

# Common genetic associations between age-related diseases

Handan Melike Dönertaş<sup>1\*</sup>, Daniel K. Fabian<sup>1</sup>, Matías Fuentealba Valenzuela<sup>1,2</sup>, Linda Partridge<sup>2,3</sup>, Janet M. Thornton<sup>1\*</sup>

<sup>1</sup>European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge, CB10 1SD, UK.

<sup>2</sup>Institute of Healthy Aging, Department of Genetics, Evolution and Environment, University College London, London, UK

<sup>3</sup>Max Planck Institute for Biology of Aging, Cologne, Germany

\*Correspondence: [donertas.melike@gmail.com](mailto:donertas.melike@gmail.com); [thornton@ebi.ac.uk](mailto:thornton@ebi.ac.uk)

## Abstract

Age is a common risk factor in many diseases, but the molecular basis for this relationship is elusive. In this study we identified 4 disease clusters from 116 diseases in the UK Biobank data, defined by their age-of-onset profiles, and found that diseases with the same onset profile are genetically more similar, suggesting a common etiology. This similarity was not explained by disease categories, co-occurrences or disease cause-effect relationships. Two of the four disease clusters had an increased risk of occurrence from age 20 and 40 years respectively. They both showed an association with known aging-related genes, yet differed in functional enrichment and evolutionary profiles. We tested mutation accumulation and antagonistic pleiotropy theories of aging and found support for both. We also identified drug candidates for repurposing to target multiple age-dependent diseases with the potential to improve healthspan and alleviate multimorbidity and polypharmacy in the elderly.

## Keywords:

Aging, age-related disease, GWAS, UK Biobank, mutation accumulation, antagonistic pleiotropy

## Introduction

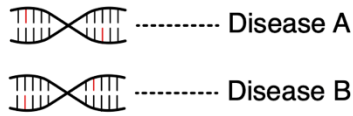
Aging is associated with a time-dependent decrease in the functional integrity of organisms and an increase in susceptibility to pathologies<sup>1</sup>. The worldwide increase in lifespan has not been matched by an increase in healthspan, and there is a growing period of loss of function, and disease at the end of life<sup>2</sup>. Aging thus poses a significant global challenge, because it is the major risk factor for chronic conditions, including cardiovascular disease, cancer, and dementia<sup>3</sup>. Although these diseases involve different organs and pathologies, they all show a strong dependence on age<sup>4</sup> and could, therefore share common etiologies based upon the underlying mechanisms of aging. It is therefore important to understand if the aging process itself leads to different age-related conditions through common pathways, or if the age-dependency of different diseases has independent, time-dependent causes.

Despite the negative impact of aging on organismal fitness and functionality, it is widespread in the animal world as well as in humans<sup>5</sup> and has therefore been described as an evolutionary paradox<sup>6</sup>. Aging can nonetheless evolve as the force of natural selection weakens with age due to extrinsic hazard. Mutations that are deleterious only in later life can accumulate in the population through mutation pressure, because the force of natural selection eliminating them from the population declines with the age of onset of their effects (mutation accumulation theory of aging)<sup>7</sup>. Pleiotropic variants that are beneficial during early life but detrimental later in life can also become prevalent in populations through natural selection (antagonistic pleiotropy theory of aging)<sup>8</sup>. Thus, genome-wide germline genetic variants that increase the risk of diseases at old age may not be pruned by natural selection or may be associated with beneficial phenotypes earlier in life.

The risk of many age-related diseases (ARDs) is influenced by genetic variation. Genome-wide association studies (GWAS) have identified genetic variants that alter complex traits. Pleiotropy, where variants or genes influence multiple traits, is more prevalent than previously thought<sup>9–12</sup>, indicating that different traits share common causal pathways<sup>13</sup>. Pleiotropy within the disease classification system<sup>12</sup> and in certain disease classes, such as immune-related diseases<sup>14,15</sup> and cancer<sup>16</sup>, has been studied, but the understanding of pleiotropy in ARDs more broadly is limited. Some studies have investigated the common pathways between manually curated age-related traits<sup>17–19</sup>. Despite the challenges of combining results from different published datasets, these studies provided the first clues that at least some ARDs share common pathways, which are also related to a significant but limited number of longevity-regulating genes in model organisms. In this study, using disease age-of-onset profiles, we extend the previous efforts by providing the first data-driven classification of a large number of diseases according to their age-profile, followed by a genetic analysis in one of the largest and most comprehensive cohorts available. In this way, we provide a comparative genetic analysis between ARDs and non-ARDs and also between ARDs with different age profiles.

The UK Biobank (UKBB)<sup>20</sup> includes genetic and health-related data for almost half a million participants. We extracted age-of-onset profiles for 116 diseases and identified unbiased clusters to define the relationship between disease incidence and age. We identified variants associated with each disease and compared the genetic associations between diseases based on these clusters. We first found that diseases with the same age profile share genetic associations, which cannot be explained by disease categories, co-occurrences, or mediated pleiotropy, and thus reflects a common etiology (Figure 1). We further characterized these shared associations compared to previously known longevity-associated genes and biological functions. We next performed a drug repurposing study to find drugs that could target multiple late-onset diseases simultaneously. Finally, we compared the variants associated with diseases that start to occur at different ages and identified different evolutionary characteristics, supporting the mutation accumulation and antagonistic pleiotropy theories of aging.

## I. Independent Genetic Associations



## II. Shared Genetic Associations

### a. Common Etiology



### b. Mediated Pleiotropy

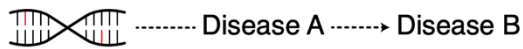


Figure 1: Summary of different models explaining the associations between diseases. Independent genetic associations (I) reflects the case where the number of shared genetic associations between diseases is not more than expected by chance. If the overlap is more than expected (II), it could either reflect (a) common etiology, which reflects shared causes, or (b) mediated pleiotropy, which suggests a common genetic factor influencing only one disease, which in turn increases the risk of a second disease.

## Results

### Data

We used self-reported diseases and age-at-diagnosis covering 484,598 participants, and their genotypes in the UKBB<sup>20</sup>. Details of the UKBB data, quality control steps, and exploratory analyses are given in the Supplementary Information and Supplementary Figures S1-8. Self-reported diseases in the UKBB are hierarchically structured and the top nodes; such as cardiovascular or endocrine diseases, were considered as *disease categories* (Figure S6). We only analyzed common diseases (*i.e.* with at least 2,000 cases) that were not sex-limited ( $n=116$  in 472 self-reported diseases). Importantly, we did not include cancer in our analyses as the interaction between genetic and environmental contributors is likely different from non-cancer illnesses, even though they may have a similar age-of-onset profile (for details, see Methods).

### Age-of-onset clusters

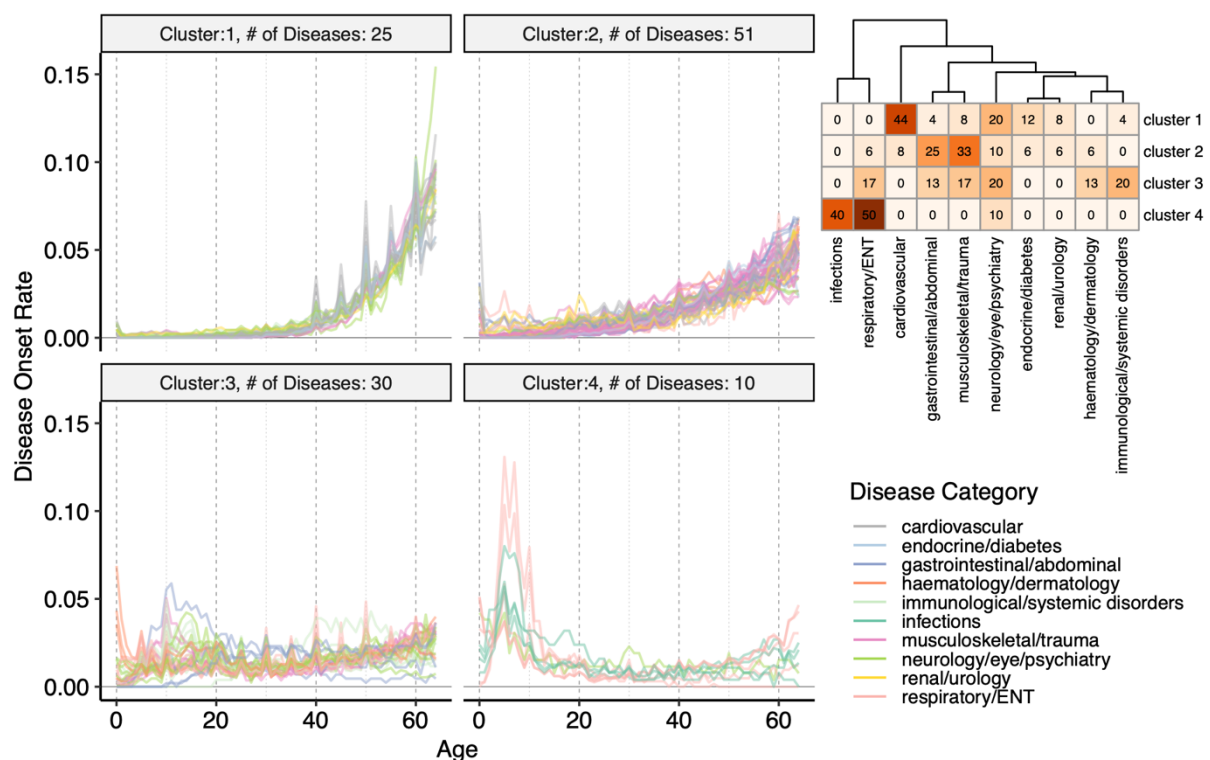


Figure 2: Age-of-onset profiles clustered by the PAM algorithm, using dissimilarities calculated with temporal correlation measure (CORT). The y-axis shows the number of individuals who were diagnosed with the disease at a certain age, divided by the total number of people having that disease. Values were calculated by taking the median value of 100 permutations of 10,000 people in the UKBB (see Methods). The x-axis shows the age-of-onset in years. Each line denotes one disease and is colored by disease categories. The heatmap in the right upper corner shows the percent overlap between categories and clusters. Numbers give the % of an age-of-onset cluster belonging to each category. Supplementary Figures S9-18 shows the distributions for each disease separately.

Age is associated with increased risk of many diseases. In order to characterize the association between age and different diseases we first used age-at-diagnosis as a proxy to disease onset and derived disease age-of-onset profiles (Figure S9-18). On average, cardiovascular and endocrine diseases had a high median age-of-onset, while infections had the lowest age-of-onset (Figure S19). We then clustered diseases into 4 clusters (the optimum number determined by the gap statistic) using the PAM algorithm and disease dissimilarities calculated using CORT distance<sup>21</sup> (Figure 2, Table S1). Cluster 1 diseases ( $n=25$ ) showed a rapid increase with age after the age of 40; 11 were cardiovascular diseases, but the cluster

also included other diseases such as diabetes, osteoporosis, and cataract. Cluster 2 (n=51) diseases started to increase in the population at an earlier age of 20, but had a slower rate of increase with age; the diseases in this cluster were the most diverse, including 17 musculoskeletal, 13 gastrointestinal diseases, as well as others such as anemia, deep venous thrombosis, thyroid problems, depression. Cluster 3 diseases (n=30) showed a low age dependency with a mostly uniform distribution across ages, but with slight increases around the ages of 10 and 60 years. This category included similar numbers of immunological, neurological, musculoskeletal, gastrointestinal and respiratory diseases but all have an 'immune' component even if not classified in this way by the UKBB (e.g., inflammatory bowel disease (gastrointestinal), asthma (respiratory), psoriasis (dermatology)). Cluster 4 (n=10) had a peak at around 0-10 years of age and included respiratory diseases (n=5) and infections (n=4). Notably, all infectious diseases were in this cluster.

### **Diseases in the same age-of-onset cluster show higher genetic similarity**

Using linear mixed models implemented in BOLT-LMM<sup>22</sup>, we performed GWAS for case versus control on each disease separately and included approximately 10 million common variants that pass quality control (see Methods). Considering associations with the literature-standard p-value lower than  $5 \times 10^{-8}$  as significant<sup>23,24</sup>, we next quantified the associations for each disease, category, and age-of-onset cluster (Figure S20). The major histocompatibility complex (MHC) region is excluded from all analyses, as in the literature, because of its unusually high effect sizes and LD patterns (chr6: 28,477,797 - 33,448,354)<sup>25,26</sup>. Out of 116 diseases, 36 had no significant association and the total number of polymorphic sites with at least one significant association was 93,817. The maximum number of significant associations per disease was 35,001 (hypertension) and the median and mean were 13.5 and 1389.3, respectively. We also checked if diseases from different age-of-onset clusters vary in the number of associations. Cluster 4 had hardly any significant associations (the disease with the maximum number of associations had only 3 significant variants). Although cluster 1 had the highest number of significant associations on average, the values across clusters 1, 2, and 3 were not significantly different (Figure S20b). Moreover, endocrine, immunological, cardiovascular diseases had the highest number of associations and infections had the lowest (Figure S20c). Only 1% of the significant polymorphisms (n=932) were in coding regions, and of these 49% (n=452) were missense and only 1% (n=10) were nonsense. We further found that 47% of significant variants (n=43,810) were associated with multiple diseases, but only ~9% were associated with multiple diseases from different categories (n=8,048) and again ~9% with different age-of-onset clusters (n=8,801) (Figure S21).

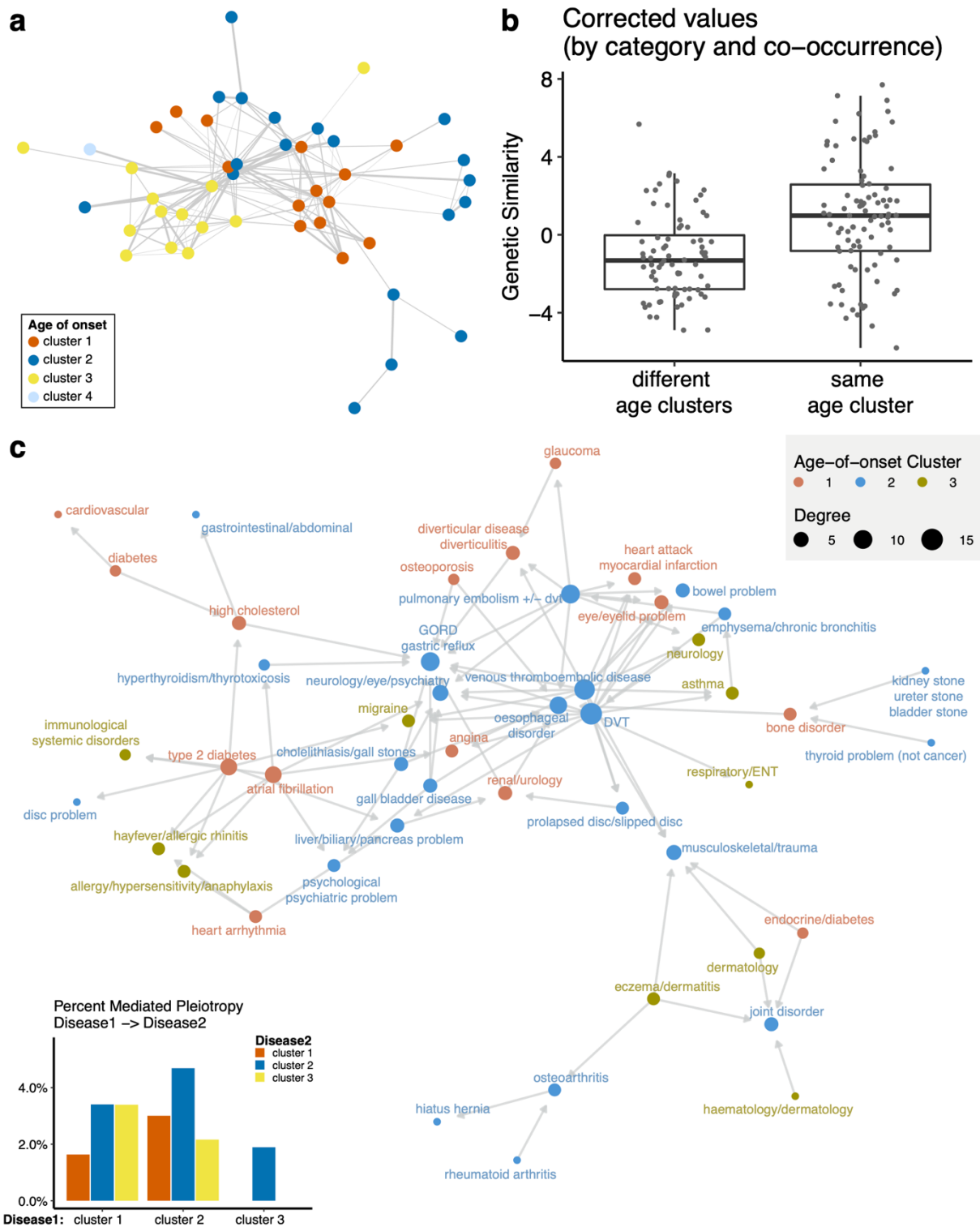


Figure 3: (a) Network representation of the genetic similarities. Nodes ( $n=47$ ) show diseases with a significant genetic similarity to at least one disease and are colored by the age-of-onset cluster. Edges ( $n=167$ ) are weighted by the genetic similarity corrected by disease categories and co-occurrences. (b) The difference between genetic similarity within and across the age-of-onset clusters. The y-axis shows genetic similarity corrected by category and co-occurrence (raw values are available in Figure S22). The x-axis groups similarities into different or same age-of-onset clusters. (c) Network representation of the causal relationships between diseases. Each node ( $n=48$ ) shows a disease, colored by the age-of-onset cluster. Size of the nodes represent the number of significant causal relationships between diseases, including both in and out degrees. Arrows show the causal relationship between pairs with FDR corrected  $p \leq 0.01$  and GCP  $> 0.6$ . The inset bar plot shows the percent significant causal

relationships among all possible pairs (y-axis) between disease 1 (x-axis) and disease 2 (bars colored by the age-of-onset).

We next sought to characterize the genetic similarities between diseases using a score that shows the excess of overlapping associations between diseases, given the number of significant associations for each disease (see Methods). Importantly, we calculated genetic similarities between 80 diseases that have at least one significant association, excluding the pairs that are vertically connected (*i.e.* ancestors to child) in the disease hierarchy (*e.g.* essential hypertension and hypertension). We found 47 significant overlaps and diseases with similar age-of-onset profiles showed a higher genetic similarity, even when controlled for disease categories and co-occurrences (F-test  $p=1.19 \times 10^{-8}$ , Figure 3a-b). Moreover, this trend was reproducible when each cluster was analyzed separately (Figure S23). While correcting for the disease categories and co-occurrences, some true positive signals may be removed from the analysis. However, this correction is necessary, as we used the same cohort for multiple diseases and, thus, diseases that co-occur use the same set of samples. Nevertheless, we retained a significant signal even after this correction, demonstrating that diseases with a similar age-of-onset profile show increased genetic similarity compared to those with different profiles, suggesting shared genetic associations (Figure 1).

We further confirmed the results using 1,703 previously defined LD blocks<sup>27</sup> instead of considering all SNPs as independent. There was no significant genetic similarity between diseases from different age-of-onset clusters (Figure S24) and the similarities within the same age-of-onset cluster were not explained by the disease categories ( $p=0.89$ ) and co-occurrences ( $p=0.15$ ).

### **Mediated pleiotropy does not explain higher genetic similarities within age-of-onset clusters**

Following the models described in Figure 1, we next asked if mediated pleiotropy, rather than a common etiology, may explain higher similarity within age-of-onset clusters. Using a recent methodology developed by O'Connor & Price, we tested for partial or fully causal relationship between diseases<sup>28</sup>. In particular, the method identifies if a latent causal variable (LCV) mediates the genetic correlation between diseases. Using a genetic causality proportion, it assigns a causal relationship if one of the diseases is more strongly correlated with the LCV. The authors report that, unlike mendelian randomization, this method can distinguish between the correlation due to common etiology and causation. We tested for potential causation between 60 diseases, excluding the ones with less than 10 significant genetic variants and low heritability estimates ( $Z_h < 7$ )<sup>28</sup>. Also, similar to genetic similarities, we did not calculate the causation between diseases that are vertically connected in the disease hierarchy. Following the same significance criteria proposed in the methods article (FDR corrected  $p \leq 0.01$  and mean Genetic Causality Proportion (GCP)  $> 0.6$ ), we found significant evidence for full or partial genetic causality in 91 disease pairs between 48 out of 116 diseases in our analysis (Figure 3c, Table S2). Using Fisher's exact test, we tested if mediated pleiotropy was more common between diseases in certain age-of-onset clusters but did not find any significant difference (FDR corrected  $p > 0.1$  for all comparisons, inset bar plot in Figure 3c). Thus, although we detected mediated pleiotropy between some diseases, higher genetic similarities within the same age-of-onset clusters (Figure 2a-b) were not explained in this way and were more likely to be driven by common etiologies.

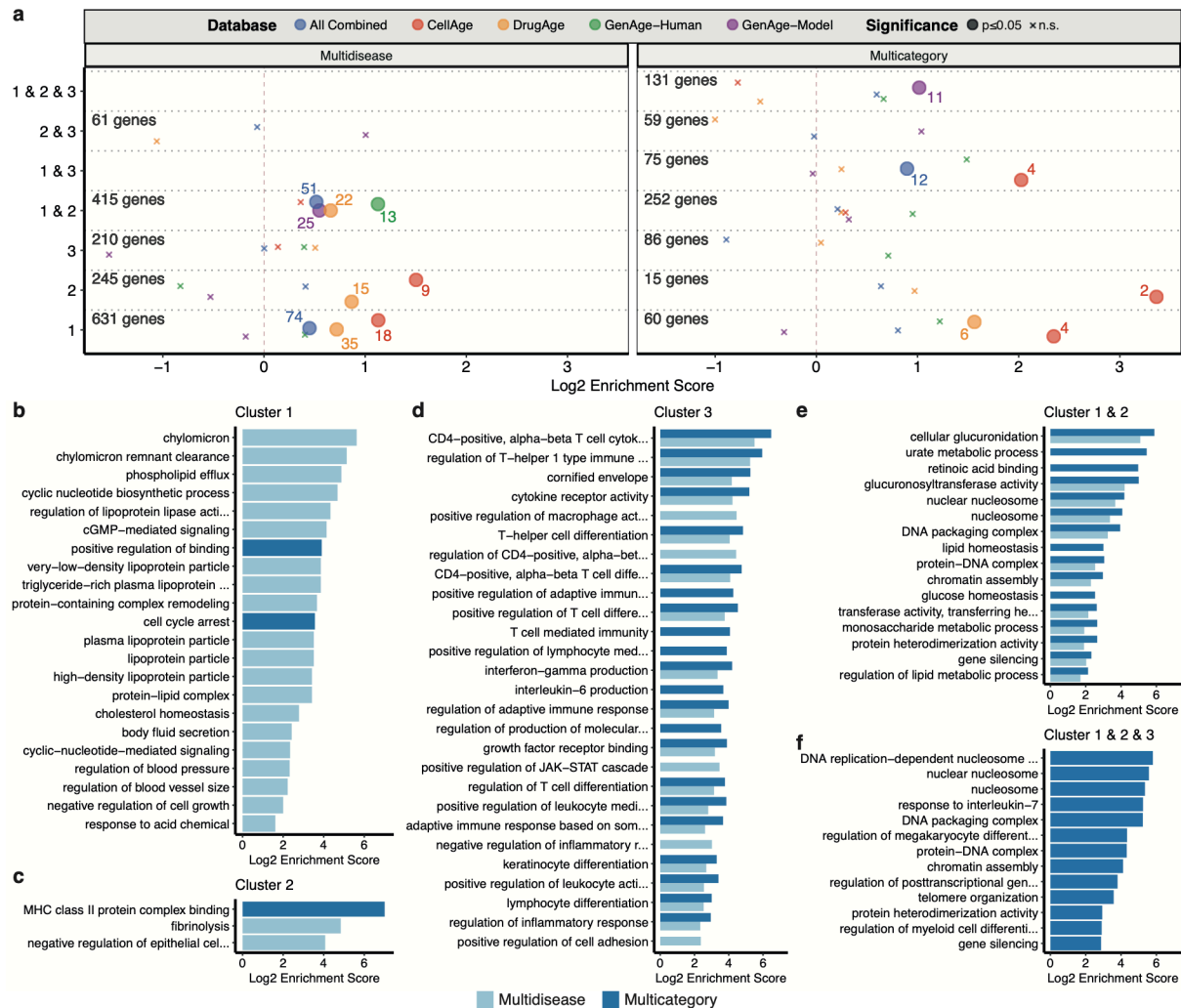
We also investigated the diseases with the highest involvement in mediated pleiotropy. DVT ( $n=14$ ), venous thromboembolic disease ( $n=13$ ), and pulmonary embolism ( $n=9$ ) had the highest number of *out* degrees, meaning they were found as causal for multiple diseases, including all 3 age-of-onset clusters and 5 different disease categories. Gastro-esophageal reflux (GORD)/gastric reflux ( $n=10$ ) and esophageal disorder ( $n=8$ ), on the other hand, had the highest number of *in* degrees, meaning there are multiple diseases detected as causal. These causal diseases spanned 5 disease categories and age-of-onset clusters 1 and 2.

### **Genes associated with the age-dependent disease clusters overlap with known cellular senescence and longevity modulators**

We next mapped all variants to genes based on proximity or known eQTLs using the GTEx eQTL associations<sup>30</sup> (see Methods). To assess the reproducibility of the genes identified, we compared the significant hits with all those reported in the GWAS Catalog. We verified that most of the diseases had significant overlaps with the same or associated traits in the GWAS-Catalog (e.g. osteoporosis and bone density), confirming that our results were reproducible with independent data (Table S3). We next compiled the genes associated with multiple diseases and multiple categories and grouped them based on the age-of-onset cluster of the associated diseases (Table S4). In particular we created two sets of genes, 'multidisease' and 'multicategory', for clusters 1, 2, and 3. We excluded cluster 4 because the number of variants significantly associated with this cluster was low (n=7 associated with 5 diseases), mapping to only 2 genes (*ZBP2*, *NPC1L1*). We also compiled genes associated with multiple diseases or categories in combinations of different age-of-onset clusters. Importantly, genes associated with multiple clusters are not in the gene sets for individual clusters as the latter only include the genes specific to individual clusters, *i.e.* cluster 1, cluster 2 and cluster 1 & 2 genes were all include mutually exclusive sets.

Using these lists, we sought to understand if the common genes between diseases with the same age-of-onset profile had previously been associated with aging. We compared the *multidisease* and *multicluster* gene lists with the literature-based aging databases: GenAge human (genes associated lifespan in humans or closely related species), human orthologs in GenAge model organism (genes modulating lifespan in model organisms), CellAge (genes regulating cellular senescence), DrugAge targets (drugs modulating lifespan in model organisms), and all databases combined<sup>31-33</sup> (Figure 4a). In general, genes associated with clusters 1 and 2, but not Cluster 3, showed significant enrichment with known aging-related genes. The list of overlapping genes is given in Table S5. The CellAge database showed the largest number of significant overlaps, with genes associated with clusters 1, 2, and '1 & 3'. DrugAge targets had a significant overlap with clusters 1, 2, and '1 & 2'. GenAge Human only had significant association with genes associated with cluster '1 & 2'. GenAge model organism data significantly overlapped with genes associated with cluster '1 & 2' and all clusters (1 & 2 & 3). In conclusion, although the association is established through a small subset of genes as also reported in the literature<sup>17,18</sup>, the clusters 1 and 2, constituting age-dependent profiles, shared a significant genetic component with known longevity- and senescence-modulators, while cluster 3 did not.





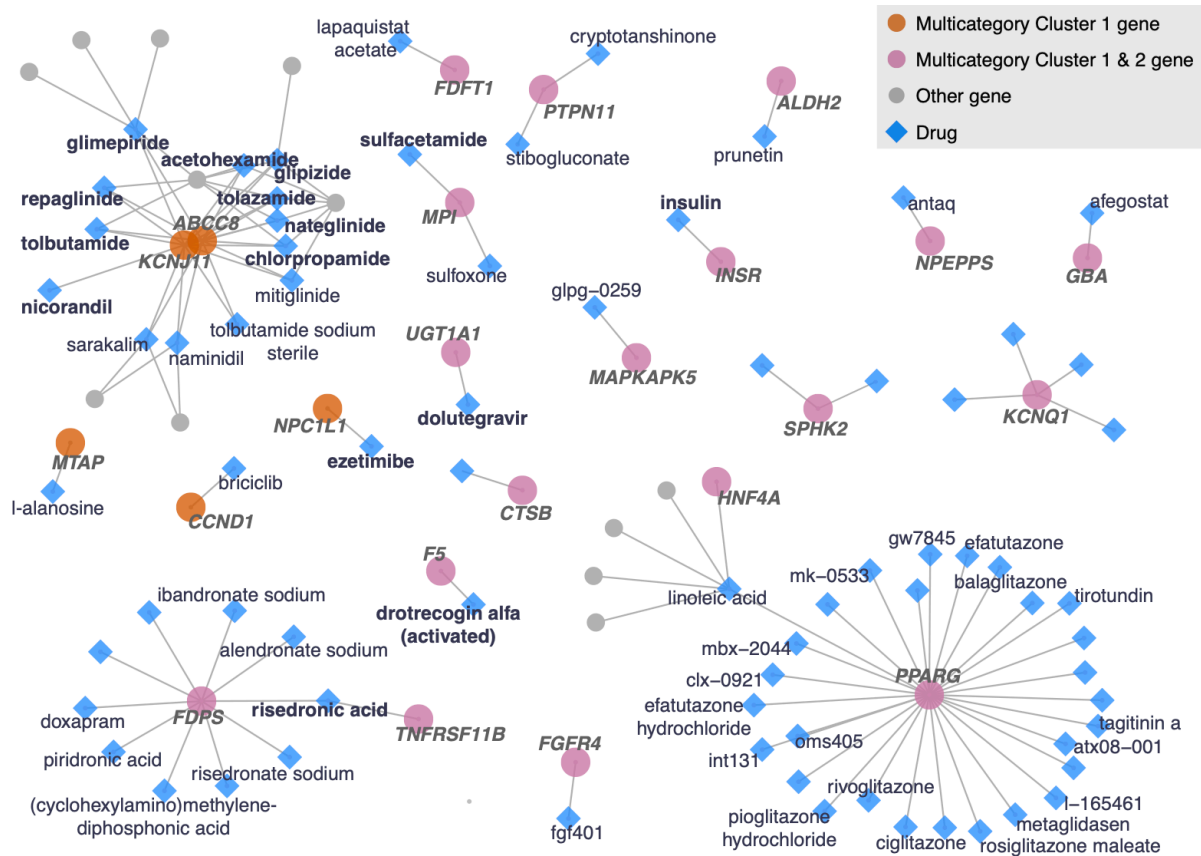
**Figure 4:** a) Overlap between known aging-related genes in databases and genes associated with diseases in different clusters. The x-axis shows log<sub>2</sub> enrichment score, and the y-axis shows the age-of-onset clusters. The numbers of genes in each cluster (for both Multidisease and Multicategory genes) are given. The size of the points shows the statistical significance (large points show marginal  $p$ -value  $\leq 0.05$ , small 'x' indicate non-significant overlaps) and the color shows different databases. The colored numbers near the points show the numbers of overlapping genes. b-f) Gene Ontology (GO) Enrichment results for genes associated with diseases in b) Cluster 1, c) Cluster 2, d) Cluster 3, e) Cluster '1 & 2', f) Cluster '1 & 2 & 3'. Representative GO categories for significantly enriched categories (BY-adjusted  $p$ -value  $\leq 0.05$ ) are listed on the y-axis (see Methods). Log<sub>2</sub> enrichment scores are given on the x-axis. The color of the bar shows the result for multidisease and multicategory genes. There was no significant enrichment for cluster 1 & 3 and 2 & 3.

### Genes associated with different age-of-onset clusters have different functions

Gene Ontology (GO) enrichment analyses were applied to the gene lists, including Biological Process (BP), Molecular Function (MF) and Cellular Component (CC) categories. Cluster 1 was associated with many lipoprotein-related categories, cellular signaling, cellular response, cell cycle arrest, and blood pressure (Figure 4b). Cluster 2 showed association to MHC class II binding, fibrinolysis, and negative regulation of epithelial cell (Figure 4c). Cluster 3 had associations to many immune-related categories and cell adhesion (Figure 4d). Genes in clusters '1 & 3', and in '2 & 3' did not have any significant associations. Genes associated with cluster '1 & 2' were related to nucleosome complex, gene silencing, glucose homeostasis, retinoic acid binding (Figure 4e). Genes associated with at least one disease in all clusters ('1 & 2 & 3') showed association with interleukin-7 response, differentiation, telomere as well as nucleosome complex and gene silencing (Figure 4f). Since cluster 3 did not have an age-dependent profile, the association with gene silencing and nucleosome complex could

represent pleiotropic genes in general. Here we listed the categories that are representative to all other significant functional groups. The full list is given in Table S6, and the procedure of selecting representatives is described in the Methods. Overall, these results suggest that, although cluster 1 and cluster 2 genes were both linked to previously identified aging-related genes, they have distinct functional profiles.

## Drug repurposing to improve healthspan

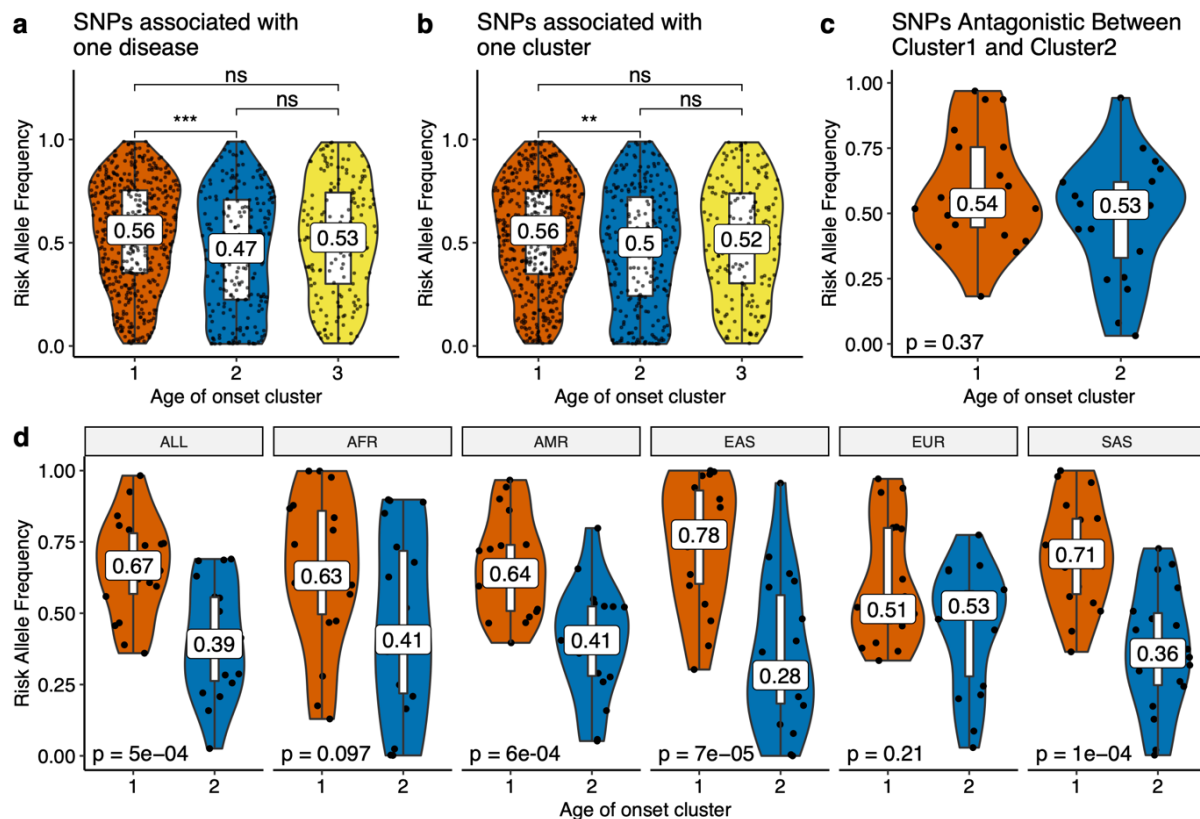


**Figure 5:** 'Drug-target gene' interaction network for the drugs that specifically target multicategory cluster 1, cluster 2 or cluster '1 & 2' genes as determined by Fisher's exact test. Blue diamonds show the drugs with significant association or targeting only one gene in these gene groups. Diamonds without written names are only represented with the ChEMBL IDs in the datasets and did not have names. Drug labels written in bold are drugs approved for different conditions. Circles represent the genes targeted by the significant hits, colored by their age-of-onset cluster. Gray circles show the genes targeted by these drugs but are not among the gene set of interest.

Identification of drugs that can target the multicategory genes associated with diseases in clusters 1 and 2 could enable the treatment of many diseases simultaneously and improve healthspan in the elderly. Thus, we investigated if there are drugs that target these genes specifically ( $p \leq 0.01$  or having only one specific target, Figure 5). We found drugs targeting multicategory cluster 1 genes i) *ABCC8* and *KCNJ11*, which code for parts of K-ATP channels, ii) *CCND1*, iii) *MTAP*, iv) *NPC1L1*. There were also several drugs targeting multicategory genes associated with both cluster 1 and 2 diseases, such as *PPARG*, *INSR*, *FGFR4*, *MAPKAPK5*, *ALDH2*, *PTPN11*. One of the drugs we identified, prunetin (targeting *ALDH2*), was previously shown to increase the lifespan of male *Drosophila melanogaster*<sup>34</sup>. Importantly, the significant hits included approved drugs for 14 conditions, including diabetes, hyperlipidemia, osteoporosis, cardiovascular diseases (list of all drugs and indications available in Table S7). Although the majority of these conditions are age-related, drugs used to treat these conditions do not necessarily target the multicategory genes we identified

(Figure S25), and thus, the drugs identified here offer new possibilities to prevent polypharmacy in the elder population if their use is prioritized to treat multiple diseases. Moreover, some of these drugs are already considered for multiple diseases from different categories. For example, acetohexamide, which targets the K-ATP channel, is in use for diabetes mellitus and is undergoing clinical trials for cataracts<sup>35</sup>.

## Evolution of aging and age-related diseases



**Figure 6:** Risk allele frequency distributions (y-axis) for different age-of-onset clusters (x-axis) in the UKBB for a) SNPs associated with one disease (excluding antagonistic associations), b) SNPs specific to one cluster (excluding antagonistic associations) (ns:  $p > 0.05$ , \*:  $p \leq 0.05$ , \*\*:  $p \leq 0.01$ , \*\*\*:  $p \leq 0.001$ , \*\*\*\*:  $p \leq 0.00001$ ), and c) SNPs that have antagonistic association with cluster 1 and 2 (excluding agonists between cluster 1 and 2). d) The same as panel c but for different 1000 Genomes super-populations (ALL: complete 1000 Genomes cohort, AFR: African, AMR: Ad Mixed American, EAS: East Asian, EUR: European, SAS: South Asian).

Lastly, we sought to understand the abundance of disease-associated variants in the population and their relationship with the evolutionary theories of aging. We first hypothesized that, according to the mutation accumulation theory of aging, SNPs associated with later-onset diseases (Cluster 1) would have a higher frequency than the SNPs associated with diseases that occur at earlier ages (Clusters 2 and 3), which are presumably under stronger selection pressure. Following a similar methodology to Rodriguez et al. to compare allele frequencies associated with different traits<sup>36</sup>, we compared the allele frequencies associated with different age-of-onset clusters. As SNPs which are close together in the genome are expected to have similar allele frequencies due to linkage, we calculated the median risk allele frequency for SNPs within previously defined LD blocks<sup>27</sup>. Supporting the mutation accumulation theory of aging, diseases of cluster 1 had significantly higher risk allele frequencies than cluster 2, both for the SNPs associated with one disease (Figure 6a, Wilcoxon test  $p = 0.00033$ ) or with one cluster (Figure 6b, Wilcoxon test  $p = 0.0068$ , also confirmed by bootstrapping  $n = 100$  loci for  $B = 1,000$ ; Figure S26). We further confirmed that this trend is not specific to the UK population, as we obtained comparable results in all super-populations of the 1000 Genomes Project<sup>37</sup>

(Figure S27-28). Variants associated with cluster 3, which includes immune-related diseases, were not significantly different from those associated with cluster 1, although cluster 3 diseases can occur even at an earlier age. Moreover, although the difference in the median allele frequencies was not significant, the shapes of distributions were different, with a significant shift towards higher risk allele frequencies only in cluster 1 ( $p_{cl1}=0.05$ ,  $p_{cl3}=0.86$  calculated using 10,000 permutations). High minor allele frequencies, thus a higher variation we observed in cluster 3 is in line with the previous suggestion that immune-related genes are under long-term balancing selection in humans<sup>38</sup>, although positive selection also influences immunity<sup>39-41</sup>.

To test the antagonistic pleiotropy theory (AP), we first asked if the diseases with different onsets have an excess of antagonistic SNPs. Similar to a previous study<sup>36</sup>, we defined a pleiotropic biallelic SNP as *agonistic* if the risk allele is the same for different diseases, and as *antagonistic* if opposite alleles are associated with increased risk for different diseases. If one of these diseases is under a stronger negative selection, then the risk allele of the other disease could increase over time. Comparing the proportion of agonist and antagonist SNPs within and between the age-of-onset clusters, we found that there is an excess of antagonistic pleiotropy between diseases with different age-of-onset profiles (Fisher's exact test  $p<0.001$ , Table S8). Next, we tested the differences in risk allele frequencies between the clusters, as AP predicts a higher risk allele frequency for late-onset diseases. Interestingly, the difference between the risk allele frequencies for cluster 1 and cluster 2 was not significant for the UKBB population (Figure 6c). However, all 1000 Genome super-populations except for Europeans had higher risk allele frequencies for cluster 1 diseases (Figure 6d). We hypothesized that this is mainly due to false positive disease associations in the UKBB, due to increased power when testing the antagonistic associations with frequency closer to 0.5. We thus investigated the allele frequency differences for the significant variants with increased effect sizes. Indeed, associations with a larger effect size showed the expected differences in allele frequencies, although the number of independent loci was limited (Figure S29). We also examined the type of diseases and genes associated with antagonistic pleiotropy. The main driver of the pattern was the loci with *ABCG8* and *ABCG5* genes, showing antagonistic relationship for high cholesterol (cluster 1) and other lipid-related diseases in cluster 2, such as gallbladder disease and cholelithiasis. Another locus included variants that show antagonistic relationship with cardiovascular disease (cluster 1) and the cluster 2 diseases gout (*ADH1B*), osteoarthritis and joint disorder (*SLC39A8*), and osteoarthritis (*BANK1*). Another potential candidate was a locus associated with hypertension (cluster 1) and musculoskeletal diseases (cluster 2), but this locus included multiple candidate genes (Table S9). Nevertheless, our comparison is between common diseases that occur after the age of 20 and 40, which are both after the average age at first reproduction and therefore the start of the decrease in the force of natural selection<sup>42</sup>. Thus, a better comparison would include the mutations causing rare developmental diseases, which are not available in the UKBB.

## Discussion

The number and the incidence of diseases increase with age. In this study, we explored whether this results from a common genetic component among ARDs, which might also be linked to aging. We compared genetic associations and age-of-onset distributions of 116 self-reported diseases in the UKBB and found shared variants, genes and pathways, which were also associated with aging.

Using an unsupervised, data-driven approach to classify diseases based on their age-of-onset profiles, we found 4 main clusters; i) diseases that rapidly increase after 40 years of age, ii) diseases that increase after 20 years of age, iii) diseases with no age-related pattern, and iv) diseases that peak at around 10 years of age. Notably, unlike previous studies<sup>18,43</sup>, by using this unsupervised approach, we detect a distinction between cluster 1 and cluster 2, which

both show age-dependency but distinct age-of-onset distributions. These two clusters were associated with genes with different functional and evolutionary characteristics, although they both overlap with known aging-related genes.

Based on genetic associations, diseases with similar age-of-onset profiles showed a higher genetic similarity on average, compared to diseases in other clusters, even when controlled for disease categories and co-occurrences (Figure 2a-b). Moreover, this similarity within age-of-onset clusters was not explained by mediated pleiotropy, in which one of the diseases is causal for the other one, suggesting instead a common etiology. We then studied the genes involved and found that genes associated with clusters 1 and 2 (both constituting ARDs) are enriched with known longevity- and senescence-modulators, while genes associated with cluster 3, which does not show an age-dependent profile, did not show this enrichment. In addition, we found that genes associated with different age-of-onset clusters have different functions. Comparing the risk allele frequencies of variants associated with different age-of-onset profiles, we found support for both mutation accumulation and antagonistic pleiotropy theories of aging, although the number of independent loci supporting the second was limited. Lastly, we identified drugs that can target the common genetic component between ARDs, which may also limit the multimorbidity and polypharmacy associated with late life.

In this study, we had a limited age range, covering individuals up to 65 years old and thus, could not analyze diseases of later ages. Neither did we consider the cancers or changes in regulation of gene expression, which are affected not only by aging, but also various environmental or intrinsic factors. Future cohorts with a broader age range and spanning multi-omics data, somatic mutations, health outcomes, and lifestyle information, will enable a better understanding of the genetic mechanisms of age-of-onset determination and establishing the causal link with candidate genes. Despite these limitations, we present a novel approach to study ARDs using an unbiased, data-driven approach and show that ARDs share common genetic associations linked to aging. We suggest that targeting the common pathways between multiple ARDs could offer compression of late life multimorbidity as well as alleviating the effects of polypharmacy.

### **Software Availability**

All the code used to perform analyses is available in GitHub:  
[https://github.com/mdonertas/ukbb\\_ageonset](https://github.com/mdonertas/ukbb_ageonset)

### **Data Availability**

The full set of GWAS results from this study can be accessed using BioStudies (S-BSST407) and all results generated in the analysis are provided as Supplementary Datasets and Tables.

### **Ethics Statement**

#### **Conflict of interest**

The authors declare that they have no conflict of interest.

### **Author Contributions**

H.M.D conceived and designed the study with contributions from L.P. and J.M.T.. H.M.D analyzed the data with the help of D.K.F and M.F.V.. H.M.D. interpreted the results and wrote the manuscript with contributions from all authors. All authors read, revised and approved the final version of this manuscript.

### **Acknowledgment**

This research has been conducted using the UK Biobank Resource (application no. 30688). The authors thank the GWAS-Catalog team for providing the list of studies using UK Biobank data; James Stephenson and Roman Laskowski for their help in running VarMap tool; and

Mehmet Somel, Susan Ozanne, Pedro Beltrao, and Wolfgang Huber for fruitful discussions. H.M.D. is a member of Darwin College, University of Cambridge.

### Funding Statement

This work is funded by EMBL (H.M.D., J.M.T.), the Wellcome Trust (098565/Z/12/Z; L.P., J.M.T), and Comisión Nacional de Investigación Científica y Tecnológica - Government of Chile (CONICYT scholarship; M.F.V.).

### References

1. López-Otín, C., Blasco, M. A., Partridge, L., Serrano, M. & Kroemer, G. The hallmarks of aging. *Cell* **153**, 1194–1217 (2013).
2. Crimmins, E. M. Lifespan and Healthspan: Past, Present, and Promise. *Gerontologist* **55**, 901–911 (2015).
3. Partridge, L., Deelen, J. & Slagboom, P. E. Facing up to the global challenges of ageing. *Nature* **561**, 45–56 (2018).
4. Niccoli, T. & Partridge, L. Ageing as a risk factor for disease. *Curr. Biol.* **22**, R741-52 (2012).
5. Flatt, T. & Partridge, L. Horizons in the evolution of aging. *BMC Biol.* **16**, 93 (2018).
6. Medvedev, Z. A. An attempt at a rational classification of theories of ageing. *Biol. Rev. Camb. Philos. Soc.* **65**, 375–398 (1990).
7. Medawar, P. B. Unsolved problem of biology. *Med. J. Aust.* **1**, 854–855 (1953).
8. Williams, G. C. Pleiotropy, Natural Selection, and the Evolution of Senescence. *Evolution* **11**, 398–411 (1957).
9. Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases and traits. *Nat. Genet.* **47**, 1236–1241 (2015).
10. Pickrell, J. K. *et al.* Detection and interpretation of shared genetic influences on 42 human traits. *Nat. Genet.* **48**, 709–717 (2016).
11. Cross-Disorder Group of the Psychiatric Genomics Consortium *et al.* Genetic relationship between five psychiatric disorders estimated from genome-wide SNPs. *Nat. Genet.* **45**, 984–994 (2013).
12. Cortes, A., Albers, P. K., Dendrou, C. A., Fugger, L. & McVean, G. Identifying cross-disease components of genetic risk across hospital data in the UK Biobank. *Nat. Genet.* **52**, 126–134 (2020).
13. Solovieff, N., Cotsapas, C., Lee, P. H., Purcell, S. M. & Smoller, J. W. Pleiotropy in complex traits: challenges and strategies. *Nat. Rev. Genet.* **14**, 483–495 (2013).
14. Ellinghaus, D. *et al.* Analysis of five chronic inflammatory diseases identifies 27 new associations and highlights disease-specific patterns at shared loci. *Nat. Genet.* **48**, 510–518 (2016).
15. Parkes, M., Cortes, A., van Heel, D. A. & Brown, M. A. Genetic insights into common pathways and complex relationships among immune-mediated diseases. *Nat. Rev. Genet.* **14**, 661–673 (2013).
16. Bien, S. A. & Peters, U. Moving from one to many: insights from the growing list of pleiotropic cancer risk genes. *Br. J. Cancer* **120**, 1087–1089 (2019).
17. Johnson, S. C., Dong, X., Vijg, J. & Suh, Y. Genetic evidence for common pathways in human age-related diseases. *Aging Cell* **14**, 809–817 (2015).
18. Fernandes, M. *et al.* Systematic analysis of the gerontome reveals links between aging and age-related diseases. *Hum. Mol. Genet.* **25**, 4804–4818 (2016).
19. Wang, J., Zhang, S., Wang, Y., Chen, L. & Zhang, X.-S. Disease-aging network reveals significant roles of aging genes in connecting genetic diseases. *PLoS Comput. Biol.* **5**, e1000521 (2009).
20. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
21. Chouakria, A. D. & Nagabhushan, P. N. Adaptive dissimilarity index for measuring time



- series proximity. *Adv. Data Anal. Classif.* **1**, 5–21 (2007).
22. Loh, P.-R. *et al.* Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat. Genet.* **47**, 284–290 (2015).
  23. Pe'er, I., Yelensky, R., Altshuler, D. & Daly, M. Estimation of the Multiple Testing Burden for Genomewide Association Studies of Common Variants. *Nature Precedings* (2007) doi:10.1038/npre.2007.359.1.
  24. Panagiotou, O. A., Ioannidis, J. P. A. & Genome-Wide Significance Project. What should the genome-wide significance threshold be? Empirical replication of borderline genetic associations. *Int. J. Epidemiol.* **41**, 273–286 (2012).
  25. MHC region of the human genome - Genome Reference Consortium. <https://www.ncbi.nlm.nih.gov/grc/human/regions/MHC?asm=GRCh37>.
  26. Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291–295 (2015).
  27. Berisa, T. & Pickrell, J. K. Approximately independent linkage disequilibrium blocks in human populations. *Bioinformatics* **32**, 283–285 (2016).
  28. O'Connor, L. J. & Price, A. L. Distinguishing genetic correlation from causation across 52 diseases and complex traits. *Nat. Genet.* **50**, 1728–1734 (2018).
  29. Mungan, Z. & Pınarbaşı Şimşek, B. Which drugs are risk factors for the development of gastroesophageal reflux disease? *Turk. J. Gastroenterol.* **28**, S38–S43 (2017).
  30. Gamazon, E. R. *et al.* Using an atlas of gene regulation across 44 human tissues to inform complex disease- and trait-associated variation. *Nat. Genet.* **50**, 956–967 (2018).
  31. Tacutu, R. *et al.* Human Ageing Genomic Resources: new and updated databases. *Nucleic Acids Res.* **46**, D1083–D1090 (2018).
  32. Avelar, R. A., Ortega, J. G., Tacutu, R., Tyler, E. & Bennett, D. A Multidimensional Systems Biology Analysis of Cellular Senescence in Ageing and Disease. *BioRxiv* (2019).
  33. Barardo, D. *et al.* The DrugAge database of aging-related drugs. *Aging Cell* vol. 16 594–597 (2017).
  34. Piegholdt, S., Rimbach, G. & Wagner, A. E. The phytoestrogen prunetin affects body composition and improves fitness and lifespan in male *Drosophila melanogaster*. *FASEB J.* **30**, 948–958 (2016).
  35. Compound: Acetohexamide.  
[https://www.ebi.ac.uk/chembl/compound\\_report\\_card/CHEMBL1589/](https://www.ebi.ac.uk/chembl/compound_report_card/CHEMBL1589/).
  36. Rodríguez, J. A. *et al.* Antagonistic pleiotropy and mutation accumulation influence human senescence and disease. *Nat Ecol Evol* **1**, 55 (2017).
  37. 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
  38. Bitarello, B. D. *et al.* Signatures of Long-Term Balancing Selection in Human Genomes. *Genome Biol. Evol.* **10**, 939–955 (2018).
  39. Kosiol, C. *et al.* Patterns of positive selection in six Mammalian genomes. *PLoS Genet.* **4**, e1000144 (2008).
  40. Nielsen, R. *et al.* A scan for positively selected genes in the genomes of humans and chimpanzees. *PLoS Biol.* **3**, e170 (2005).
  41. Shultz, A. J. & Sackton, T. B. Immune genes are hotspots of shared positive selection across birds and mammals. *Elife* **8**, (2019).
  42. Fisher, R. A. The genetical theory of natural selection. **272**, (1930).
  43. Wolfson, M., Budovsky, A., Tacutu, R. & Fraifeld, V. The signaling hubs at the crossroad of longevity and age-related disease networks. *Int. J. Biochem. Cell Biol.* **41**, 516–520 (2009).

## Methods

### UK Biobank Data

Data was downloaded using bash and following the guidelines provided by the UK Biobank.

#### Sample quality control

After excluding all samples from individuals who have withdrawn their data from UK Biobank, we first filtered out all samples without genotypes (N = 14,248). Then, we used the following criteria for the remaining 488,295 samples.

Discordant sex: Data includes two entries for sex: 1) self-reported and 2) genetic sex determined using the call intensities on sex chromosomes. There are multiple reasons why these two entries may not correspond, such as sample mishandling, errors in data input, transgender individuals, and sex chromosome aneuploidies<sup>1,2</sup>. Since we used sex as a covariate in our GWAS model, we preferred to be cautious about this issue and excluded all cases where the genetic sex and self-reported sex did not correspond and all cases where sex chromosome aneuploidy was detected. Specifically, we used the fields '31-0.0' (Sex) and '22001-0.0' (Genetic sex) to compile discordant information. There were 235 self-reported males being identified as female by the genetics, and 143 self-reported females being identified as males by the genetics. We excluded these 378 cases, 0.077% of the data. Moreover, field '22019-0.0' (Sex chromosome aneuploidy) is used to exclude cases with sex chromosome aneuploidy. There were 651 cases of aneuploidy, 0.133% of all data. 181 of these cases (27.80% of aneuploidy cases) were also detected as discordant information in the first step. This corresponds to 47.88% of discordant sex cases. Overall, we identified 848 samples to be excluded based on this criterion.

Genotype call rate & Heterozygosity: Genotype missingness and heterozygosity are widely used as a measure of DNA sample quality. For quality filtering based on missingness and heterozygosity we only used the suggested exclusions by UK Biobank. Specifically, we used the field '22010-0.0' (Recommended genomic analysis exclusions) and determined the cases with 'poor heterozygosity/missingness' (N = 469). We next used the field '22018-0.0' (Genetic relatedness exclusions) and noted down the cases with 'Participant self-declared as having a mixed ancestral background' (N = 692), and the cases with 'High heterozygosity rate (after correcting for ancestry) or high missing rate' (N = 840). Lastly, there were 968 cases that are suggested as outliers for heterozygosity or missing rate, field '22027-0.0' (Outliers for heterozygosity or missing rate). We then checked the scatter plots for logit(Missingness) vs. Heterozygosity for each Ethnic Background, in accordance with the identification of samples to exclude by the UK Biobank<sup>2</sup> (Figure S30). Logit transformation is used to linearize sigmoidal distribution of missingness. Investigation of heterozygosity can detect DNA sample contamination, inbreeding, or mixed ethnicity<sup>1</sup>. This quality check reveals when people with a mixed ethnicity tend to have a higher heterozygosity, even after correcting for PCs. We confirmed these are in accordance with the original article and excluded the samples suggested by the UK Biobank.

Overall, there were 3,697 samples excluded based on these two criteria. Please note that the numbers presented above may not add up to this number, because there were some samples excluded based on multiple criteria. The percent overlap across multiple criteria is given in Figure S31.

#### Preparing the trait data

Using the samples that passed the quality control (N=484,598), we subsetted the data so that it included only the baseline visit. Apart from the data that is already available in UK Biobank, we calculated some other values: 1) *BMI*: Using the columns for 'Weight' and 'Standing height' we calculated BMI as:  $\text{Weight} / (\text{Standing Height} / 100)^2$ , 2) *Parent Age at Death - Minimum*:



The youngest age at which either parent died. 3) *Parent Age at Death - Maximum*: The age of death for the parent who lived longest. 4) *Parent Age at Death - Average*: The average age of death for the two parents. If neither of the parents died, or if the data was unavailable, these values (2-4) were set to be NA. If only one parent died, we use the corresponding age as both the minimum, maximum, and average. 5) *The number of self-reported non-cancer diseases*: The number of unique self-reported non-cancer illnesses each participant recorded in the baseline recruitment. 6) *The number of self-reported cancers*: The number of unique self-reported cancers each participant recorded in the baseline recruitment. 7) *Self-reported diseases after taking the disease hierarchy into consideration (Propagated disease data)*: The self-reported diseases in UK Biobank are not independent, but rather are organized in a hierarchical manner. Using the relationship information between diseases, we propagated disease-participant associations, upwards, including terms higher up the tree. For example, if a person reports having “essential hypertension”, we also annotate that person with “hypertension”, and “cardiovascular disease”. 8) *Age at diagnosis for the self-reported diseases after taking the disease hierarchy into consideration (Propagated age at diagnosis data)*: We re-defined age at diagnosis using the minimum age at diagnosis for all the diseases that were child term for a particular disease in the disease hierarchy. 9) *The number of self-reported non-cancer diseases after taking the disease hierarchy into consideration (Propagated number of non-cancer diseases)*: The number of unique self-reported diseases each participant records after taking into account the data propagation. 10) *Age when the last deceased person died*: We calculated the age of each person when the last death entry in the UKBB happened. This value is used to calculate the proportion of people who died at a certain age interval in Figure S1c.

#### Selecting diseases to analyze

We calculated the disease occurrences for all self-reported diseases in UK Biobank. Specifically, among the cohort we used, we calculated how many participants and what proportion of males and females reported each disease. Since we analyzed the same set of SNPs that have  $MAF \geq 0.01$  across multiple diseases, to decrease the false positive rate in GWAS, we limited the diseases to a subset with at least 2,000 cases ( $n = 129$  out of 472). Moreover, we only focused on diseases that were common and not sex-limited, *i.e.* we only considered diseases that are seen in 1 in every 1,000 males and females ( $n = 189$  out of 472). The intersection of these two conditions was 116 diseases and we excluded all others.

We only analyzed self-reported non-cancer diseases (field ‘20002’) and did not combine self-reported cancers (field ‘20001’), mainly because i) the number of cases is low (45,224 compared to 384,906 for other diseases), ii) cancer is thought as a result of a complex interaction between germline and somatic mutations<sup>3,4</sup>, whereas the evidence for the effect of somatic mutations in other diseases is limited to rare and neurological disorders<sup>5,6</sup>, iii) the relationship between cancer and aging is complex, *e.g.* while telomere attrition and cellular senescence are thought to be evolved as a tumor suppressor mechanisms; aging-related changes in epigenomic landscape and genomic instability contribute to cancer occurrence<sup>7</sup>. Thus, although a similar analysis using cancers would be interesting, we only focused on non-cancer self-reported diseases in this study. Since we did not exclude the individuals with cancer, we also checked if there is a significant overlap in individuals with cancer with the other diseases we analyzed (Figure S32). However, there was no such association.

## Disease co-occurrence calculations

Table 1: Contingency table for disease comorbidities.

	Disease B	No disease B	Total
Disease A	$N_{ab}$	$N_{anb}$	$T_a$
No disease A	$N_{nab}$	$N_{nanb}$	$T_{na}$

### Relative risk (RR) score

Relative risk is an estimate of having the disease A, when already affected by disease B. Overall it measures if disease A co-occurs with disease B more frequently than expected if these diseases were independent in the population. It is calculated as a fraction between the number of patients diagnosed with both diseases and a random expectation based on disease prevalence<sup>8</sup>. Mathematically it can be expressed as follows, using the values from Table 1:

$$P_{exposed} = \frac{N_{ab}}{T_a}, P_{notexposed} = \frac{N_{nab}}{T_{na}}$$

$$RR = \frac{P_{exposed}}{P_{notexposed}}$$

$$CI = \ln RR \pm 1.96 \sqrt{\frac{\frac{T_a - N_{ab}}{N_{ab}}}{T_a} + \frac{\frac{T_{na} - N_{nab}}{N_{nab}}}{T_{na}}}$$

### $\phi$ value (Pearson correlation for binary variables)

The  $\phi$  value measures the robustness of the association between diseases based on co-occurrences<sup>9</sup>. Mathematically, it can be expressed as:

$$\phi_{AB} = \frac{C_{AB}N - P_A P_B}{\sqrt{P_A P_B (N - P_A)(N - P_B)}}$$

$N$ : the total number of individuals

$P_A$ : Prevalence of disease A

$C_{AB}$ : Number of patients with both diseases

$\phi$  ranges between -1 and 1, where the sign indicates the type of association.

## Disease age-of-onset

### Disease dissimilarity measure

Temporal correlation: In order to calculate dissimilarities among diseases, we use CORT<sup>10</sup> distance as included in R package TSclust<sup>11</sup>. Euclidean distance and dynamic time warping<sup>12</sup> are the two most widely used proximity measures for time series proximity. However, they are both calculated based on the closeness of the values and disregard the growth behavior. Correlation-based measures are also used to calculate the similarity between time series.

However, Pearson correlation overestimates the similarity because of the underlying temporal dependency and Spearman correlation fails to consider the growth rate as it is based on ranks. Chouakria et al., on the other hand, suggested a measure that also considers the proximity-based on growth behavior,  $CORT^{10}$ . Temporal correlation between two time series objects  $S_1=(u_1, u_2, \dots, u_p)$  and  $S_2=(v_1, v_2, \dots, v_p)$  is calculated as follows:

$$CORT(S_1, S_2) = \frac{\sum_{i=1}^{p-1} (u_{(i+1)} - u_i)(v_{(i+1)} - v_i)}{\sqrt{\sum_{i=1}^{p-1} (u_{(i+1)} - u_i)^2} \sqrt{\sum_{i=1}^{p-1} (v_{(i+1)} - v_i)^2}}$$

CORT ranges between -1 and 1. A value of  $CORT = 1$  implies that two time series increase or decrease simultaneously with the same growth rate, whereas a value of -1 shows the same growth rate but in opposite direction. If the value is 0, it means there is no temporal correlation between the series.

Dissimilarity Index: The dissimilarity index suggested by Chouakria et al.<sup>10</sup>, is calculated based on an automatic adaptive tuning function and considers similarity based on both values and behavior, *i.e.* the strength of monotonicity and closeness of the growth rates as calculated by CORT measure introduced in the previous section. They suggest a dissimilarity index D as follows:

$$D(S_1, S_2) = f(CORT(S_1, S_2)) \cdot \delta_{conv}(S_1, S_2)$$

Where  $f(x)$  is an exponential adaptive tuning function:

$$f(x) = \frac{2}{1 + \exp(kx)}, k \geq 0$$

As  $k$  increases, the contribution of behavior increases. We use  $k = 2$  and as a result behavior (CORT) contributes 76.2% to D and values ( $\delta_{conv}$ ) contribute 23.8%. For  $\delta_{conv}$  we used conventional Euclidean distance.

#### Clustering diseases by age-of-onset

We clustered data using 'partition around medoids (PAM)' algorithm<sup>13</sup> based on the distance measure calculated using the previous step. The aim of this algorithm is to minimize the average distance (based on any dissimilarity measure) between the objects and their closest selected medoid object. It works very similarly to k-means, except instead of defining arbitrary points as the means, it defines medoids among the objects. Thus, it can incorporate any distance measure instead of just using the mean distance between points (*i.e.*, euclidean distances). The algorithm first searches for  $k$  number of objects that represent the structure of the data (Here the number  $k$  is assumed to be known a priori but see the next section for the determination of  $k$ ). After finding a set of  $k$  medoids,  $k$  clusters are constructed by assigning each observation to the nearest medoid. Overall, the goal is to find  $k$  representative objects such that the sum of dissimilarities of the observations to their closest representative is as small as possible. After each assignment, medoid and non-medoid data points are swapped and a cost (sum of distances of points to the new medoid) is calculated. If the total cost of configuration is decreased, then the new configuration is maintained, otherwise, it is reversed. We used 'pam' function in the 'cluster' package<sup>14</sup> in R to apply this algorithm.

#### Choosing the optimum number of clusters

The clustering algorithm we used, PAM, clusters data into  $k$  clusters, which is determined by the user. So, even if there is no real structure in data, as we increase the number of clusters, we can get more and more clusters. A potential way to decide on the number of clusters is

using the gap statistic<sup>15</sup>. This value is calculated by comparing the logarithm of the within-sum-of-squares (WSS) to averages from simulated data without any structure.

$$WSS_k = \sum_{l=1}^k \sum_{x_i \in C_l} d^2(x_i, \bar{x}_l)$$

$k$ : number of clusters

$C_l$ : objects in the  $l$ -th cluster

$\bar{x}$ : the average point.

Calculating only WSS, however, is not enough as it would be minimized when each point has its own cluster. Thus, we use the gap statistic which suggests calculating  $\log(WSS_k)$  for a range of values of  $k$  and compare it to that obtained by WSS calculated based on simulated data. So, after WSS is calculated for various values of  $k$ , the algorithm involves generating  $B$  (we choose  $B=1,000$ ) reference datasets, using Monte Carlo sampling from a homogeneous distribution and re-calculate WSS for all  $k$  values. Using these values  $\text{gap}(k)$  statistic is calculated as:

$$\text{gap}(k) = \bar{l}_k - \log(WSS_k)$$

$$\bar{l}_k = \frac{1}{B} \sum_{b=1}^B \log(W_{kb}^*)$$

If the clustering is good (i.e. WSS is small) we expect  $\bar{l}_k$  to be higher than  $\log(WSS)$ . Thus, gap statistic is mostly positive and we are interested in the highest value. Tibshirani et al.<sup>15</sup> suggests using the smallest  $k$  such that,

$$\text{gap}(k) \geq \text{gap}(k+1) - s'_{k+1}$$

where

$$s'_{k+1} = sd_{k+1} \sqrt{1 + \frac{1}{B}} \text{ and } sd_k^2 = \frac{1}{B-1} \sum_{b=1}^B (\log(W_{kb}^*) - \bar{l}_k)^2$$

Using this approach, we determined  $k = 4$ .

## Genome wide association study

### Preparing the files required for GWAS

*Fixing FAM files:* In UK Biobank FAM files, the column for 'phenotype' includes batch that is coded with characters. In order to use BOLT-LMM<sup>16</sup>, we updated all the entries in this column to numeric values<sup>17</sup>.

*'Remove' files for BOLT-LMM:* BOLT-LMM accepts a list of individuals to be removed from the analysis as an input. These files are called 'remove' files and are in the FAM format. We prepared these files for i) withdrawn samples ( $n = 51$ ), ii) samples that failed the quality control ( $n = 3,779$ ), iii) samples that have information in PLINK files but lack BGEN files ( $n = 968$ ).

*Calculating the SNP statistics:* In order to apply a quality filter for SNPs, using PLINK<sup>1</sup>, we calculated i) p-values for each SNP showing whether it deviates from Hardy-Weinberg equilibrium, and ii) Minor allele frequencies (MAF).

*SNP Quality Control:* We excluded SNPs that deviate from Hardy-Weinberg equilibrium ( $p \leq 1e-6$ ,  $n = 202,473$ ) or with a minor allele frequency (MAF) smaller than 0.01 ( $n = 127,969$ ). In total, we discarded 314,697 SNPs (Note that the numbers do not add up as these SNPs can overlap), resulting in 9,886,868 sites.

*Phenotype File:* We created a phenotype file that can be used as an input for BOLT-LMM, including the following fields: sex, age when attended assessment center, calculated BMI, assessment center, ethnicity, batch, first 20 PCs, and self-reported diseases (one column per disease).

### GWAS run using BOLT-LMM

For each disease, we run GWAS separately using BOLT-LMM with the following inputs:

- We remove the samples that are in plink files but now in bgen; samples that did not pass our QC; samples from the individuals who have withdrawn their data from the UKBB
- We excluded the SNPs that deviate from Hardy-Weinberg equilibrium, and have minor allele frequency lower than 0.01.
- We used Sex, Age, BMI, assessment center, ethnicity, batch, and the first 20 PCs as covariates.
- To run the mixed-model, a reference LD score table is required. We used LD scores generated using 1000 Genomes European-ancestry samples, which is provided with the BOLT-LMM download.
- Genetic map for hg19 file provided in the BOLT-LMM website.
- We set 'bgenMinMAF' argument to  $1e-2$  and 'bgenMinINFO' parameter to 0.5 to only include SNPs that pass these criteria.

### GWAS Results

We removed MHC region (chr6: 28,477,797 - 33,448,354) from the analysis and considered positions with a p-value lower than  $5 \times 10^{-8}$  as a significant association.

### **Coding Variants**

We used VarMap<sup>18</sup> to map variants to proteins and domains. VarMap provides detailed information about coding variants, including annotations for the missense, synonymous, and nonsense variations. In our analysis, if a variant is not annotated as a coding variant in VarMap output, we assumed it is non-coding.

### **Genetic similarities between diseases**

In order to calculate the overlap between diseases we used the number of SNPs that are significantly associated with both diseases, but corrected by the number that is expected by chance, if two diseases are independent:

$$Genetic\ Similarity = \frac{N_{common}}{N_{d1} \times N_{d2}} \times N_{total}$$

$N_{common}$ : Number of SNPs in common.

$N_{dx}$ : Number of SNPs associated with disease X.

$N_{total}$ : Total number of SNPs analyzed in the study.

The statistical significance of these genetic similarities is calculated using the binomial test, and the similarity is only considered for downstream analysis if  $p \leq 0.01$ . Moreover, the value is only calculated if two diseases do not have any hierarchical relationships in the disease hierarchy.

In order to assess the genetic similarity within age-of-onset clusters, we further used linear regression to correct log<sub>2</sub> genetic similarity value by disease co-occurrences (risk ratios) and disease categories (binary data showing whether two diseases are of the same category). The 'corrected genetic similarity' is the residuals from this linear model.

#### LD Blocks

In order to assess the similarity between different diseases we use overlaps across significant associations and thus preferred not to do fine mapping. However, a significant challenge is that genomic variations are not independent but instead linked in the genome. To understand the effect of linkage disequilibrium or overcome it, we made use of linkage disequilibrium blocks previously defined for human genome<sup>19</sup>. We repeated the analysis for genetic similarity after collapsing all positions within an LD block and thus creating independent genomic loci ( $n = 1,703$ ). We use binary information for LD blocks, i.e. blocks with at least one significant association are considered as a hit, and the rest are not.

#### **Analysis of mediated pleiotropy between diseases**

Using the LCV method developed by O'Connor & Price, we tested the causal relationships between diseases<sup>20</sup>. We used the R functions developed by the authors and provided on GitHub ([github.com/lukejoconnor/LCV/](https://github.com/lukejoconnor/LCV/)). We calculated the genetic causality proportion (GCP) between each disease pair, if the diseases have at least 10 significant variants and a significant heritability estimate as suggested by the developers ( $Z_h \geq 7$ ). We only calculated GCP if the diseases are not vertically connected on the disease hierarchy. Following the criteria applied by the developers, we considered pairs with FDR corrected  $p \leq 0.01$  and mean  $GCP > 0.6$  as significant.

#### **SNP to gene mapping**

We map all SNPs analyzed in GWAS to genes based on proximity and eQTL results.

#### Using proximity

Using VariantAnnotation<sup>21</sup>, TxDb.Hsapiens.UCSC.hg19.knownGene<sup>22</sup>, and GenomicRanges<sup>23</sup> packages in R, we mapped the genomic coordinates for each SNP to genes. Specifically, if a gene is within the coding region, intron, 5' or 3' UTR, or 1kb down- or up-stream of the transcription start site, we annotated that SNP to the gene. As a result, we had 4,443,872 SNP-gene associations for 4,236,176 SNPs and 22,228 Entrez gene IDs. We used the Ensembl biomaRt<sup>24</sup> package in R to retrieve HGNC symbols (17,994), Ensembl Gene IDs (20,507), and gene descriptions for the Entrez gene IDs obtained from TxDb.Hsapiens.UCSC.hg19.knownGene database.

#### Using GTEx eQTL data

Using SNP-gene associations based on GTEx v7 eQTL data (accessed on 04.09.2018)<sup>25</sup>, we associated SNPs with the genes they could potentially regulate. We generated a combined tissue list, which associates SNP to the gene if there is at least one tissue in which there is a significant ( $p \leq 5e-8$ ) association. As a result, there are 2,166,300 unique SNPs associated with 15,312 Ensembl Gene IDs. We used the biomaRt<sup>24</sup> package in R to retrieve HGNC Symbols (12,292), Entrez IDs (10,163), and gene descriptions.

#### Comparison of proximity and eQTL based mapping

Instead of only focusing on disease-associated SNPs, we first mapped all SNPs that we analyzed to discover if there is a bias for certain genes (e.g. some genes could have many more SNPs because they are longer, or because they are already associated with certain

traits and the chip is designed in that way). There were as much as 19,195 SNPs mapped to one gene (CSMD1) by proximity, whereas there were 82 SNPs per gene on average (median). The number of SNPs per gene was on average, higher for the mappings by eQTL (Figure S33a). The maximum was 8473 SNPs for HLA-C gene and the median number of SNPs per gene was 218. However, we did not consider MHC region in our downstream analysis and thus this region is also excluded. The correlation between the number of SNPs per gene was low ( $\rho = 0.13$ , Figure S33b). Since the proximity-based mapping is by definition dependent on the gene length, we also tested if there is a significant correlation between the number of SNPs per gene and gene length. While the correlation is low for gene mappings by eQTL (Spearman's correlation  $\rho = 0.03$ ,  $p = 1.073e-4$ ), mappings by proximity show a high correlation as expected (Spearman's correlation  $\rho = 0.87$ ,  $p < 2.2e-16$ ). This also explains the low correlation between eQTL and proximity-based mappings. We next checked the correlation between the number of SNPs per gene mapped by proximity but only to promoter region. The correlation between the number of SNPs and gene length decreased ( $\rho = 0.21$ ), and the correlation with the number of SNPs by eQTL slightly increased but was still low ( $\rho = 0.08$ ). Overall, we concluded that both eQTL data and proximity-based mapping could capture different information and decided to use both for the downstream analyses.

### **GWAS Catalog analysis**

We accessed the GWAS Catalog on 30-07-2019 and used v1.0.2 e96 dataset<sup>26</sup>. We excluded all studies which used UK Biobank dataset ( $n = 190$ , data courtesy of GWAS Catalog team). Using the associations with a p-value lower than  $5 \times 10^{-8}$ , we compiled significant associations between MAPPED\_GENEs and MAPPED\_TRAITs. We use GWAS catalog analysis to check if our GWAS hits are supported by previous studies and applied a Fisher test between all traits in GWAS catalog and the diseases in our study. P-values are corrected for multiple testing using FDR correction.

### **Analysis of the association with aging**

We downloaded GenAge human, GenAge model organism<sup>27</sup> and DrugAge<sup>28</sup> data on Aug 13, 2019 and CellAge<sup>29</sup> data on Oct 02, 2019 (CellAge data is kindly provided by Avelar et al.). We used HGNC Symbols for GenAge and CellAge genes. In order to compile genes that are targeted by the drugs in DrugAge database, using the drug names in DrugAge data, we first compiled PubChem IDs using PubChem REST API<sup>30</sup>. Using UniChem<sup>31</sup>, we mapped PubChem IDs to ChEMBL IDs<sup>32</sup>. Next, using DGIdb<sup>33</sup>, we compiled the genes targeted by these ChEMBL IDs. As a result, we had 307 genes from GenAge human database, 902 genes from GenAge model organism database, 279 genes from CellAge database, and 714 genes targeted by DrugAge drugs. We next calculated the overlaps between these databases and the genes associated with multiple diseases or multiple categories in different age-of-onset clusters. To calculate the expected values and statistical significance, we used 10,000 permutations calculating the overlap for the same number of random genes among genes that can be detected by GWAS. Then, an odds ratio is calculated by dividing the observed value to the mean of expected values.

### **Functional Enrichment Test**

Using the goseq package in R<sup>34</sup>, which takes the gene length bias into account, we performed a functional analysis of the genes associated with different age-of-onset clusters. Using GO categories with more than 10 and less than 500 annotated genes, we applied an enrichment test for the Gene Ontology (GO)<sup>35,36</sup> Biological Process (BP), Molecular Function (MF), and Cellular Compartment (CC) categories. BY correction<sup>37</sup> is applied to the p-values for all tests for all clusters and 3 GO Categories (BP, MF, and CC) combined. We considered associations with a BY-corrected p-value lower than 0.05 as significant. For the ease of visualization and comprehension we selected representative categories for significant associations as follows: For each cluster and GO Ontology (*i.e.* BP, MF, CC) separately; i) Jaccard similarity index (*i.e.* number of genes in common divided by the number of unique genes in each category combined) is calculated between all significantly associated GO Categories; ii) Jaccard indices

are hierarchically clustered and cut to k number of groups, where k is the minimum number of clusters which ensures median Jaccard similarity within a cluster is above 0.5; iii) The category with the highest average similarity to other categories in the same cluster is assigned as the representative.

### **Drug Repurposing**

We searched for the drugs that specifically target multicategory genes in cluster 1, cluster 2, or cluster 1 and 2. Using the Fisher's exact test, we compiled the drugs in DGIdb<sup>33</sup> that specifically target these genes ( $p \leq 0.01$ ) and drugs that target only one gene in one of these clusters. Importantly, we excluded all non-specific drugs (*i.e.* targeting more than 10 genes) from the analyses. The interaction data is compiled from DGIdb, and the names, indications and phases of the drugs are obtained from ChEMBL REST API<sup>32</sup>.

### **Evolutionary Analysis**

In order to test the mutation accumulation and antagonistic pleiotropy theories of aging we used the risk allele frequencies in UK Biobank and 1000 Genomes super-populations<sup>38</sup>. A risk allele is an allele that shows positive association with a disease. Since the SNPs are not independent and have similar allele frequencies in a given LD block, we analyzed LD blocks instead of individual SNPs and used the median risk allele frequency for a given LD block. We used only the biallelic SNPs for this analysis. Allele frequencies for UK Biobank are calculated using BOLT-LMM and the allele frequencies for 1000 Genome super-populations are obtained from the vcf file provided on the 1000 Genomes project website. To test the antagonistic pleiotropy excess, we calculated the proportion of antagonistic vs. agonist SNPs within the same vs. different age-of-onset clusters using Fisher's exact test. We considered pleiotropic SNPs as agonist if the risk allele for two or more diseases are the same, and antagonist if the risk alleles are opposite. We only tested the risk allele frequency differences between cluster 1 and cluster 2. Also, we excluded any SNPs that are antagonistic within an age-of-onset cluster and agonist between clusters.



## References

1. Anderson, C. A. *et al.* Data quality control in genetic case-control association studies. *Nat. Protoc.* **5**, 1564–1573 (2010).
2. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
3. Kanchi, K. L. *et al.* Integrated analysis of germline and somatic variants in ovarian cancer. *Nat. Commun.* **5**, 3156 (2014).
4. Khurana, E. *et al.* Role of non-coding sequence variants in cancer. *Nat. Rev. Genet.* **17**, 93–108 (2016).
5. Poduri, A., Evrony, G. D., Cai, X. & Walsh, C. A. Somatic mutation, genomic variation, and neurological disease. *Science* **341**, 1237758 (2013).
6. Zhang, L. & Vijg, J. Somatic Mutagenesis in Mammals and Its Implications for Human Disease and Aging. *Annu. Rev. Genet.* **52**, 397–419 (2018).
7. Finkel, T., Serrano, M. & Blasco, M. A. The common biology of cancer and ageing. *Nature* **448**, 767–774 (2007).
8. Sanchez-Valle, J. *et al.* Unveiling the molecular basis of disease co-occurrence: towards personalized comorbidity profiles. *bioRxiv* 431312 (2018) doi:10.1101/431312.
9. Gutiérrez-Sacristán, A. *et al.* comoRbidity: an R package for the systematic analysis of disease comorbidities. *Bioinformatics* **34**, 3228–3230 (2018).
10. Chouakria, A. D. & Nagabhushan, P. N. Adaptive dissimilarity index for measuring time series proximity. *Adv. Data Anal. Classif.* **1**, 5–21 (2007).
11. Montero, P. & Vilar, J. TSclust: An R Package for Time Series Clustering. *Journal of Statistical Software, Articles* **62**, 1–43 (2014).
12. Berndt, D. J. & Clifford, J. Using Dynamic Time Warping to Find Patterns in Time Series. in *Proceedings of the 3rd International Conference on Knowledge Discovery and Data Mining* 359–370 (AAAI Press, 1994).
13. Partitioning Around Medoids (Program PAM). in *Finding Groups in Data* (eds. Kaufman,

- L. & Rousseeuw, P. J.) 68–125 (John Wiley & Sons, Inc., 1990).
14. Maechler, M., Rousseeuw, P., Struyf, A., Hubert, M. & Hornik, K. *cluster: Cluster Analysis Basics and Extensions*. (2019).
  15. Tibshirani, R., Walther, G. & Hastie, T. Estimating the number of clusters in a data set via the gap statistic. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* vol. 63 411–423 (2001).
  16. Loh, P.-R. *et al.* Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat. Genet.* **47**, 284–290 (2015).
  17. Loh, P.-R. *BOLT-LMM v2. 3.1 User Manual*.  
<https://data.broadinstitute.org/alkesgroup/BOLT-LMM/> (2017).
  18. Stephenson, J. D., Laskowski, R. A., Nightingale, A., Hurles, M. E. & Thornton, J. M. VarMap: a web tool for mapping genomic coordinates to protein sequence and structure and retrieving protein structural annotations. *Bioinformatics* (2019)  
doi:10.1093/bioinformatics/btz482.
  19. Berisa, T. & Pickrell, J. K. Approximately independent linkage disequilibrium blocks in human populations. *Bioinformatics* **32**, 283–285 (2016).
  20. O'Connor, L. J. & Price, A. L. Distinguishing genetic correlation from causation across 52 diseases and complex traits. *Nat. Genet.* **50**, 1728–1734 (2018).
  21. Obenchain, V. *et al.* VariantAnnotation: a Bioconductor package for exploration and annotation of genetic variants. *Bioinformatics* **30**, 2076–2078 (2014).
  22. Carlson, M. & Maintainer, B. P. TxDb.Hsapiens.UCSC.hg19.knownGene: Annotation package for TxDb object(s). (2015).
  23. Lawrence, M. *et al.* Software for computing and annotating genomic ranges. *PLoS Comput. Biol.* **9**, e1003118 (2013).
  24. Durinck, S., Spellman, P. T., Birney, E. & Huber, W. Mapping identifiers for the integration of genomic datasets with the R/Bioconductor package biomaRt. *Nat. Protoc.* **4**, 1184–1191 (2009).
  25. Gamazon, E. R. *et al.* Using an atlas of gene regulation across 44 human tissues to

- inform complex disease- and trait-associated variation. *Nat. Genet.* **50**, 956–967 (2018).
26. Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* **47**, D1005–D1012 (2019).
  27. Tacutu, R. *et al.* Human Ageing Genomic Resources: new and updated databases. *Nucleic Acids Res.* **46**, D1083–D1090 (2018).
  28. Barardo, D. *et al.* The DrugAge database of aging-related drugs. *Aging Cell* **16**, 594–597 (2017).
  29. Avelar, R. A., Ortega, J. G., Tacutu, R., Tyler, E. & Bennett, D. A Multidimensional Systems Biology Analysis of Cellular Senescence in Ageing and Disease. *BioRxiv* (2019).
  30. Kim, S. *et al.* PubChem 2019 update: improved access to chemical data. *Nucleic Acids Res.* **47**, D1102–D1109 (2019).
  31. Chambers, J. *et al.* UniChem: a unified chemical structure cross-referencing and identifier tracking system. *J. Cheminform.* **5**, 3 (2013).
  32. Gaulton, A. *et al.* The ChEMBL database in 2017. *Nucleic Acids Res.* **45**, D945–D954 (2017).
  33. Cotto, K. C. *et al.* DGldb 3.0: a redesign and expansion of the drug–gene interaction database. *Nucleic Acids Res.* **46**, D1068–D1073 (2018).
  34. Young, M. D., Wakefield, M. J., Smyth, G. K. & Oshlack, A. Gene ontology analysis for RNA-seq: accounting for selection bias. *Genome Biol.* **11**, R14 (2010).
  35. The Gene Ontology Consortium & The Gene Ontology Consortium. The Gene Ontology Resource: 20 years and still GOing strong. *Nucleic Acids Research* vol. 47 D330–D338 (2019).
  36. Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat. Genet.* **25**, 25–29 (2000).
  37. Benjamini, Y. & Yekutieli, D. The control of the false discovery rate in multiple testing under dependency. *Ann. Stat.* **29**, 1165–1188 (2001).

38. 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).