (GIP1, or simply the first principal 108 109 component) as an approximation of shared genetic background across the considered pain types. 110 In our recent study we developed the novel SHAHER framework, [25] allowing more accurate 111 decomposition of genetic background of correlated traits in order to evaluate not only their shared, but unshared genetic factors as well. Here we aimed to employ SHAHER to investigate the 112 complex genetic architecture of CBP by analyzing its shared genetic background across several 113 114 chronic musculoskeletal and non-musculoskeletal pain conditions (back, neck, hip, knee, stomach and head pain) and its unshared genetic background, particular to chronic pain in the back. We 115 conducted an extensive bioinformatic annotation of shared and unshared genetic background of 116 117 CBP to reveal functional differences and identify genes associated with each aspect. To find the associated genes, we additionally performed a gene-based association analysis. Finally, we 118 119 performed a set of PRS analyses among people with CBP to check whether the division between 120 shared and unshared genetic background might be beneficial for CBP subphenotyping.

## 121 **Results**

#### 122 Overview of the study design

This study was carried out using the results of genome-wide association studies (GWAS) for six 123 chronic pain sites (pain in the back, neck, hip, knee, stomach, and headache). All data ( $N_{total} =$ 124 456,000) were provided by UK Biobank under projects #18219 and #59345. The sample was split 125 126 into discovery (265,000 individuals of European descent) and replication (a total of 191,000 127 individuals in three subsamples of African, South Asian and European descent) samples. Details 128 of the phenotypes definition are available in Supplementary Methods. Sample characteristics (size, 129 sex and age structure, pain type prevalence and BMI distribution) are provided in Supplementary Table 1. 130

131 The study design included four stages. In the first stage, we decomposed the genetic background 132 of CBP into shared and unshared ones, and calculated the GWAS summary statistics for two new 133 traits: SGIT (shared genetic impact trait) which is controlled predominantly by shared genetic background common for all traits analyzed, and UGIT (unshared genetic impact trait) which is 134 135 controlled by genetic background specific for CBP. These calculations were performed for each 136 subsample using the SHAHER framework (Figure 1, Materials and Methods). Then we performed 137 two meta-analyses for SGIT and UGIT using inverse-variance-weighted method: one combining the summary statistics from the two European samples (European meta-analysis) and another one 138 139 combining summary statistics from the three replication samples (Replication meta-analysis). 140 Thus, four sets of GWAS summary statistics both for SGIT and UGIT were obtained in the first stage: the "Discovery sample", the "European replication sample", the "European meta-analysis" 141 and the "Replication meta-analysis" (Figure 1). 142



Figure 1. The first step of the study. The black frames label the analyzed sample and the blue frames indicate the input data; the black text in bold shows intermediate results; the text written in blue bold highlights the output data, and the black text written in italic represents the type of analysis.

148

In the second stage, we identified loci and genes associated with SGIT and UGIT. We did this using two approaches. The first approach (Figure 2) was identification of associated loci in the discovery sample, followed by selection of independent association signals within these loci utilizing conditional analysis. For those loci that were previously observed and replicated in our recent work utilizing the same data [15,16], we did not perform additional analyses. For other

154 association signals, we conducted replication using the Replication meta-analysis. For loci replicated in the Replication meta-analysis we carried out gene prioritization procedure using the 155 European meta-analysis. Gene prioritization included prediction of SNP effects in replicated loci 156 (VEP and FATHMM), analysis of colocalization with gene expression effects (SMR-HEIDI), and 157 gene prioritization using DEPICT and FUMA. Details are available in Supplementary Methods. 158 159 The second approach utilizing which we identified genetic factors associated with shared and unshared backgrounds of CBP (Figure 2) was a gene-based association analysis (SKAT-O [26], 160 PCA [27], ACAT-V [28], and ACAT-O [28] methods) using the GWAS summary statistics 161 obtained from the discovery sample, followed by a conditional analysis and replication utilizing 162 the gene-based association analysis results from the European replication sample and European 163 meta-analysis (Figure 2, Materials and Methods). 164



165

Figure 2. The second step of the study. The black frames label the analyzed sample and the blueframes indicate the input data; the black text in bold shows intermediate results; the text written

in blue bold highlights the output data, and the black text written in italic represents the type of
 analysis.

170

At the third stage (Figure 3), we investigated the genetic architecture of CBP using functional 171 bioinformatic analyses. We conducted a gene set and cell type/tissue enrichment analyses of the 172 identified genes and SNPs in the associated loci in order to characterize the shared and unshared 173 genetic background of CBP. Using GWAS summary statistics obtained for SGIT and UGIT, we 174 175 estimated the genetic correlations of these traits with a large number of complex human traits from 176 the GWAS-MAP database [29,30]. Then we calculated polygenic risk scores (PRS) of SGIT and 177 UGIT using individual genotype data and estimated their role in disease/medical intervention 178 prediction.



Figure 3. The third and the fourth steps of the study. The black frames label the analyzed sampleand the blue frames indicate the input data; the black text in bold shows intermediate results; the

text written in blue bold highlights the output data, and the black text written in italic represents
the type of analysis.

184

Finally, we investigated the potential of SGIT and UGIT PRSs for CBP subphenotyping (Figure 3, Materials and Methods) using participants with CBP from the entire European sample (not divided into discovery and replication). We built binary-coded PRSs for SGIT and UGIT by splitting the PRS distribution into two parts according to the lowest or highest decile (see Materials and Methods). Then we tested if the binary-coded PRSs can be used for CBP subphenotyping.

## 190 Decomposition of genetic background of CBP

191 We estimated the heritability, the genetic correlation and phenotypic correlation for six chronic 192 pain traits from the discovery sample (see Supplementary Methods and Supplementary Tables ST2a, b). Using these estimates and applying the SHAHER framework, we calculated the 193 194 coefficients of optimal linear combinations for SGIT and UGIT building (Supplementary Table ST2a) and obtained the summary statistics for SGIT and UGIT. Heritability estimates and 195 correlation coefficients of the original pain traits, SGIT and UGIT are depicted in Figure 4. 196 Notably, the SNP-based heritability of SGIT is almost two times higher than those of the original 197 traits (0.07 versus 0.01 - 0.04). The phenotypic and genetic correlations between SGIT and the 198 original traits are higher in comparison to those between the original traits. UGIT shows positive 199 200 correlation only with CBP.



202

Figure 4. Heatmap visualization of the heritability, phenotypic correlations and genetic correlations of the original pain traits, SGIT and UGIT. Inscriptions name anatomic sites of chronic pain, UGIT refers to unshared genetic background of CBP. The upper triangle of the matrix represents pairwise phenotypic correlations. The lower triangular of the matrix contains pairwise genetic correlations. Statistically insignificant (p-value > 0.05/28) genetic correlation coefficients are crossed out. Diagonal elements of the matrix correspond to SNP-based heritability (h<sup>2</sup>) estimates of the particular traits.

210

## 211 Gene identification

#### 212 Identification of loci associated with SGIT and UGIT

The results of SGIT and UGIT GWASs obtained from the discovery sample are visualized at Manhattan and QQ plots (Supplementary.Figures 1 - 4). We identified five loci statistically

| 215 | significantly (p-value < 8.3e-09) associated with SGIT (see Table 1). Conditional analysis showed    |
|-----|--|
| 216 | that each of these loci contains only one independent signal. Four loci (tagged with the following   |
| 217 | lead SNPs: rs9436127, rs7652179, rs12705966 and rs73581580) were previously observed and             |
| 218 | replicated in our recent study [16] in the first genetically independent phenotype (GIP1). The locus |
| 219 | with the corresponding lead SNP rs11079993 was the new one. Similarly, we identified one locus       |
| 220 | with rs1271351 as the leading SNP significantly associated with UGIT. Conditional analysis           |
| 221 | showed a single independent association signal in this region. Since this locus has been already     |
| 222 | reported in [15], we interpreted it as known. Detailed information on significant loci and their     |
| 223 | replication is available in Supplementary Table ST3.   |

**Table 1**. Loci associated with SGIT and UGIT.

| Lead SNP   | Chr:position   | RefA/<br>EffA  | Nearest<br>gene   | Discovery sample  |   |   | Replication meta-analysis   |   |   |  |
|------------|--|--|---|---|---|---|---|---|---|--|
|            |  |  |   | β   | se  | p-value   | EAF   | β   | se  | p-value  |
| rs9436127  | 1:150490565  | G/A  | ECM1  | -0.018  | 0.003   | 1.72e-10  | 0.40  | -0.012  | 0.003   | 5.98e-04   |
| rs7652179  | 3:49808618   | T/A  | AMIGO3  | 0.023   | 0.004   | 2.30e-10  | 0.18  | 0.011   | 0.004   | 1.00e-02   |
| rs12705966 | 7:114248851  | G/A  | FOXP2   | 0.017   | 0.003   | 6.55e-09  | 0.67  | 0.011   | 0.004   | 2.16e-03   |
| rs73581580 | 9:140251458  | G/A  | EXD3  | 0.026   | 0.004   | 1.29e-09  | 0.12  | 0.029   | 0.005   | 5.36e-09   |
| rs11079993 | 17:50301552  | T/G  | snoZ178   | -0.017  | 0.003   | 4.98e-09  | 0.62  | -0.013  | 0.003   | 7.94e-05   |
| rs1271351  | 10:73798873  | T/C  | CHST3   | -0.018  | 0.003   | 1.13e-10  | 0.43  | -0.014  | 0.003   | 3.40e-05   |
|            | Lead SNP rs9436127 rs7652179 rs12705966 rs73581580 rs110799933 rs1271351 | Lead SNPChr:positionrs94361271:150490565rs76521793:49808618rs127059667:114248851rs735815809:140251458rs1107999317:50301552rs127135110:73798873 | Lead SNP         Chr:position         RefA/<br>EffA           rs9436127         1:150490565         G/A           rs7652179         3:49808618         T/A           rs12705966         7:114248851         G/A           rs73581580         9:140251458         G/A           rs12703993         17:50301552         T/G           rs1271351         10:73798873         T/C | Lead SNPChr:positionRefA/<br>EffANearest<br>geners94361271:150490565G/AECMIrs76521793:49808618T/AAMIGO3rs127059667:114248851G/AFOXP2rs735815809:140251458G/AEXD3rs1107999317:50301552T/GsnoZ178rs127135110:73798873T/CCHST3 | Lead SNP         Chr:position         RefA/<br>EffA         Nearest<br>gene         β           rs9436127         1:150490565         G/A         ECMI         -0.018           rs7652179         3:49808618         T/A         AMIGO3         0.023           rs12705966         7:114248851         G/A         FOXP2         0.017           rs73581580         9:140251458         G/A         EXD3         0.026           rs1271351         10:73798873         T/C         CHST3         -0.018 | Lead SNP         Chr:position         RefA/<br>EffA         Nearest<br>gene         Discover<br>β           rs9436127         1:150490565         G/A         ECMI         -0.018         0.003           rs7652179         3:49808618         T/A         AMIGO3         0.023         0.004           rs12705966         7:114248851         G/A         FOXP2         0.017         0.003           rs73581580         9:140251458         G/A         EXD3         0.026         0.004           rs11079993         17:50301552         T/G         snoZ178         -0.018         0.003           rs1271351         10:73798873         T/C         CHST3         -0.018         0.003 | Lead SNP         Arrest         RefA/EffA         Nearest         Image of the set | Lead SNPChr:positionRefA/<br>EffANearest<br>geneEiscover-samplers94361271:150490565G/A <i>ECM1</i> <b>Bp-valueEAF</b> rs76521793:49808618T/A <i>AMIGO3</i> 0.0230.0041.72e-100.401rs127059667:114248851G/A <i>FOXP2</i> 0.0170.0036.55e-090.671rs735815809:140251458G/A <i>EXD3</i> 0.0260.0041.29e-090.122rs127135110:73798873T/C <i>CHST3</i> 0.0180.0031.13e-100.431 | Lead SNPArrow Price<br>PriceNearest Break<br>geneDiscovery sampleReplicers9436127Li150490565G/AECM1βsep-valueEAFβrs76521793:49808618T/AAMIGO30.0230.0042.30e-100.400.011rs127059667:114248851G/AFOXP20.0170.0036.55e-090.6120.011rs735815809:140251458G/AEXD30.0260.0041.29e-090.120.029rs127135110:73798873T/CCHST30.0180.0031.13e-100.430.014 | Lead SNPArrest<br>FfrRefar<br>geneNearest<br>geneIsiscover sampleReplice sampleReplice samplers94361271:150490565G/AECM1-Bp-valueEAFB.0003rs76521793:49808618T/AAMIGO30.0230.0042.30e-100.180.0110.004rs12705967:114248851G/AFOXP20.0170.0036.55e-090.670.0110.004rs735815809:140251458G/AEXD30.0260.0041.29e-090.120.0120.005rs127135110:73798873T/GSmoZ178-0.0180.0031.13e-100.43-0.0140.003 |

The new replicated locus is shown in bold. Associations shown in grey were observed in previous

studies [15,16]. Lead SNP – the SNP with the lowest p-value in the locus; Chr:position – genomic

227 coordinates in a format "chromosome number : base pairs" (according to the GRCh37 genomic 228 build); RefA – reference (not effective) allele; EffA – effective allele;  $\beta$  – effect size of SNP 229 counted for effective allele; se – standard error of  $\beta$ ; p-value – p-value of association between SNP 230 and a trait after correction for genomic control; EAF – effective allele frequency.

231

## 232 Prioritization of genes in the new associated locus

233 We conducted gene prioritization for the rs11079993 locus associated with SGIT utilizing several approaches: literature-based prioritization of the genes located in this locus; prediction of 234 pathogenicity of SNP effects in the locus using VEP [31], FATHMM-XF [32] and FATHMM-235 236 INDEL [33]; DEPICT [34] and FUMA [35] gene prioritization and estimation of pleiotropic 237 effects on gene expression using SMR-HEIDI [36]. Gene prioritization suggested two genes in the rs11079993 locus (see Supplementary Results and Supplementary tables ST4-6 for more details): 238 CA10 (prioritized by literature-based annotation and FUMA) and LINC01982 (prioritized by 239 240 literature-based annotation). Although the latter is well studied it provided less evidence for 241 prioritization than CA10, so we consider it less likely to be causal. The CA10 gene may have an 242 effect on chronic pain through processes in the central nervous system, cancer or disease of bones. Despite the fact that *snoZ178* is the nearest gene to rs11079993, it did not provide strong arguments 243 244 for being causal.

#### 245 Gene-based association analysis

We performed the gene-based association analysis using different SNP annotation sets within genes. For SGIT we detected 9, 18, and 43 genome-wide significant gene-based signals for nonsynonymous, protein coding and protein non-coding SNP sets, respectively (Supplementary Table ST7). After conditional analysis, 10 out of 43 for non-coding SNP set and one for each

- coding and nonsynonymous SNP sets (both in *SLC39A8*) retained significance. We replicated five
  of 12 genes using the European replication sample and European meta-analysis (Table 2). For
  UGIT we detected two significant signals for the non-coding SNP set, and one of them (*CHST3*)
  passed conditional analysis and replication.
- **Table 2**. Results of gene-based association analysis of SGIT and UGIT

|       | Gene    | Chr                | Position  | SNP<br>set | Gene-based p-value  |                         |                                   |                               |  |
|-------|---------|--------------------|-----------|------------|---------------------|-------------------------|-----------------------------------|-------------------------------|--|
| Trait |         |                    |           |            | Discovery<br>sample | Conditional<br>analysis | European<br>replication<br>sample | European<br>meta-<br>analysis |  |
| UGIT  | CHST3   | 10                 | 71964395  | ncod       | 7.81e-08            | 7.81e-08*               | 1.42e-06                          | 7.93e-13                      |  |
|       | SLC39A8 | 4                  | 102251080 | nsyn       | 5.68e-10            | 5.68e-10*               | 4.73e-09                          | 4.38e-17                      |  |
|       | SLC39A8 | 4                  | 102251080 | cod        | 1.17e-07            | 1.17e-07*               | 1.62e-06                          | 2.54e-13                      |  |
|       | IP6K1   | 3                  | 49724294  | ncod       | 1.32e-10            | 1.32e-10*               | 0.0083                            | 2.66e-12                      |  |
|       | GRM3    | 7                  | 86643909  | ncod       | 1.66e-07            | 1.66e-07*               | 0.0935                            | 6.55e-08                      |  |
|       | FOXP2   | 7                  | 114086327 | ncod       | 2.00e-07            | 2.00e-07*               | 7.90e-04                          | 8.22e-11                      |  |
| SGIT  | MAML3   | 4                  | 139716753 | ncod       | 3.17e-07            | 3.17e-07*               | 6.40e-05                          | 1.30e-12                      |  |
|       | PTBP1   | 19                 | 797075    | ncod       | 3.64e-07            | 3.64e-07*               | 0.8760                            | 2.33e-07                      |  |
|       | TCF20   | 22                 | 42160013  | ncod       | 1.93e-06            | 1.93e-06*               | 3.77e-04                          | 8.44e-10                      |  |
|       | SOX6    | 11                 | 15966449  | ncod       | 2.01e-06            | 2.01e-06*               | 0.0921                            | 4.17e-07                      |  |
|       | FZD10   | FZD10 12 130162459 |           | ncod       | 2.03e-06            | 2.22e-06                | 0.2626                            | 4.36e-06                      |  |
|       | ERICH2  | 2                  | 170766878 | ncod       | 2.25e-06            | 2.25e-06*               | 0.1304                            | 3.19e-05                      |  |
|       | GABRB2  | 5                  | 161288429 | ncod       | 2.08e-06            | 2.40e-06                | 5.20e-05                          | 3.69e-11                      |  |

<sup>255 \*</sup> These genes harbored SNPs with the lowest p-values within 5 Mb from their borders, thus no conditional

258 number of genes retained significance after the conditional analysis.

<sup>256</sup> SNPs were selected, and the gene-based p-values remained unchanged.

<sup>257</sup> The significance threshold for replication sample was set at p-value < 4.2e-03 = 0.05/12, where 12 is the

259

## 260 Summary of gene identification

Using prioritization of genes in the associated loci and gene-based association analysis, we 261 262 identified 18 genes associated with SGIT, with 13 of them being previously reported in our recent 263 works [15,16] (SLC39A8, FOXP2, ECM1, AMIGO3, BSN, RBM6, FAM212A, RNF123, UBA7, MIR7114, NSMF, NOXA1, GRIN1), and 5 of them being new (CA10, MAML3, TCF20, GABRB2, 264 LINC01982). One of the identified genes, namely FOXP2, was found using both gene 265 prioritization and gene-based analysis. The SLC39A8 gene was identified using gene-based 266 analysis and has been prioritized in our previous study using the genetically independent 267 268 phenotype (GIP) approach [16].

For UGIT, we identified the *CHST3* gene utilizing gene-based analysis, and prioritized two genes
in the UGIT-associated locus: *CHST3* and *SPOCK2*, failing to give preference to one of them [15].

## 271 Analysis of CBP genetic architecture

#### 272 Gene set and tissue/cell type enrichment

We performed two enrichment analyses based on DEPICT for SNPs associated with SGIT and 273 274 UGIT in European meta-analysis, and FUMA analysis for the genes identified for SGIT. In 275 DEPICT analysis at p-value < 5e-06 we found SGIT to be enriched in genes involved in the 276 BTBD2 PPI subnetwork and decreased cochlear coiling (Supplementary Table ST8a). We also 277 detected enrichment for genes expressed in the central nervous and neurosecretory systems, retina and neural stem cells (Supplementary Table ST8b). Similar analyses of SGIT-associated variants 278 279 at p-value < 2.5e-08 threshold provided no significant findings. For UGIT we were able to conduct 280 DEPICT analysis only for SNPs with p-value < 5e-06, but no statistically significant results were

obtained. However, we observed a tendency to enrichment for genes expressed in the digestive,
nervous, musculoskeletal, cardiovascular, hemic and immune systems.

We did not obtain statistically significant results in tissue specificity and gene set enrichment 283 analyses using FUMA, however, we have shown that SGIT-associated genes have different 284 expression patterns. For instance, there is a group of genes (NSMF, PTBP1, RBM6, IP6K1, UBA7, 285 ECM1, NOXA1, RNF123) expressed almost in every organ, and a group of genes expressed 286 287 predominantly in the brain (GRIN1, GABRB2, GRM3, BSN, and CA10) (Supplementary Figure 5). These findings are in line with those obtained for SGIT loci utilizing DEPICT. FUMA extends 288 DEPICT findings for SGIT, by accounting for genes identified in gene-based analysis as well and 289 290 providing more detailed information on expression of particular genes.

#### 291 Genetic correlation between SGIT, UGIT and complex traits

292 We estimated the genetic correlation between each of SGIT and UGIT, and 730 complex human 293 traits from the GWAS-Map database. SGIT was statistically significantly genetically correlated 294 with almost all of the 322 preselected complex traits, and UGIT provided significant correlations with only 14 of them (see Supplementary Table ST9). In total we grouped all traits in 11 clusters 295 (see Supplementary Table ST9, Figure 5). There were some distinctions between SGIT and UGIT 296 297 correlation patterns. First, UGIT was significantly correlated with sitting height which represents the length of the spine, and the second, third and fourth genetically independent phenotypes (GIP2-298 4) from our recent study [16], while SGIT was not (Supplementary Table ST9). Second, SGIT had 299 300 a positive direction of genetic correlation with self-reported chronic knee pain and headache, when UGIT was negatively correlated with them. Regardless of statistical significance, SGIT and UGIT 301 302 demonstrated opposite patterns of genetic correlation (see Figure 5). While SGIT was mostly 303 positively correlated with traits from such clusters as: respiratory illness and smoking, injuries,

304 socio-economic and family status, psychometric traits, osteoarthritis and other musculoskeletal





Figure 5. Genetic correlation of SGIT and UGIT with complex human traits. Only 322 traits
 providing statistically significant genetic correlations (*rg*) with either SGIT or UGIT and *rg* magnitude greater than 0.25 were considered. These traits were grouped into 11 manually
 annotated clusters. From each cluster we selected one trait (top-trait) providing the smallest p value of genetic correlation among all pairwise correlations with SGIT and UGIT within this

312 313 cluster and depicted this value on a heatmap. Genetic correlations not passing the significance threshold of p-value < 9.41e-07 are crossed out.

314

## 315 Polygenic risk score for SGIT and UGIT. Their role in disease/medical intervention prediction

We developed a model of PRS estimation using a training sample. For a testing sample, we 316 317 calculated SGIT and UGIT PRSs using individual genotypes. Then we examined association between these PRSs and medical codes available for members of testing sample (see 318 Supplementary Table ST10). We detected 92 ICD10 codes and 57 OPCS4 codes statistically 319 320 significantly associated with at least one of the studied PRSs. SGIT PRS was associated with all of them. ICD10 codes associated with SGIT and/or UGIT PRSs related to a wide range of 321 disorders, such as disorders of musculoskeletal, nervous, cardio-vascular, digestive and excretory 322 323 systems, metabolic and skin disorders, diabetes, respiratory and neurological diseases (Supplementary Table ST10, Figure 6). The associated OPCS4 codes represented mostly surgical 324 325 manipulations on joints and bones and diagnostic endoscopic examination of gastrointestinal and 326 upper respiratory tract. Also, we observed SGIT PRS to be associated with all of the chronic pain phenotypes and these associations were characterized by the lowest p-values (Supplementary 327 328 Table ST10). Unlike SGIT PRS, UGIT PRS provided only three significant associations, which had been detected for SGIT PRS as well: associations with chronic back and knee pain and with 329 A52 OPCS4 code denoting therapeutic epidural injection, which is due to the association with 330 331 chronic back pain (Supplementary Table ST10, Figure 6).



Figure 6. Associations of SGIT and UGIT PRSs with ICD10 and OPCS4 codes and chronic pain
 phenotypes. Only 92 ICD10 and 57 OPCS4 codes providing statistically significant associations
 with either CBP, SGIT or UGIT PRS were considered. Medical codes were grouped into 19
 clusters. From each cluster we selected one code providing the smallest p-value of association

- with PRSs and depicted its effect sizes on a heatmap. Associations not passing the significance
   threshold of p-value < 8.25e-05 were crossed out.</li>
- 339

## 340 CBP subphenotyping among people with CBP

We investigated the potential of SGIT and UGIT binary-coded PRSs for CBP subphenotyping 341 using only the 65,011 participants with CBP. For each individual and each of two traits, SGIT and 342 343 UGIT, we formed two types of binary-coded PRSs defined by the lowest and highest deciles of 344 risk (see "Polygenic risk score calculation" in Supplementary Methods section). Then we performed an analysis of associations between each of the four binary-coded PRSs (with two 345 reflecting the lowest decile of the PRS for SGIT and UGIT, respectively, and the other two 346 347 reflecting the highest decile- for each of SGIT and UGIT) and the ICD10/OPCS4 medical codes and identified 69 ICD10 and 35 OPCS4 codes statistically significantly associated with at least 348 one of binary-coded PRSs (Supplementary Table ST11). 349

Being in the lowest decile of the SGIT PRS (low genetic predisposition to CBP through shared genetic background) provided a "unique" (not observed for other binary PRS traits) negative association with acute myocardial infarction (I21 ICD10 code, OR = 0.60, p-value = 1.09e-06). At the same time, the lowest decile of UGIT PRS (low genetic predisposition to CBP through the unshared genetic background) was uniquely positively associated with unspecified diabetes mellitus (E14 ICD10 code, OR = 1.64, p-value = 3.05e-05).

The highest decile of the SGIT PRS (high genetic predisposition to CBP through the shared genetic background) showed 102 statistically significant positive associations with the rest medical codes and 57 of them were associated only with this binary PRS trait. Among these 57 associated codes

359 were those related to diagnostic endoscopic examination of joints and rheumatoid arthritis, surgical manipulation on muscles, joints and bones (such as repair of muscle, bone excision, replacement 360 of joint etc.), conditions of the spine and nervous system, disorders of digestive, genitourinary and 361 362 endocrine systems and neurological disorders. No specific association was found with the highest 363 decile of the UGIT PRS, but differences in association pattern from other binary PRS traits were 364 found. For example, the SGIT PRS was positively associated with gastrointestinal tract examination procedures (may be because these people are generally prone to pain, so they have 365 unexplained abdominal pain as well or they may take NSAIDs), obesity, arthrosis of knee and hip, 366 367 and knee surgery, while the UGIT PRS was negative associated with these traits.

Generally, the lowest decile of SGIT PRS and the highest decile of UGIT PRS manifested protective effects on medical states being diagnosed and treated, whereas the lowest decile of UGIT PRS and the highest decile of SGIT PRS had the opposite effect (Supplementary Table ST11).

## 373 **Discussion**

In this work, we first applied the decomposition of genetic background of CBP into shared and unshared ones and showed that they differ in their functions. The shared genetic background is common to different chronic pain conditions while the unshared genetic background is related only to CBP. We built two traits, SGIT and UGIT, corresponding to the shared and unshared genetic background, respectively, and analyzed their properties.

We identified five loci associated with SGIT and one locus related to UGIT. Among them, only 379 380 one was new – the locus on chromosome 17 associated with CBP through SGIT. We prioritized 381 two genes (CA10 and LINC01982) near it which are potentially involved in brain development, synapse formation and carcinogenesis, and found associations with gastroesophageal reflux, bone 382 disease, multisite chronic pain and various known risk factors of chronic back pain (educational 383 384 attainment, smoking, depression). The other three loci found associated with SGIT had previously been reported for GIP1 [16], as expected given our assumption that GIP1 is an approximation of 385 the shared genetic background of the four chronic pain phenotypes studied before. In past work, 386 we have prioritized 12 genes (MIR7114, NOXA1, GRIN1, NSMF, FOXP2, BSN, AMIGO3, RBM6, 387 FAM212A, RNF123, UBA7, ECM1) in these three loci associated with SGIT, with six of them 388 389 (NOXA1, GRIN1, NSMF, FOXP2, BSN, AMIGO3) being related to the nervous system (brain 390 development and recovery, synapse plasticity, signal transduction, neuropathic and inflammatory pain). Some other genes were related to musculoskeletal (MIR7114 and ECM1, functional in 391 392 osteoarthritis and osteogenesis, respectively) and immune (UBA7) system processes.

In addition to five loci significantly associated with SGIT, there were four more loci suggestively
associated (p-value < 8.3e-08) with SGIT (Supplementary Table ST12). Among these four loci</li>

395 there was a locus of SLC39A8 gene reported in our previous study of GIP [16] and there was one new replicated under p-value < 8.3e-08 threshold locus on chromosome 13 with the lead SNP 396 rs2587363, located in the AL356295.1 gene (Supplementary Figure 6). We cannot consider this 397 novel locus as statistically significant in the present study, however, it may nevertheless be a target 398 for further investigation. This locus has been reported as associated with multisite chronic pain in 399 400 females [37], major depressive disorder, osteoarthritis, post-traumatic stress disorder [38,39], and chronic fatigue syndrome [40]. Moreover, this locus contained polymorphism rs2587363 that was 401 classified as pathogenic according to the FATHMM-XF (Supplementary Table ST13) and 402 403 demonstrated a pleiotropic effect on the OLFM4 gene expression (Supplementary Table ST14) in peripheral blood. This gene encodes olfactomedin 4, which is an antiapoptotic factor promoting 404 405 tumor growth.

406 Additionally, we performed gene-based association analyses and identified five genes (SLC39A8, 407 FOXP2, MAML3, TSF20, GABRB2) significantly associated with SGIT. As can be seen, FOXP2 was observed in this study both in GWAS and in the gene-based analysis. The SLC39A8 gene has 408 been reported and replicated for GIP1 [16]. The product of the SLC39A8 gene is a metal transporter 409 with a role in manganese (Mn) homeostasis, and the missense rs13107325 in this gene is among 410 411 the top pleiotropic SNPs identified in GWAS. Specifically, it has previously been associated with 412 increased risk of osteoarthritis [41] and severe adolescent idiopathic scoliosis [42]. The gene is known to be mutated in congenital disorder of glycosylation, SLC39A8-CDG, with clinical 413 414 features including osteopenia, broadened long bone epiphysis and joint hypermobility [43,44]. The 415 dietary manganese deficiency is known to lead to bone and connective tissue disease in animals 416 [45]. MAML3 has been shown to contribute to several pathways significantly associated with CBP 417 [46]. This gene enables transcription coactivator activity; it is involved in the Notch signaling

418 pathway and positive regulation of transcription by RNA polymerase II. GABRB2 encodes a GABA (gamma-aminobutyric acid) type A receptor beta subunit. The gene has a pivotal role in 419 the central nervous system and associates with various neuropsychiatric disorders [47]. TCF20 420 421 encodes a transcription factor that recognizes the platelet-derived growth factor-responsive element in the matrix metalloproteinase 3 promoter. The encoded protein is thought to be a 422 423 transcriptional coactivator, enhancing the activity of transcription factors such as JUN and SP1. Mutations in this gene are associated with autism spectrum disorders (according to NCBI-RefSeq 424 https://www.ncbi.nlm.nih.gov/gene?Db=gene&Cmd=DetailsSearch&Term=6942). Recently, it 425 426 has been shown that TCF20 is essential for neurogenesis during embryonic brain development in mouse. TCF20 dysfunction leads to deficits in neurogenesis, which further results in the 427 428 development of autism spectrum disorders [48].

429 Recently, more has been paid to the study of the role of rare genetic variants in the control of different traits. In the current study, we also tested the association of SGIT with exome sequencing 430 431 data from UK Biobank (Supplementary Methods) and performed gene-based association analysis, which identified the SLC13A1 gene (Supplementary Tables ST15a, b). The effect of this gene was 432 explained by its loss of function (LoF) and missense variants. The SLC13A1 gene has been 433 434 previously detected as associated with back pain-related traits due to LoF variants [14,49]. This gene encodes a protein that functions as a high-affinity sodium-dependent sulfate transmembrane 435 436 transporter [50]. A deficiency of the SLC13A1 protein is associated with a reduced blood sulfate 437 level, which plays a key role in the underlying processes leading to painful intervertebral disc disorders [14]. The association of a gene with disease through LoF variants presents several 438 439 advantages as a drug target because it provides a clear functional link, target specificity, potential 440 for gene therapy, predictive value, and personalized medicine opportunities. Until now, this gene

has been considered as a potential target for treatment against back pain. In the light of our results, *SLC13A1* can be considered as a potential target for treatment against chronic pain irrespective of
location.

For UGIT, the SPOCK2 and CHST3 genes were detected using both gene-based and gene 444 prioritization analysis. The association of SPOCK2 with back pain has been described previously 445 [15]. This gene is located close to carbohydrate sulfotransferase 3 (CHST3) which codes for an 446 447 enzyme that catalyzes proteoglycan sulfating and has been identified as a susceptibility gene for lumbar disc degeneration [51]. The association of CHST3 with back pain-related phenotypes has 448 also been described [14]. This gene can be considered as most probably causal for back pain in 449 450 this region, since cumulative evidence L2G scores for CHST3 and SPOCK2 are 0.81 and 0.23, respectively, https://genetics.opentargets.org/study-locus/NEALE2\_6159\_4/10\_72001257\_A\_G, 451 [52]. 452

Information on genes associated with SGIT is potentially useful in terms of discovering biomarkers and developing multi-target drugs against chronic pain at several anatomic sites at once, because these genes are most likely involved in the development of chronic pain in general. In contrast, genes related to the unshared genetic background of CBP may be more pertinent to the development of chronic pain at back precisely, probably by affecting chondrogenesis or degenerative processes in spine. So UGIT associated genes may be of interest for the development of drugs for treatment of chronic pain specifically in the back.

Based on enrichment analysis and functional annotation of associated loci, the shared genetic background common to different chronic pain conditions contributes to these phenotypes in general through processes in the central nervous and neurosecretory systems. SGIT genetic correlations provided evidence of genetically predisposed psychosocial aspects of CBP. The SNP-

464 based heritability for UGIT was much smaller than that of SGIT, resulting in much less statistical power for all analyses. However, we still detected some significant results. A few findings from 465 functional annotation demonstrated that UGIT was inversely associated with knee pain and 466 headache and positively associated with back pain and medical treatments for back pain such as 467 epidural injection. UGIT also tended to be negatively correlated with respiratory illnesses, 468 469 musculoskeletal disorders and injuries, demographic parameters and psychometric disorders, whereas SGIT is positively associated with all of these traits. These findings taken together support 470 that the unshared genetic background of CBP may determine the development of pain in the back 471 472 specifically, whereas the shared genetic background controls non-specific pain manifestations and processes accompanying different pain conditions. 473

Genetic correlation analysis and PRS analysis are two ways to investigate the influence of genetic 474 475 background of CBP. Although the lists of the traits included in these analyses overlap, there are 476 many traits presented only in one of them. Thus, genetic correlation analysis adds information on association with self-reported and anthropometric traits (e.g., UGIT is significantly correlated with 477 height), and PRS analysis provides more information on medical diagnoses and interventions (for 478 instance, both SGIT and UGIT PRSs are positively associated with therapeutic epidural injection). 479 480 Findings from the two methods complement each other. For example, genetic correlation analysis 481 of SGIT reveals genetic predisposition to greater alcohol intake and smoking, and in PRS analysis, we observe its association with mental disorders due to alcohol intake and smoking. Similarly, 482 483 SGIT analysis shows positive genetic correlation with toothache, and PRS analysis highlights 484 association with tooth removal. For group of traits included in both analyses (such as lower 485 respiratory and urinary system disorders, anxiety and depression, intestine and cardio-vascular 486 diseases, sleep and eye disorders, arthritis, spine and knee problems, etc.) results suggest positive

487 association with SGIT in both methods, but they should not be interpreted in the same way. While genetic correlation analysis provides information on pleiotropic effects or linkage in genome [53], 488 PRS analysis reveals associations between the genetic background of CBP and phenotypic traits 489 that can be explained not only by linkage or pleiotropic effects of genes, but also by causal effects 490 of CBP on other traits. Comparing results of genetic correlation analysis with those obtained from 491 492 PRS analysis can tell us more about probable reasons of observed association between two phenotypes [54]. For instance, when we observe modest genetic correlation between traits along 493 with high magnitude of regression coefficient from PRS analysis, this may indicate the role of 494 495 causal or confounding effects, rather than effects of genetic pleiotropy or linkage. Trait pairs with such a pattern may be an interesting subject for further studies utilizing Mendelian randomization. 496 497 For example, we observed this pattern for SGIT and several traits, including obesity (rg with BMI was 0.31, beta from PRS analysis for obesity was 0.25), diabetes (rg was approximately 0.31 for 498 set of related traits, PRS beta = 0.21), and migraine (rg = 0.37, PRS beta = 0.22). Causal effects of 499 increased BMI on back pain, and back pain on type II diabetes, have already been reported in our 500 studies [55,56], while the association with migraine is yet to be examined. In contrast, if the 501 magnitude of pairwise genetic correlation coefficient is quite high it is likely that the phenotypic 502 503 association between two traits is explained by similarity of their genetic background and not by 504 causality. For pairs of traits with these characteristics, information on the similarity of underlying 505 genetic mechanisms can be valuable for development or repurposing of medications affecting both 506 traits via the same pathway. From this prospective, further in-depth studies of the high genetic correlation between irritable bowel syndrome and SGIT (rg = 0.78) might be worthwhile. 507

The results of genetic analysis of SGIT and UGIT allowed us to estimate their polygenic risk scores
for participants with CBP and show that they are associated with different ICD10 and OPCS4

510 coded conditions. We introduced the binary-coded PRSs for SGIT and UGIT to characterize low and high genetic risk of CBP related to its shared and unshared genetic background, respectively. 511 We showed that these binary-coded PRSs can be potentially helpful for subphenotyping of 512 513 individuals with CBP, because participants with low or high genetic risk of CBP defined by the shared or unshared genetic background are characterized with different patterns of association with 514 515 medical interventions and disorders. This can be crucial for forecasting pain trajectories and choosing an adequate treatment. Low genetic liability to CBP explained by shared genetic 516 background turned out to be protective against diagnosis of acute myocardial infarction (OR =517 518 0.60, p-value = 1.09e-06), but low genetic risk of CBP mediated through its unshared genetic background was associated with higher predisposition to unspecified diabetes mellitus (OR = 1.64, 519 p-value = 3.05e-05). High PRS of SGIT significantly increased the risk of diverse diseases of 520 521 musculoskeletal, nervous, digestive, genitourinary and endocrine systems (ORs varying from 1.28) to 2.08). Genetic liability to CBP through shared genetic background also showed positive 522 523 association with gastrointestinal tract examination procedures, this may be because people are generally prone to pain, so they have unexplained abdominal pain as well or they may take non-524 steroidal anti-inflammatory drugs, which have side effects on gastrointestinal system. All these 525 526 results should be interpreted with caution, since negative/positive association between PRS and 527 medical code does not necessarily mean protective/detrimental effect on disease development itself but on the disease being diagnosed. For instance, low OR of myocardial infarction among people 528 529 with low SGIT PRS could potentially indicate that their low genetic liability to pain regardless of its site (including chest pain) can make myocardial infarction much harder to detect. Such a 530 531 phenomenon was observed for diabetics, who often do not have chest pain, so their myocardial 532 infarctions do not get diagnosed [57,58].

533 Our study has some limitations and restrictions. First, the discovery analyses were performed in Northern Europeans, so the obtained linear combination coefficients for SGIT and UGIT cannot 534 strictly be applied to individuals of African and Asian ancestry, included in replication sample. 535 536 Secondly, both discovery and replication sample were based on UK Biobank participants, which means that other cohorts are needed to perform independent replication of the results. Thirdly, here 537 538 we refer to UGIT as pertaining to unshared genetic background specific to CBP, however, the SHAHER approach does not necessarily accurately divide the genetic background of the traits, so 539 UGIT may have genetic correlations with some of the other chronic pain traits but not with all of 540 541 them. Finally, we had limited statistical power in UGIT functional analyses, which resulted in few significant findings. This indicates that larger sample size may be helpful in further studies. 542

Overall, the current work demonstrates that genetic background of CBP can be split into shared 543 544 between different pain conditions and specific for CBP background. The former is likely to 545 maintain pain mechanisms and manifestations in general through the central nervous and immune systems, functioning of joints, general processes of bone formation and remodeling, while the 546 latter is responsible for processes specific to back pain. Polygenic risk scores separately accounting 547 for shared and unshared genetic background of CBP identify subphenotypes characterized by 548 549 different genetic predisposition to a range of medical conditions making them possible tools for 550 prognosis in the healthcare system. Twenty genes prioritized for CBP may provide further opportunities for advances in chronic pain management, such as the discovery of new biomarkers 551 552 and drug targets for chronic pain regardless of its site and drug targets that are specific to CBP.

## 553 Materials and Methods

#### 554 **Data description**

We used the UK Biobank GWAS data ( $N_{total} = 456,000$ ) relating to four chronic musculoskeletal 555 pain traits (pain in the back, neck, hip, knee) from our recent work (see [16] for the detailed 556 description of the phenotype definition, sample size and sample characteristics). We obtained the 557 558 remaining two chronic pain phenotypes (stomach pain and headache) from the same UK Biobank 559 [59] sample and conducted GWAS according to the protocol we have described previously 560 (Supplementary Methods, [16]). For each pain trait we split the whole sample of UK Biobank 561 participants into a discovery (N = 265,000, European ancestry individuals) and a replication sample (three samples of African, N = 7,541, South Asian, N = 9,208, and European descent 562 563 individuals, N = 174,831). Details of the samples (size, sex and age structure, pain type prevalence and BMI distribution) are available in Supplementary Table ST1. Data were uploaded to the 564 GWAS-Map database [29,30] for quality control, unification and further functional analysis of 565 566 GWAS summary statistics. The same genotype and phenotype data were used for making and 567 testing PRS models.

## 568 Decomposition of genetic background of chronic back pain

To decompose the genetic background of CBP we implemented the SHAHER framework [25]. This approach is based on the concept that genetic background of each of the genetically related traits can be decomposed into two parts reflecting common for all traits (shared genetic background) and specific for particular trait background. To identify shared and unshared (or specific) background of CBP, SHAHER composes two new traits: SGIT, which condenses the shared genetic background, and UGIT, which is controlled by specific for CBP genetic

background. Both SGIT and UGIT are considered as linear combinations of the original traits with
specific coefficients. The coefficients for SGIT are calculated by maximizing the proportion of
CBP shared genetic background in SGIT genetic background. To calculate the coefficients for
UGIT, we estimated the residual genetic background of CBP after adjustment for SGIT. Using
these coefficients and the GWAS summary statistics for all original traits, the GWAS summary
statistics for SGIT and UGIT of CBP were calculated.

To estimate the SGIT coefficients we utilized the heritability estimates, phenotypic correlation matrices, and genetic correlation matrices (see details in Supplementary Methods) assessed beforehand for the original pain traits in the discovery sample. Further, we used the coefficients obtained using discovery sample in SHAHER analysis of all replication samples (Figure 1).

#### 585 Meta-analysis of GWAS summary statistics using METAL

For both SGIT and UGIT we performed two GWAS meta-analyses: (i) a meta-analysis comprising 586 all replication samples (Replication meta-analysis, N<sub>total</sub> = 191,580); and (ii) a meta-analysis of 587 588 two European samples (European meta-analysis,  $N_{total} = 439,831$ ). The first meta-analysis served 589 as a replication set for the GWAS associations observed in the discovery set, while the second one 590 was utilized for further functional annotation. In a gene-based analysis we utilized European metaanalysis for additional replication of the signals from the discovery sample. All the meta-analyses 591 were conducted using an inverse-variance-weighted approach and fixed-effects model 592 593 implemented in METAL software, version 2011-03-25 [60].

## 594 Gene identification

#### 595 Loci identification and replication

596 In order to identify loci statistically significantly associated with either SGIT or UGIT we carried out the conditional and joint analysis using GCTA-COJO software (version 1.90.0beta [61]) on 597 598 the discovery sample. We chose the following settings: minor allele frequency (MAF) not less than 599 2e-04; stepwise model selection procedure to select independently associated SNPs; significance threshold (applied to statistics after correction for genomic control using the LD Score regression 600 intercept) p-value = 5e-08/6 = 8.3e-09, where the denominator corresponds to the number of 601 602 analyzed traits (SGIT and UGIT from the current work, as well as four genetically independent phenotypes [GIP1-4] from our recent study [16] which were used for comparison); and  $\pm 5,000$  kb 603 604 window for linkage disequilibrium assessment. A linkage disequilibrium matrix was computed using PLINK version 1.9 (https://www.cog-genomics.org/plink/1.9/) using 100,000 randomly 605 selected individuals from the discovery set. Loci associated with either SGIT or UGIT were 606 607 defined as genomic regions of  $\pm 250$  kb from the lead SNPs identified in COJO analysis.

For both SGIT-associated and UGIT-associated loci we performed a two-step replication 608 procedure. First, we checked whether associations were observed in previous pain studies. We 609 compared SGIT-associated regions with those reported for the first genetically independent 610 phenotype (GIP1) [16] and collated UGIT-associated signal with results from an earlier back pain 611 612 GWAS [15]. Loci demonstrating the same direction of the effect in the current and previous studies were considered as known. Second, we focused on new signals (not reported previously) and 613 examined them in the replication meta-analysis. We assumed new signals to be *replicated* if two 614 615 criteria were met: 1) association in the replication meta-analysis was statistically significant after Bonferroni correction [62] for multiple testing (p-value < 0.05/n, where *n* is a number of new loci 616

associated with either SGIT or UGIT); and 2) direction of the effect was the same in the discoverysample and Replication meta-analysis.

#### 619 Gene prioritization in associated loci

We prioritized genes in the replicated locus using a protocol we have described previously in our 620 work on GIP [16]. For gene prioritization we applied a series of methods including (1) literature-621 based annotation of all genes within the locus using various databases; (2) prediction of SNP 622 effects in the locus utilizing Ensembl Variant Effect Predictor (VEP) [31], FATHMM-XF [32] and 623 624 FATHMM-INDEL [33] tools; (3) gene prioritization embedded into Data-driven Expression 625 Prioritized Integration for Complex Traits (DEPICT) software [34], and Functional Mapping and Annotation tool (FUMA) [35]; (4) analysis of pleiotropic effects of the replicated loci on gene 626 expression in various tissue types (see Supplementary Table ST16) using an instrument combining 627 Mendelian randomization (Summary data-based Mendelian Randomization, SMR) with 628 heterogeneity testing (Heterogeneity in Dependent Instruments, HEIDI) [36]. Genes provided 629 630 more evidence for prioritization were considered as more likely to be causal. Details are available in the Supplementary Methods. 631

#### 632 *Gene-based association analysis*

We conducted a gene-based association analysis of SGIT and UGIT utilizing GWAS summary statistics (z-scores and effect sizes) from three datasets: the discovery sample (N = 265,000), the European replication sample (N = 174,831), and data from the European meta-analysis (N<sub>total</sub> = 439,831). We preliminarily calculated the matrices of correlations between genotypes of all SNPs within a gene using individual genotypes of non-relatives from the entire European sample of UK Biobank participants (N = 315,599) and the PLINK software v2.00a3.7LM (<u>https://www.cog-</u> genomics.org/plink/2.0/) with options –maf 5e-05 --geno 0.02. For replication, we restricted to the

European sample because discovery and replication gene-based analyses had to be performed 640 using the same matrix of genotype correlations, which is ancestry-specific. In the discovery sample 641 we set the Bonferroni adjusted significance level [62] for the total number of genes (20,000) at 642 2.5e-06. We defined two criteria for replication of a gene-based signal: 1) the p-value in European 643 replication sample had to be less than 0.05/k, where k corresponds to the number of statistically 644 645 significant signals in the discovery sample; and 2) the p-value in European meta-analysis had to be less then the p-value in the discovery sample. See Supplementary Methods for more details on 646 genomic regions studied, methods of gene-based and conditional analyses. 647

#### 648 Analysis of CBP genetic architecture

#### 649 Gene set and tissue/cell type enrichment

We conducted DEPICT tissue or cell type enrichment analyses for GWAS loci associated with SGIT or UGIT in European meta-analysis. We used DEPICT software, version 1.1 rel194 [34], with default settings. For each of the traits we analyzed two SNP lists: a set of variants associated with the trait under p-value < 2.5e-08 significance threshold, and a set of SNPs with p-value < 5e-06. The significance thresholds corresponds to those recommended by program developers with correction for multiple testing for two traits applied. Variants were selected using GCTA-COJO instrument with settings described above in 'Loci identification and replication' section.

For all genes identified for SGIT we ran the gene set enrichment and tissue specificity analyses utilizing FUMA function GENE2FUNC. As the input data we used: 1) all genes prioritized for SGIT in the GWAS loci (both from the new replicated locus and from the loci previously replicated for GIP1); 2) all genes found in gene-based analysis and replicated in European meta-analysis. We analyzed SGIT-associated genes only, because there were not enough genes identified for UGIT to perform the analysis. The GTEx v8 dataset for 30 general tissues was used for tissue specificity

analysis. A set of the input genes was tested against each of the sets of differentially expressed genes (DEG) using a hypergeometric test and Bonferroni multiple testing correction. Statistical significance was determined at an adjusted p-value < 0.05. We estimated the overrepresentation of the identified genes in the gene sets of the GWAS catalog. Hypergeometric tests were performed to check if the genes of interest were overrepresented in any of the pre-defined gene sets. In this type of analysis, we used FDR (Benjamini – Hochberg method) for multiple testing correction as was recommended by FUMA. Statistical significance was determined at a q-value < 0.05.

#### 670 Genetic correlation between SGIT, UGIT and complex traits

We assessed genetic correlations of SGIT and UGIT with 730 complex human traits from the 671 GWAS-Map database using LD Score regression tool embedded in the platform. The detailed 672 673 protocol of selection of 730 complex traits included in the analysis and the full list of them were provided in a recent study by Timmers et al. [63]. For further interpretation we restricted our 674 analysis to 322 complex traits which had statistically significant genetic correlations with SGIT 675 676 and/or UGIT with magnitude greater than 0.25 (p-value < 3.42e-5 = 0.05/(730\*2), where 730 represents the total number of complex traits and 2 refers to the number of traits). To simplify the 677 678 perception and visualization of the results of genetic correlation analysis we performed a hierarchal clustering of 322 complex traits. This was done utilizing the ward.D2 method from the standard 679 hclust() R function on a squared matrix of genetic correlations transformed to the Euclidian 680 681 distances (as.dist() R function) as described previously by Tsepilov et al. [16]. We set the arbitrary threshold of h = 1.8 to cut the dendrogram of hierarchical clustering by height and obtained 11 682 683 clusters, each of which was then annotated manually. The full list of 322 traits grouped by clusters 684 is provided in Supplementary Table ST9.

Heatmap visualization of the genetic correlations was made using the heatmap.2() function from the gplots R package, version 3.1.1. From each cluster, we depicted one trait having the lowest pvalue of genetic correlation among all pairwise correlations with SGIT and UGIT within the cluster. The threshold for statistical significance was p-value > 9.41e-07 = 0.05/(51,681 + 730\*2), where 51,681 is the number of unique genetic correlation coefficients in the squared matrix for 322 complex traits, and 730\*2 is the total number of genetic correlations between SGIT, UGIT and 730 complex traits.

## 692 Polygenic risk scores of SGIT and UGIT. Their role in disease/medical intervention prediction

In addition to other functional analyses, we examined the prognostic power of SGIT and UGIT 693 694 polygenic risk scores (PRS) and assessed their role in disease/medical intervention prediction. We 695 calculated the PRSs for SGIT and UGIT utilizing data from the entire European sample (discovery and European replication samples,  $N_{total} = 439,762$ ) using a three-step algorithm. We used the 696 SBayesR method [64] and the protocol from our recent work on CBP PRS [24], with some 697 698 modifications. We divided the entire European sample into a training set (discovery sample, N =699 265,000), a validation set (subsample from the European replication sample, N = 30,000), and a 700 test set (the other 144,831 Northern Europeans from the European replication sample). The first set was used to develop PRS models, the second one was used to validate these models in 701 individual-level data and select the optimal models for SGIT and UGIT, and the last dataset was 702 703 used to assess quality metrics of the optimal models. At the final step, PRS values were calculated 704 for all individuals from entire European sample using optimal models. More details on PRS calculation available in Supplementary Methods. 705

When the PRS for SGIT and UGIT were estimated for the entire European sample of 439,831, we
normalized the PRS (the variance was taken to 1) and used them in a number of generalized linear

708 models (GLMs) considering PRS as a predictor for a particular trait. List of the traits and corresponding phenotypic data were taken from medical histories and questionnaires obtained 709 from UK Biobank participants of European descent from the test dataset mentioned in the 710 paragraph above, using non-relatives only (N = 120,200). Non-relatives were defined by UK 711 712 Biobank data-field 22021 details (for more see https://biobank.ndph.ox.ac.uk/showcase/coding.cgi?id=682). 713

714 Medical codes were combined to the second level, meaning that participants with a specific medical code of the second level and/or codes of a lower level (the third and/or the fourth) nested 715 under the second level code were all considered to be cases of the second level medical code (for 716 717 more information on ICD10 and OPCS4 coding in the UK Biobank see https://biobank.ndph.ox.ac.uk/ukb/field.cgi?id=41202 718 and

719 <u>https://biobank.ndph.ox.ac.uk/ukb/field.cgi?id=41200</u>). Only ICD10 codes from chapters I-XVII

720 were analyzed, OPCS4 codes from chapters X, Y, Z were excluded, and then the remaining codes were filtered by prevalence (> 0.5% and < 99.5%). The final list of the traits (see Supplementary 721 Table ST10) comprised the six self-reported chronic pain phenotypes (questionnaire-based) 722 723 considered in the study, 165 ICD10 and 132 OPCS4 medical codes served as proxies for the 724 corresponding diseases/medical interventions. To perform GLM-analyses we used a logistic 725 regression from a standard glm() function in R and included sex, age, batch number and the first ten genetic principal components (PC1-10) provided by UK Biobank as covariates in addition to 726 727 the PRS predictor. The general formula is the following: medical code  $\sim$  age + sex + batch + PC1 728  $+ \dots + PC10 + PRS$ . Finally, we filtered out the GLM results not passing the significance threshold 729 of p-value < 8.25e-05 = 0.05/(2\*(165 + 132 + 6)), where 165 and 132 are the numbers of ICD10 730 and OPCS4 codes, respectively, 6 is a number of pain traits and 2 corresponds to the number of

731 PRSs. For further heatmap visualization of the significant GLM-results we performed a hierarchical clustering of medical codes based on their phenotypic correlation matrix as described 732 above in the genetic correlation analysis. We set a h = 1.25 threshold for cutting the hierarchical 733 734 clustering dendrogram by height and defined 19 clusters (see Supplementary Table ST10 for more details). Similarly, for the genetic correlation analysis we selected one medical code from each 735 736 cluster, which provided the lowest p-value of association with SGIT and UGIT polygenic risk scores. Alongside medical codes, representing particular clusters, we added six pain traits to the 737 heatmap plot. The visualization was made as described in the genetic correlation analysis section. 738

## 739 CBP subphenotyping among people with CBP

To investigate the potential of SGIT and UGIT PRSs for back pain subphenotyping, we analyzed 740 741 only CBP cases from the entire European sample, focusing on non-relatives only (N = 65,011). 742 The decision to switch to the entire European sample instead of working with the test sample only 743 was motivated by the limited statistical power of the test sample. We normalized SGIT and UGIT 744 PRS values and recoded each of the vectors of normalized PRSs as binary traits in two different 745 ways. First, we coded the normalized PRS values as 1 if a participant with back pain was attributed 746 to the lowest decile of the normalized PRS vector for SGIT or UGIT, respectively. Normalized 747 PRS values falling out of this range were coded as 0. By such a coding ("yes/no in the lowest decile of the normalized PRS values for SGIT"; "yes/no in the lowest decile of the normalized 748 749 PRS values for UGIT") we highlighted individuals with the lowest genetic predisposition to CBP 750 through SGIT and UGIT. Then, we coded as 1 the normalized PRS estimates referring to the 751 highest decile of the SGIT or UGIT PRS distribution, respectively, and coded them as 0 otherwise. 752 In this binary coding ("yes/no in the highest decile of the normalized PRS values for SGIT"; 753 "yes/no in the highest decile of the normalized PRS values for UGIT") we marked individuals with

| 754 | the highest genetic predisposition to CBP through SGIT and UGIT. For each of these binary traits         |
|-----|--|
| 755 | we performed GLM analyses using the model described above. As a dependent variable, we                   |
| 756 | considered ICD10 and OPCS4 codes described above. In this sample of 65,011 non-related                   |
| 757 | participants with CBP we selected 199 ICD10 and 165 OPCS4 codes with prevalence $> 0.5\%$ and            |
| 758 | < 99.5%. The significance threshold for this analysis was set at $p = 3.43e-05 = 0.05/(4*(199 + 10.05))$ |
| 759 | 165)), where 4 is the number of binary PRS traits, 199 is the number of ICD10 codes, and 165 is          |
| 760 | the number of OPCS4 codes.   |
|     |  |

761

## 763 Data and code availability

Full project code is available at <u>https://github.com/ElizavetaElgaeva/BP-SH project\_code</u>.
GWAS and EWAS summary statistics for all pain traits along with SGIT and UGIT PRS model
weights can be found here (the link will be added later).

## 767 Acknowledgements

- This study was conducted using the UK Biobank Resource (project #18219 and #59345). The work
- of ANT is an output of a research project implemented as part of the Research Program at the MSU
- 770 Institute for Artificial Intelligence. FMKW was supported by Versus Arthritis, grant #22467. We
- thank Dr. Alexandra S. Shadrina for providing template for Figure 1 and Dr. Sodbo Z. Sharapovfor advices on study design.

## 773 Funding

IVZ, YAT, GRS, and AVK were supported by the budget project of the Institute of Cytology and
Genetics FWNR-2022-0020. The work of TIA and EEE was supported by the Russian Science
Foundation (grant #22-15-20037) and Government of the Novosibirsk region.

## 777 Authors' contributions

EEE participated in conceptualization, data curation, formal analysis, investigation, project administration, software, supervision, validation, visualization, and original draft preparation. IVZ contributed to formal analysis, investigation, software, validation, visualization, and original draft preparation. AVN took part in data curation, formal analysis, software, methodology, validation, visualization, and original draft preparation. DAV carried out formal analysis, software and

783 investigation. EST participated in data curation, formal analysis and software. ANT contributed in formal analysis and software. AVK provided data curation, formal analysis, resources and 784 software. GRS contributed in methodology and software. MBF, FMKW, PS and YSA carried out 785 786 data curation. TIA, YAT, YSA, FMKW, and PS participated in conceptualization, investigation, 787 project administration, and supervision. TIA contributed to visualization. YAT provided funding acquisition. All co-authors discussed the results and contributed to review and editing while 788 preparing the final version of the manuscript. All authors have read and agreed to the published 789 version of the manuscript. 790

# 791 **Conflict of interest**

YSA is a cofounder and a co-owner of PolyOmica and PolyKnomics, private organizations
providing services, research, and development in the field of computational and statistical
genomics.

## 796 Supplementary note

### 797 Supplementary Methods

798 Extended description of the Materials and Methods section.

#### 799 Supplementary Results

800 Detailed description of the results of gene prioritization.

## 801 Supplementary Figure 1

802 Manhattan plot for SGIT in discovery sample.

## 803 Supplementary Figure 2

804 Quantile-quantile (QQ) plot for SGIT in discovery sample.

## 805 Supplementary Figure 3

806 Manhattan plot for UGIT in discovery sample.

## 807 Supplementary Figure 4

808 Quantile-quantile (QQ) plot for UGIT in discovery sample.

#### 809 Supplementary Figure 5

810 Gene expression patterns of 23 genes identified for SGIT.

## 811 Supplementary Figure 6

- 812 Regional association plot for rs2587363 locus associated with SGIT under suggestive significance
- 813 threshold p-value < 8.3e-08.

## 814 Supplementary Table ST1

- 815 Description of the sample.
- 816 Supplementary Table ST2
- 817 Characteristics of pain traits.

## 818 Supplementary Table ST3

- Top loci associated with SGIT and UGIT at a study-level threshold of statistical significance (p-
- 820 value < 8.3e-09).

## 821 Supplementary Table ST4

822 Results of SNP effect prediction using Variant Effect Predictor and FATMM tools.

## 823 Supplementary Table ST5

- Results of gene prioritization using DEPICT for SNPs associated with SGIT in European meta-
- 825 analysis.

## 826 Supplementary Table ST6

- 827 Results of FUMA gene prioritization for SNPs associated with SGIT in European meta-analysis
- 828 at p-value < 2.5e-08 and p-value < 5e-06.
- 829 Supplementary Table ST7
- 830 Results of gene-based and conditional analyses.
- 831 Supplementary Table ST8
- 832 Results of DEPICT analysis for SNPs associated with SGIT in European meta-analysis.

## 833 Supplementary Table ST9

834 Analysis of genetic correlations.

## 835 Supplementary Table ST10

836 Results of PRS analysis in the test sample, non-relatives only.

## 837 Supplementary Table ST11

838 Results of binary-coded PRS traits analysis in participants with back pain.

#### 839 Supplementary Table ST12

Top loci associated with SGIT and UGIT under suggestive significance threshold p-value < 8.3e-

**841** 08.

## 842 Supplementary Table ST13

843 Results of SNP effect prediction using FATHMM-XF for the locus tagged by rs2587363.

## 844 Supplementary Table ST14

Results of SMR-HEIDI analysis of pleiotropic effects on gene expression for the locus tagged by

846 rs2587363.

## 847 Supplementary Table ST15

848 Results of gene-based association analysis in exome data.

## 849 Supplementary Table ST16

List of tissues used for SMR-HEIDI analysis of pleiotropic effects on gene expression.

# 852 **References**

| 853 | 1. | Manchikanti L, Singh V, Falco FJE, Benyamin RM, Hirsch JA. Epidemiology of Low            |
|-----|----|---|
| 854 |    | Back Pain in Adults. Neuromodulation Technol Neural Interface. 2014;17: 3–10.             |
| 855 |    | doi:https://doi.org/10.1111/ner.12018   |
| 856 | 2. | Vos T, Allen C, Arora M, Barber RM, Bhutta ZA, Brown A, et al. Global, regional, and      |
| 857 |    | national incidence, prevalence, and years lived with disability for 310 diseases and      |
| 858 |    | injuries, 1990–2015: a systematic analysis for the Global Burden of Disease Study 2015.   |
| 859 |    | Lancet. 2016;388: 1545–1602. doi:10.1016/S0140-6736(16)31678-6                            |
| 860 | 3. | Foster NE, Anema JR, Cherkin D, Chou R, Cohen SP, Gross DP, et al. Prevention and         |
| 861 |    | treatment of low back pain: evidence, challenges, and promising directions. Lancet.       |
| 862 |    | 2018;391: 2368–2383. doi:10.1016/S0140-6736(18)30489-6                                    |
| 863 | 4. | Reddy K, Sinha P, O'Kane CM, Gordon AC, Calfee CS, McAuley DF. Subphenotypes in           |
| 864 |    | critical care: translation into clinical practice. The Lancet Respiratory Medicine. 2020. |
| 865 |    | doi:10.1016/S2213-2600(20)30124-7   |
| 866 | 5. | Phillips CJ. Economic burden of chronic pain. Expert Rev Pharmacoecon Outcomes Res.       |
| 867 |    | 2006;6: 591–601. doi:10.1586/14737167.6.5.591   |
| 868 | 6. | Loeser JD. Chapter 2 Pain as a disease. Handb Clin Neurol. 2006. doi:10.1016/S0072-       |
| 869 |    | 9752(06)80006-0   |
| 870 | 7. | Gatchel RJ. The biopsychosocial model of chronic pain. Futur Med. 2013; 5–17.             |
| 871 |    | doi:10.2217/EBO.13.469  |
| 872 | 8. | Scholz J. Mechanisms of chronic pain. Mol Pain. 2014; (Suppl 1):O15. doi:10.1186/1744-    |

- 873 8069-10-s1-o15
- 874 9. Crofford LJ. Chronic Pain: Where the Body Meets the Brain. Trans Am Clin Climatol
  875 Assoc. 2015;126: 167–183.
- Battie MC, Videman T, Levalahti E, Gill K, Kaprio J. Heritability of Low Back Pain and
  the Role of Disc Degeneration 53rd Annual Meeting of the Orthopaedic Research Society
  Paper No : 0317. Pain. 2007;131: 272–280.
- 11. Junqueira DRG, Ferreira ML, Refshauge K, Maher CG, Hopper JL, Hancock M, et al.
- 880 Heritability and lifestyle factors in chronic low back pain: Results of the Australian Twin
- Low Back Pain Study (The AUTBACK study). Eur J Pain (United Kingdom). 2014.
  doi:10.1002/ejp.506
- 883 12. Suri P, Palmer MR, Tsepilov YA, Freidin MB, Boer CG, Yau MS, et al. Genome-wide
- meta-analysis of 158,000 individuals of European ancestry identifies three loci associated
  with chronic back pain. PLoS Genet. 2018. doi:10.1371/journal.pgen.1007601
- 13. Johnston KJA, Adams MJ, Nicholl BI, Ward J, Strawbridge RJ, Ferguson A, et al.
- 887 Genome-wide association study of multisite chronic pain in UK biobank. PLoS Genet.
- 888 2019. doi:10.1371/journal.pgen.1008164
- 889 14. Bjornsdottir G, Stefansdottir L, Thorleifsson G, Sulem P, Norland K, Ferkingstad E, et al.
- 890 Rare SLC13A1 variants associate with intervertebral disc disorder highlighting role of
- sulfate in disc pathology. Nat Commun. 2022;13: 634. doi:10.1038/s41467-022-28167-1
- 15. Freidin MB, Tsepilov YA, Palmer M, Karssen LC, Suri P, Aulchenko YS, et al. Insight
- into the genetic architecture of back pain and its risk factors from a study of 509,000
- individuals. Pain. 2019. doi:10.1097/j.pain.00000000001514

- 16. Tsepilov YA, Freidin MB, Shadrina AS, Sharapov SZ, Elgaeva EE, Zundert J van, et al.
- Analysis of genetically independent phenotypes identifies shared genetic factors
- associated with chronic musculoskeletal pain conditions. Commun Biol. 2020;3.
- 898 doi:10.1038/s42003-020-1051-9
- 17. Li S, Brimmers A, Van Boekel RLM, Vissers KCP, Coenen MJH. Systematic Review and
- 900 Meta-Analysis A systematic review of genome-wide association studies for pain,
- nociception, neuropathy, and pain treatment responses. 2023 [cited 12 Feb 2024].
- 902 doi:10.1097/j.pain.000000000002910
- 903 18. Vehof J, Zavos HMS, Lachance G, Hammond CJ, Williams FMK. Shared genetic factors
  904 underlie chronic pain syndromes. Pain. 2014. doi:10.1016/j.pain.2014.05.002
- 905 19. Williams FMK, Spector TD, MacGregor AJ. Pain reporting at different body sites is

explained by a single underlying genetic factor. Rheumatology. 2010.

- 907 doi:10.1093/rheumatology/keq170
- 20. Jensen RK, Jensen TS, Kjaer P, Kent P. Can pathoanatomical pathways of degeneration in
- lumbar motion segments be identified by clustering MRI findings. BMC Musculoskelet

910 Disord. 2013;14. doi:10.1186/1471-2474-14-198

- 21. Ranoux D, Attal N, Morain F, Bouhassira D. Botulinum toxin type A induces direct
- analgesic effects in chronic neuropathic pain. Ann Neurol. 2008;64.
- 913 doi:10.1002/ana.21427
- 22. Tagliaferri SD, Angelova M, Zhao X, Owen PJ, Miller CT, Wilkin T, et al. Artificial
- 915 intelligence to improve back pain outcomes and lessons learnt from clinical classification
- approaches: three systematic reviews. npj Digital Medicine. 2020. doi:10.1038/s41746-

- 917 020-0303-x
- 918 23. Shendure J, Findlay GM, Snyder MW. Genomic Medicine–Progress, Pitfalls, and
  919 Promise. Cell. 2019;177: 45–57. doi:10.1016/j.cell.2019.02.003
- 920 24. Tsepilov YA, Elgaeva EE, Nostaeva A V., Compte R, Kuznetsov IA, Karssen LC, et al.
- 921 Development and Replication of a Genome-Wide Polygenic Risk Score for Chronic Back
- Pain. J Pers Med. 2023;13. doi:10.3390/jpm13060977
- 25. Svishcheva G, Tiys E, Elgaeva E, Feoktistova S, Timmers P, Sharapov S, et al. A Novel
- 924 Framework for Analysis of the Shared Genetic Background of Correlated Traits. Genes
- 925 (Basel). 2022;13: 1694. doi:10.3390/genes13101694
- 26. Lee S, Wu MC, Lin X. Optimal tests for rare variant effects in sequencing association
  studies. Biostatistics. 2012;13: 762–775. doi:10.1093/biostatistics/kxs014
- 27. Wang K, Abbott D. A principal components regression approach to multilocus genetic
- association studies. Genet Epidemiol. 2008;32: 108–118. doi:10.1002/gepi.20266
- 930 28. Liu Y, Chen S, Li Z, Morrison AC, Boerwinkle E, Lin X. ACAT: A Fast and Powerful p
- 932 Genet. 2019;104: 410–421. doi:10.1016/j.ajhg.2019.01.002
- 29. Shashkova TI, Gorev DD, Pakhomov ED, Shadrina AS, Sharapov SZ, Tsepilov YA, et al.
- 934 The GWAS-MAP platform for aggregation of results of genome-wide association studies
- and the GWAS-MAP|homo database of 70 billion genetic associations of human traits.
- 936 Vavilovskii Zhurnal Genet Selektsii. 2021;24: 876–884. doi:10.18699/VJ20.686
- 937 30. Shashkova TI, Pakhomov ED, Gorev DD, Karssen LC, Joshi PK, Aulchenko YS.
- 938 PheLiGe: an interactive database of billions of human genotype-phenotype associations.

| 333 $1000000000000000000000000000000000000$ | 939 | Nucleic Acids Res. | 2021;49: 1347-1350. | doi:10.1093/nar/gkaa1086 |
|---|-----|--------------------|---------------------|--------------------------|
|---|-----|--------------------|---------------------|--------------------------|

- 940 31. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GRS, Thormann A, et al. The Ensembl
- 941 Variant Effect Predictor. Genome Biol. 2016. doi:10.1186/s13059-016-0974-4
- 32. Rogers MF, Shihab HA, Mort M, Cooper DN, Gaunt TR, Campbell C. FATHMM-XF:
- Accurate prediction of pathogenic point mutations via extended features. Bioinformatics.
- 944 2018. doi:10.1093/bioinformatics/btx536
- 33. Ferlaino M, Rogers MF, Shihab HA, Mort M, Cooper DN, Gaunt TR, et al. An integrative
- approach to predicting the functional effects of small indels in non-coding regions of the
- 947 human genome. BMC Bioinformatics. 2017. doi:10.1186/s12859-017-1862-y
- 948 34. Pers TH, Karjalainen JM, Chan Y, Westra HJ, Wood AR, Yang J, et al. Biological
- 949 interpretation of genome-wide association studies using predicted gene functions. Nat
- 950 Commun. 2015. doi:10.1038/ncomms6890
- 951 35. Watanabe K, Taskesen E, van Bochoven A, Posthuma D. Functional mapping and
- annotation of genetic associations with FUMA. Nat Commun. 2017;8: 1826.
- 953 doi:10.1038/s41467-017-01261-5
- 36. Zhu Z, Zhang F, Hu H, Bakshi A, Robinson MR, Powell JE, et al. Integration of summary
- data from GWAS and eQTL studies predicts complex trait gene targets. Nat Genet. 2016.
- 956 doi:10.1038/ng.3538
- 37. Johnston KJA, Ward J, Ray PR, Adams MJ, McIntosh AM, Smith BH, et al. Sex-stratified
  genome-wide association study of multisite chronic pain in UK Biobank. PLoS Genet.
- 959 2021;17. doi:10.1371/journal.pgen.1009428
- 38. Zhang F, Rao S, Baranova A. Shared genetic liability between major depressive disorder

| 961 |     | and osteoarthritis. Bone Jt Res. 2022;11. doi:10.1302/2046-3758.111.BJR-2021-0277.R1    |
|-----|-----|---|
| 962 | 39. | Zhang F, Rao S, Cao H, Zhang X, Wang Q, Xu Y, et al. Genetic evidence suggests          |
| 963 |     | posttraumatic stress disorder as a subtype of major depressive disorder. J Clin Invest. |
| 964 |     | 2022;132. doi:10.1172/JCI145942   |
| 965 | 40. | Das S, Taylor K, Kozubek J, Sardell J, Gardner S. Genetic Risk Factors for ME/CFS       |
| 966 |     | Identified using Combinatorial Analysis. medRxiv. 2022; 2022.09.09.22279773.            |
| 967 |     | doi:10.1101/2022.09.09.22279773   |
| 968 | 41. | Tachmazidou I, Hatzikotoulas K, Southam L, Esparza-Gordillo J, Haberland V, Zheng J,    |
| 969 |     | et al. Identification of new therapeutic targets for osteoarthritis through genome-wide |
| 970 |     | analyses of UK Biobank data. Nature Genetics. 2019. doi:10.1038/s41588-018-0327-1       |
| 971 | 42. | Haller G, McCall K, Jenkitkasemwong S, Sadler B, Antunes L, Nikolov M, et al. A         |
| 972 |     | missense variant in SLC39A8 is associated with severe idiopathic scoliosis. Nat Commun. |
| 973 |     | 2018. doi:10.1038/s41467-018-06705-0  |
| 974 | 43. | Park JH, Hogrebe M, Grüneberg M, Duchesne I, Von Der Heiden AL, Reunert J, et al.       |
| 975 |     | SLC39A8 Deficiency: A Disorder of Manganese Transport and Glycosylation. Am J Hum       |
| 976 |     | Genet. 2015;97: 894–903. doi:10.1016/j.ajhg.2015.11.003                                 |
| 977 | 44. | Bonaventura E, Barone R, Sturiale L, Pasquariello R, Alessandrì MG, Pinto AM, et al.    |
| 978 |     | Clinical, molecular and glycophenotype insights in SLC39A8-CDG. Orphanet J Rare Dis.    |
| 979 |     | 2021;16: 307. doi:10.1186/s13023-021-01941-y  |
| 980 | 45. | Hidiroglou M, Ivan M, Bryan MK, Ribble CS, Janzen ED, Proulx JG, et al. Assessment of   |
| 981 |     | the role of manganese in congenital joint laxity and dwarfism in calves. Ann Rech Vet.  |
| 982 |     | 1990;21.  |
|     |     |   |

- 983 46. Bortsov A V., Parisien M, Khoury S, Martinsen AE, Lie MU, Heuch I, et al. Brain-
- 984 specific genes contribute to chronic but not to acute back pain. Pain Reports. 2022;7.
- 985 doi:10.1097/PR9.000000000001018
- 47. Barki M, Xue H. GABRB2, a key player in neuropsychiatric disorders and beyond. Gene.
- 987 2022;809. doi:10.1016/j.gene.2021.146021
- 988 48. Feng C, Zhao J, Ji F, Su L, Chen Y, Jiao J. TCF 20 dysfunction leads to cortical
- neurogenesis defects and autistic-like behaviors in mice . EMBO Rep. 2020;21.
- 990 doi:10.15252/embr.201949239
- 49. Ao X, Parisien M, Zidan M, Grant A V, Martinsen AE, Winsvold BS, et al. Rare variant
- analyses in large-scale cohorts identified SLC13A1 associated with chronic pain. 2023

993 [cited 17 Aug 2023]. doi:10.1097/j.pain.00000000002882

- 50. Langford R, Hurrion E, Dawson PA. Genetics and pathophysiology of mammalian sulfate
  biology. Journal of Genetics and Genomics. 2017. doi:10.1016/j.jgg.2016.08.001
- 996 51. Song YQ, Karasugi T, Cheung KMC, Chiba K, Ho DWH, Miyake A, et al. Lumbar disc
- degeneration is linked to a carbohydrate sulfotransferase 3 variant. J Clin Invest. 2013.
- 998 doi:10.1172/JCI69277
- 999 52. Mountjoy E, Schmidt EM, Carmona M, Schwartzentruber J, Peat G, Miranda A, et al. An
- 1000 open approach to systematically prioritize causal variants and genes at all published
- 1001 human GWAS trait-associated loci. Nat Genet. 2021;53. doi:10.1038/s41588-021-00945-5
- 1002 53. Chebib J, Guillaume F. Pleiotropy or linkage? Their relative contributions to the genetic
- 1003 correlation of quantitative traits and detection by multitrait GWA studies. Genetics.
- 1004 2021;219. doi:10.1093/GENETICS/IYAB159

| 1005 | 54. | Elgart M, Goodman MO, Isasi C, Chen H, Morrison AC, de Vries PS, et al. Correlations    |
|------|-----|---|
| 1006 |     | between complex human phenotypes vary by genetic background, gender, and                |
| 1007 |     | environment. Cell Reports Med. 2022;3. doi:10.1016/j.xcrm.2022.100844                   |
| 1008 | 55. | Elgaeva EE, Tsepilov Y, Freidin MB, Williams FMK, Aulchenko Y, Suri P. ISSLS Prize      |
| 1009 |     | in Clinical Science 2020. Examining causal effects of body mass index on back pain: a   |
| 1010 |     | Mendelian randomization study. Eur Spine J. 2020;29: 686–691. doi:10.1007/s00586-019-   |
| 1011 |     | 06224-6   |
| 1012 | 56. | Suri P, Elgaeva EE, Williams FMK, Freidin MB, Zaytseva OO, Aulchenko YS, et al.         |
| 1013 |     | Evidence of causal effects of blood pressure on back pain and back pain on type II      |
| 1014 |     | diabetes provided by a bidirectional Mendelian randomization study. Spine J. 2023.      |
| 1015 |     | doi:10.1016/j.spinee.2023.04.001  |
| 1016 | 57. | Elliott MD, Heitner JF, Kim H, Wu E, Parker MA, Lee DC, et al. Prevalence and           |
| 1017 |     | prognosis of unrecognized myocardial infarction in asymptomatic patients with diabetes: |
| 1018 |     | A two-center study with up to 5 years of follow-up. Diabetes Care. 2019;42: 1290–1296.  |
| 1019 |     | doi:10.2337/dc18-2266   |
| 1020 | 58. | Chiariello M, Indolfi C. Silent myocardial ischemia in patients with diabetes mellitus. |
| 1021 |     | Circulation. 1996. doi:10.1161/01.CIR.93.12.2089  |
| 1022 | 59. | Sudlow C, Gallacher J, Allen N, Beral V, Burton P, Danesh J, et al. UK Biobank: An      |
| 1023 |     | Open Access Resource for Identifying the Causes of a Wide Range of Complex Diseases     |
| 1024 |     | of Middle and Old Age. PLOS Med. 2015. doi:10.1371/journal.pmed.1001779                 |
| 1025 | 60. | Willer CJ, Li Y, Abecasis GR. METAL: Fast and efficient meta-analysis of genomewide     |
| 1026 |     | association scans. Bioinformatics. 2010. doi:10.1093/bioinformatics/btq340              |

- 1027 61. Yang J, Ferreira T, Morris AP, Medland SE, Madden PAF, Heath AC, et al. Conditional
- and joint multiple-SNP analysis of GWAS summary statistics identifies additional variants
- influencing complex traits. Nat Genet. 2012. doi:10.1038/ng.2213
- 1030 62. Curtin F, Schulz P. Multiple correlations and Bonferroni's correction. Biological
- 1031 Psychiatry. 1998. doi:10.1016/S0006-3223(98)00043-2
- 1032 63. H J Timmers PR, Tiys ES, Sakaue S, Akiyama M, J Kiiskinen TT, Zhou W, et al.
- 1033 Mendelian randomization of genetically independent aging phenotypes identifies LPA and
- 1034 VCAM1 as biological targets for human aging. Nat AGiNG |. 2022;2.
- 1035 doi:10.1038/s43587-021-00159-8
- 1036 64. Lloyd-Jones LR, Zeng J, Sidorenko J, Yengo L, Moser G, Kemper KE, et al. Improved
- 1037 polygenic prediction by Bayesian multiple regression on summary statistics.
- 1038 doi:10.1038/s41467-019-12653-0