Using Systems Genetics to Unravel the Genetics of Bone-related Traits

Basel Maher Al-Barghouthi

Ramallah, State of Palestine

Bachelor of Science, Biology, The University of Texas at Austin, 2013 Master of Science, Bioinformatics, University of Michigan, 2015 Master of Science, Biological and Physical Sciences, University of Virginia, 2017

A Dissertation presented to the Graduate Faculty of the University of Virginia in Candidacy for the Degree of Doctor of Philosophy

Department of Biochemistry and Molecular Genetics

University of Virginia

August 2021

Dr. Charles Farber

Dr. Stefan Bekiranov

Dr. Mete Civelek

Dr. Jason Papin

Dr. Stephen Rich

Dr. Jeffrey Saucerman

Abstract

Osteoporosis is a highly prevalent disease, characterized by reduced bone strength and an increased susceptibility to bone fractures, with over 10 million affected individuals in the U.S. alone. As populations age more successfully, the prevalence of osteoporosis is expected to rise; therefore, understanding the genetic basis of bone strength and related traits is of the utmost importance to the development of therapeutic interventions aimed at reducing the societal burden of osteoporosis. To this end, over the last decade, geneticists have performed genome-wide association studies (GWASs) of bone mineral density (BMD) in order to gain insight into the genetics of osteoporosis. These studies have been very successful, identifying over 1,100 independent associations to date. However, efforts to understand the genetics of bone and to discover actionable therapeutic targets have been limited due to two main shortcomings of BMD GWASs. First, GWASs in the bone field have almost exclusively focused on BMD as a trait. While BMD is a clinically relevant predictor of osteoporosis, it only explains part of the variance in bone strength. Second, progress has been limited due to the inherent difficulties in identifying the causal genes that underlie GWAS associations.

Here, we address these limitations by utilizing systems genetics approaches. Using a novel mouse population, the Diversity Outbred, we perform a GWAS of 55 bone traits and identify putatively causal genes underlying some of the GWAS associations. Furthermore, we utilize systems genetics approaches in order to inform existing BMD GWASs. The work presented in this dissertation provides a resource that will increase our understanding of the genetics of bone, and presents methodological techniques that are applicable across myriad complex traits.

Dedication

This dissertation is dedicated to my mother Reem, my father Maher, my brother Faris, and my sister Liane. Without your love and support, I wouldn't have been able to achieve this.

~Promise fulfilled~

Table of Contents

Abstract	II
Dedication	III
List of Figures	VII
Chapter 1: Introduction	1
1.1 The skeleton as a dynamic system	2
1.2 Current state of osteoporosis genetics	4
1.3 Using systems genetics to inform bone GWAS	6
1.3.1 Epigenetics-based approaches	7
1.3.2 Transcriptomics-based approaches	9
1.3.2.1 eQTL colocalization	9
1.3.2.2 Transcriptome-wide association studies	12
1.3.3 Network-based approaches	15
1.3.3.1 Biological networks	16
1.3.3.2 Co-expression networks	17
1.3.3.3 Bayesian networks	19
1.4 The mouse as a model for bone-related traits	22
1.5 Summary	25
Chapter 2: Systems Genetics Analyses in Diversity Outbred Mice Inform Identify Determinants of Bone Strength	n BMD GWAS and 28
2.1 Abstract	
2.2 Introduction	
2.3 Results	
2.3.1 Development of a resource for the systems genetic bone strength	s of 32
2.3.2 Identification of bone-associated nodes	

GWAS Assoc	iations	
3.1 Ab	ostract	
3.2 Int	troduction	
3.3 Re	sults	
	3.3.1 TWAS and eQTL colocalization identify potentially causal BMD GWAS genes	
	3.3.2 Characterization of genes identified by TWAS/eQTL colocalization	
	3.3.3 <i>PPP6R3</i> is a candidate causal gene for a GWAS association on Chr. 11	
	3.3.4 <i>PPP6R3</i> is a regulator of femoral geometry, BMD, and vertebral microarchitecture	
3.4 Di	scussion	
3.5 Me	ethods	
3.6 Da	nta availability	

3.7 Code availability	
3.8 Acknowledgements	109

Chapter 4: Concluding Remarks and Future Directions 11	<u>11</u>
4.1 Summary and conclusions11	12
4.1.1 Association mapping of bone-related traits11	12
4.1.2 Systems genetics analyses in mouse inform human BMD GWAS11	13
4.1.3 A combined TWAS/eQTL colocalization approach informs BMD GWAS11	15
4.2 Limitations11	16
4.2.1 Bone-specific –omics data11	16
4.2.2 Single-cell –omics data11	16
4.2.3 Beyond steady-state transcriptomics data	18
4.2.4 Non-European ancestries11	18
4.2.5 Investigating the genetics of other aspects of bone11	19
4.3 Future directions11	19
Appendix A: Supplementary Data 12	<u>22</u>
Appendix B: Supplemental Figures 12	<u>25</u>
References 14	<u>43</u>

List of Figures

Figure 1.1 Overview of systems genetics	3
Table 1.1 Examples of BMD GWAS	6
Figure 1.2 Systems genetics approaches for prioritizing GWAS data	8
Figure 1.3 Network analysis can reveal systems-level information	.21
Figure 2.1 Resource overview	.33
Figure 2.2 Characterization of the experimental Diversity Outbred cohort	.34
Figure 2.3 Overview of the network approach used to identify genes potentially responsible for BMD GWAS loci	.36
Figure 2.4 Identifying SERTAD4 and GLT8D2 as putative regulators of BMD	.42
Table 2.1 QTL identified for complex skeletal traits in the DO	.44
Figure 2.5 Overview of our approach to QTL fine-mapping	.46
Figure 2.6 QTL (locus1) on chromosome1	.48
Figure 2.7 Characterization of <i>Qsox1</i>	.50
Figure 2.8 Qsox1 is responsible for several chromosome 1 QTL	.53
Figure 3.1 TWAS and eQTL colocalization identify potentially causal BMD GWAS genes	.87
Table 3.1 Top 10 significant TWAS/eQTL genes	.88
Figure 3.2 TWAS and eQTL colocalization identify <i>GPATCH1</i> as a novel potentially causal BMD GWAS gene	.89
Table 3.2 Top 10 novel significant TWAS/eQTL genes (RCP ≥ 0.5)	.91
Figure 3.3 PPP6R3 is a top-10 novel eBMD gene	.94
Figure 3.4 Ppp6r3 functional validation shows an effect of genotype on bone mass	.96

Chapter 1

Introduction

Published in part in: Al-Barghouthi, B.M. and Farber, C.R. Dissecting the Genetics of Osteoporosis using Systems Approaches. *Trends in Genetics* (2019). doi:10.1016/j.tig.2018.10.004

1.1 The skeleton as a dynamic system

The human skeleton is a dynamic, adaptive, and complex system impacting a wide array of physiological processes. It provides support and protection, enables locomotion, maintains hematopoiesis, serves as a reservoir for calcium and phosphorus, and has important endocrine functions ^{1–3}. Diseases of bone, however, inhibit the ability of the skeleton to carry out these functions. The most common disease of bone is osteoporosis, a condition of low bone mineral density (BMD) and an increased risk of fracture ⁴. Osteoporosis affects over 12 million individuals in the U.S. and over 200 million worldwide⁵. Osteoporotic fractures are a serious clinical outcome associated with increased morbidity and mortality, particularly in the elderly. In fact, of the ~300,000 people in the U.S. over the age of 50 that suffer a hip fracture annually, 1 in 5 will die in the subsequent 12 months, and half of the survivors will not return to their prior independent living status ⁶. Alarmingly, the incidence of fractures is expected to rise by 50% over the next decade, as the number of individuals over the age of 50 increases ⁷.

One of the hallmarks of quantitative traits related to osteoporosis (BMD, bone size, etc.) is their high heritability ($h^2=0.5-0.8$)⁸. As a result, the development of a comprehensive understanding of bone biology necessitates a thorough understanding of the genetic factors underlying variation in bone traits. This not only includes defining the individual variants and genes contributing to osteoporosis, but also how they interact to impact molecular networks and systems-level function. Here, we discuss how systems genetics approaches are being used to accomplish these goals (**Figure 1.1**).



Figure 1.1 Overview of systems genetics. Genetic and environmental factors, as well as their interactions, influence complex bone traits, such as bone strength. The influence of these factors is mediated through impacts on molecular intermediates (transcriptomes, proteomes, metabolomes, etc.) and can be assayed via appropriate "omics" techniques. Biological information is propagated through complex molecular networks to affect bones traits and ultimately the risk of fracture.

1.2 Current state of osteoporosis genetics

The genetic analysis of osteoporosis began in the early 1990s with candidate gene studies describing associations between polymorphisms in bone-relevant genes (e.g., vitamin D receptor and type I collagen) and BMD ⁹. This was followed by a plethora of additional candidate gene investigations and linkage scans in families ¹⁰. In retrospect, little information was gained from either approach ^{11,12}. In 2007, the tide began to turn with the first of many genome-wide association studies (GWASs) of BMD¹³. In a BMD GWAS, the genotypes of millions of single nucleotide polymorphisms (SNPs) across the genome are tested for an association with BMD in thousands, now often hundreds of thousands, of individuals¹⁴. To date, over 20 primary GWAS and GWAS meta-analyses have identified hundreds of associations for BMD¹⁵⁻¹⁸. The largest GWAS to date analyzed estimated BMD (eBMD) at the heel in 426,824 individuals and identified 1,103 independent genome-wide significant associations in 518 loci (Table 1.1)¹⁹. BMD has been the primary target of GWASs, mainly because of its strong association with fracture, high heritability, and relative ease of assessment in very large cohorts²⁰. However, while BMD is the main focus of bone-related GWASs, it is important to note that BMD is not an all-encompassing bone phenotype. For example, studies have shown that only $\sim 60\%$ of the variance in bone strength, which is the most important determinant of fracture, is explained by BMD²¹. Concordant with bone strength being an emergent property of many bone traits, the remainder of the variance in bone strength can be explained by other bone traits. To date, bone traits such as bone size, bone geometry, and serum bone remodeling markers have been sparsely interrogated by GWAS ^{22–25}.

GWASs have reinforced the importance of known genes and pathways (RANK-RANKL, WNT signaling, etc.) in human bone biology. More importantly, GWASs have provided a treasure trove of loci containing only novel genes with the potential to revolutionize our understanding of the genetics, and more importantly, the biology of bone.

A limitation of current GWASs is that they have yet to fully uncover the genetic architecture of BMD ²⁶. In the heel eBMD study referenced above, the 1,103 independent associations explained only 20% of the phenotypic variance in eBMD ¹⁹. These data suggest that BMD is highly polygenic, or even omnigenic, and that much of the genetic basis of BMD remains to be discovered ^{27,28}. GWASs are ideally suited to identify associations with common variants (minor allele frequency (MAF) > 1%). Therefore, it is possible that rare variants (MAF < 1%) may explain part of the "missing heritability" ²⁹. In support of this hypothesis, recent whole genome-sequencing projects have identified rare variants with large effects on BMD ^{17,30–32}. It will likely require much larger GWASs and rare variant studies to fully dissect the genetic architecture of BMD and other bone traits.

By any standard, GWASs have been wildly successful at identifying new loci; however, to date this information has done little to increase our understanding of bone biology or disease. Of the hundreds of loci impacting BMD and other bone traits, the genes responsible for nearly all of the associations are unknown. There are many reasons for this knowledge gap, including the fact that most associations are due to non-coding variation, the lack of bone-specific "-omic" resources, and the inherent difficulties in experimentally establishing causality between variants, genes, and traits. **Table 1.1** Examples of large-scale GWAS and GWAS meta-analyses for BMD illustrate the increase in identified loci as a function of sample size. Note that the most recent discovery of 518 loci encompassed nearly all of the previously discovered loci identified in prior studies.

Study	Phenotype	Sample size	Association count
Morris et al. (2019) ¹⁹	Estimated heel BMD	426,824	1,103 independent associations (518 loci, 301 novel)
Kemp et al. (2017) ¹⁷	Estimated heel BMD	142,487	307 independent associations (203 loci, 153 novel)
Estrada et al. (2012) (meta- analysis) ¹⁶	Lumbar spine and femoral neck BMD	83,894 (32,961 discovery, 50,933 replication)	64 independent associations (56 loci, 32 novel)
Rivadeneira et al. (2009) ³³	Lumbar spine and femoral neck BMD	19,195	20 independent associations (20 loci, 13 novel)

1.3 Using systems genetics to inform bone GWAS

One approach that has the potential to increase our understanding of bone genetics is the emerging field of systems genetics^{34,35}. Systems genetics integrates the principles of systems biology with genetics to determine how genetic variation affects molecular phenotypes and cellular networks ³⁴. In the context of GWAS, systems genetics approaches have proven extremely useful for connecting associated variants with molecular functions (e.g., transcription). The "layering" of different "-omics" datasets (transcriptomics, metabolomics, proteomics, etc.) onto a set of GWAS loci is the most direct way to begin to identify the molecular consequences of disease-associated variants (**Figure 1.2**). Most importantly, it also serves to connect disease-associated variants to the genes they regulate. In this section, systems genetics approaches are discussed using three broad categories; namely epigenetic, transcriptomic and networks-based approaches.

1.3.1 Epigenetics-based approaches

An example of a systems genetics approach that has proven useful for informing GWAS is the integration of epigenetics data ³⁶. As mentioned above, the vast majority of BMD GWAS loci implicate only non-coding variation that presumably impacts gene regulation. Thus, it is likely that most causal GWAS variants reside in regulatory elements, such as promoters and enhancers, which can be identified as regions of open chromatin. In support of this hypothesis, studies have demonstrated an enrichment for GWAS variants overlapping enhancers in disease-relevant tissues ^{37–40}.

Epigenetic data have recently been used to inform BMD GWAS⁴¹. Using publicly available data, an eQTL in blood cells was identified for Long Intergenic Non-protein Coding RNA 339 (*LINC00339*) that colocalized with a BMD GWAS association on Chr. *1p36.12*. It was then found, using chromosome conformation capture (HI-C) data, that one of the eQTL SNPs (rs6426749) was located in a genomic region interacting with the promoter of *LINC00339*. Using epigenetics data from the ENCODE project, this SNP was found to overlap and influence the activity of an enhancer element in osteoblasts by altering a binding site for Transcription Factor AP-2 Alpha (*TFAP2A*)⁴².



Figure 1.2 Systems genetics approaches for prioritizing GWAS data. By integrating biological data with data from GWAS, SNPs affecting traits can be prioritized for functional follow-up. For example, transcriptomic, biophysical and epigenetic data pertaining to lists of GWAS SNPs can be leveraged in order to prioritize the most likely causal SNPs (red bar, bottom row). TFBS = transcription factor binding sites, H3K27ac = histone H3 lysine 27 acetylation.

Furthermore, alteration of *LINC00339* expression influenced the transcript levels of a nearby gene, Cell Division Control Protein 42 Homolog (*CDC42*), which plays a key role in bone modeling and remodeling ⁴³. Using a similar approach, another recent study determined that the BMD GWAS SNP rs9533090 affects the expression of Receptor Activator of Nuclear Factor kappa-B Ligand (*RANKL*), which plays a central role in osteoclastogenesis, by disrupting a nuclear factor 1 C-type (*NFIC*) binding site and enhancer activity ^{44,45}. These studies demonstrate the power of systems genetics approaches that combine multiple data types to unravel the molecular consequences of BMD-associated variants.

1.3.2 Transcriptomics-based approaches

1.3.2.1 eQTL colocalization

One of the most widely used systems genetics approaches for informing GWAS is the identification of expression quantitative trait loci (eQTL) ⁴⁶. Just like a clinical trait, GWAS can be used to identify associations for the expression of a gene ⁴⁷. These analyses identify sets of genetics variants, or eQTL, that influence transcript levels of any gene expressed in a given cell-type or tissue. There are two types of eQTL, local (*ais*) and distant (*trans*) ⁴⁸. Local eQTL influence the transcript levels of genes in close proximity; whereas distant eQTL influence gene expression over a long genomic distance. The identification of eQTL is a logical follow-up to a GWAS, given that the vast majority of GWAS loci are due to non-coding variants that presumably play a role in gene regulation.

Analyses utilizing eQTL have been greatly supported by recent efforts, such as the Genotype-Tissue Expression (GTEx) project, to provide reference datasets for gene expression across many human tissues and cell-types ^{49,50}. It is important to note that a major consideration for eQTL studies (and for that matter the generation of any other "-omics" dataset) is the cell-type or tissue used for the generation of gene expression profiles. Recently, the GTEx project demonstrated that many eQTL are tissue-specific, thus ideally

the transcriptomics data would be from a disease-relevant source ⁵¹. In the context of bone GWAS, this would be either bone tissue or bone cells, i.e., osteoblasts and osteoclasts ^{51,52}. However, the GTEx project does not include any bone-relevant data sources, and to date only three relatively small studies have generated bone relevant eQTL data. One such study generated microarray profiles on trans-iliacal bone biopsies from 84 postmenopausal women⁵³. These data were used to identify loci associated with lumbar spine BMD ²². Microarray profiles of undifferentiated osteoblasts from 95 individuals have also been used to identify eQTL and inform several bone GWASs ^{16,17,33,54}. More recently, eQTL were identified in cultured primary osteoclasts using RNA-seq profiles in 158 individuals ⁵⁵.

Recently, colocalization approaches (e.g., COLOC ⁵⁶, ENLOC ⁵⁷, eCAVIAR ⁵⁸) have been utilized in order to identify putatively causal genes underlying GWAS loci. A typical colocalization analysis consists of identifying local eQTL for genes located within a GWAS locus and then determining if the GWAS and the eQTL signals are due to the same sets of variants (referred to as colocalizing eQTL) ^{56,58}. A high probability of eQTL colocalization would then suggest that a genetic variant, or variants, are affecting the interrogated GWAS trait through the regulation of the expression of the eQTL gene. This approach has been successfully utilized across myriad complex traits including BMD. For example, the aforementioned osteoclast eQTL study used colocalization analysis to identify eight BMD loci with colocalizing eQTL ⁵⁵.

While many eQTL are tissue-specific, recent studies have also shown a high degree of shared eQTL across different tissues and cell-types, allowing for the informative utilization of eQTL data from non-trait relevant sources ^{50,51}. For example, a recent study colocalized eQTL from GTEx expression data in thyroid tissue with a BMD GWAS, in order to link the expression of Microtubule Affinity Regulating Kinase 3 (MARK3) to BMDassociated variants on Chr. 14q32.32 ^{59,60}. Evidently, eQTL colocalization approaches can be useful for the elucidation of informative biology, even when trait-relevant transcriptomic data are unavailable.

While colocalization analyses can certainly aid in the prioritization of genes underlying GWAS loci, it is important to note that identifying causal genes remains difficult. One major reason for this difficulty arises due linkage disequilibrium (LD). For example, it has been demonstrated that mismatching LD structures between GWASs and eQTL reference datasets (as is the norm) can lead to reduced power⁶¹. Furthermore, colocalizations in a locus can arise from several underlying relationships between variants, genes and traits: true causality (a variant affecting the GWAS trait by affecting the expression of a gene), linkage (variants independently affecting a GWAS trait and the expression of a gene) and pleiotropy (a variant independently affecting the GWAS trait and the expression of a gene), thereby increasing the uncertainty of the causal links between gene expression and phenotypes ⁶². Some methods have been developed to untangle these causal relationships, such as structural equation modeling and Summary-based Mendelian Randomization approaches ^{63–65}. However, these issues are further compounded by the use of summary statistics, which is the most frequently available modality of GWAS and eQTL data sources, rather than individual-level genotyping data. Another limitation of colocalization approaches is that the reliance of these approaches on GWAS associations means that variants with effect sizes that are too small to be identified by GWAS are not included in the analyses, leading to further loss of power.

While different colocalization approaches have attempted to tackle these issues by incorporating methods such as fine-mapping approaches (e.g., ENLOC) or utilizing statistical approaches to test different causal configurations (e.g., COLOC), there are no perfect methods that completely obviate these issues ^{56,57}. As is the prevailing trend in the genomics and systems genetics fields, these issues may not be resolved until larger, more ancestrally diverse GWAS and eQTL mapping studies, which also provide complete genotyping and LD structural data, are performed and made available.

1.3.2.2 Transcriptome-wide association studies

Another family of systems genetics approaches that utilize transcriptomic data to inform biology are transcriptome-wide association studies (TWAS) (e.g., FUSION ⁶⁶, PrediXcan ⁶⁷). As mentioned above, GWASs associate genomic variants with traits, and downstream analyses generally attempt to identify genes that are implicated by GWAS associations. Unlike GWAS, TWAS approaches provide insight into biology by directly associating gene expression with traits. Besides being more direct in associating genes with traits, TWAS approaches also benefit from reduced multiple testing burdens when compared to GWAS, as there are substantially fewer genes than there are genetic variants ⁶⁷. One issue with such an approach is the infeasibility of profiling gene expression across large cohorts (hundreds of thousands of individuals). To overcome this issue, TWAS approaches have relied on the use of gene expression reference datasets in order to impute (predict) gene expression in GWAS cohorts based on genotype. In other words, TWAS approaches can leverage both GWAS and eQTL reference data (e.g., GTEx) in order to impute gene expression in the GWAS cohort, thereby side-stepping the need to assay gene expression in prohibitively large cohorts. These imputed transcriptomes are then associated with biological traits. Furthermore, recent advances in the field have produced TWAS methods that apply this imputation approach by using GWAS summary statistics, thereby addressing the usual lack of individual-level genotyping data ^{66,68,69}.

Another important advance in the TWAS field is the recent generation of methods (e.g., MultiXcan⁶⁹) that can integrate imputed transcriptomes from multiple biological sources, thereby incorporating the largest components of gene expression variation across various tissues and cell-types. These approaches allow for the association of the joint effects of gene expression, from various biological sources, with biological traits. This is beneficial for multiple reasons; first, it is often more efficient to test the joint effects of gene expression variation, due to a decreased multiple testing burden 69. Second, this approach of combining multiple imputed transcriptomes performs better when there are multiple causal tissues that underlie a trait, and can be useful when expression data from causal tissues are unavailable ^{69,70}. TWAS approaches have recently been successfully utilized in the bone field in the absence of eQTL data from bone-relevant sources. For example, one study used TWAS to associate gene expression data from skeletal muscle and peripheral blood with femoral neck and lumbar spine BMD, thereby identifying 18 candidate BMD genes ⁷¹. While performing TWAS on individual non-bone tissues can be informative, approaches that integrate imputed transcriptomes from multiple sources will likely increase the power to prioritize bone-relevant genes ⁷⁰.

One important caveat regarding TWAS approaches is that TWAS rely on geneticallypredicted gene expression, as opposed to total expression. In addition to the statistical uncertainty inherent in predicting gene expression, total gene expression is also affected by non-genetic factors, such as environmental and technical factors. Therefore, TWAS approaches that rely on genetically-predicted gene expression will not reflect these factors ⁷⁰. Furthermore, TWAS approaches rely on data from common *cis*-eQTLs, which may only account for ~10% of the total genetic variance in gene expression ⁷².

TWAS and eQTL colocalization approaches are distinct, yet complementary, approaches. TWAS approaches test the correlation between genetically-imputed gene expression, while colocalization approaches calculate the probability that the same variant(s) influences both gene expression and disease risk (or a disease-associated quantitative trait). To this end, the use of both approaches simultaneously can improve the identification and prioritization of phenotype-relevant genes ⁷³. Such an approach has been recently performed in the PhenomeXcan study, where GWAS summary statistics on 4,091 traits were interrogated by both colocalization and TWAS ⁷³. While this study included eBMD and identified 76 protein-coding genes that met their significance thresholds, the GWAS summary statistics were not drawn from the largest available BMD GWAS. Furthermore, due to the breadth of their analysis, their results incurred a higher multiple-testing burden than if they had focused on a single GWAS. Another recent study performed a similar TWAS/colocalization approach to inform BMD GWAS; however, that study only used GTEx eQTL data from whole blood and skeletal muscle ⁷⁴.

In summary of this section, there are limitations to using eQTL data to inform BMD GWAS. As described above, the most powerful sets of eQTL data (e.g., GTEx) are from non-bone tissue. While such data have been informative for identifying colocalizing eQTL, it is likely that well-powered eQTL studies in bone tissue and bone cells will provide more insight. It has also become evident that tissue and cell-type specificity is a critical factor when

trying to dissect how GWAS loci influence BMD. As a result, not only do we need efforts focused on generating data in bone tissue and bone cells, but also specific bone cell populations at different stages of their lifecycle exposed to varying stimuli. It should also be noted that differences in the genetic backgrounds (with differences in LD structure) of GWASs and eQTL studies impact the accuracy and interpretability of results. This can be solved by efforts to generate both types of data from racially diverse populations.

Additionally, the use of eQTL data to inform GWAS is inherently focused on quantified gene expression, meaning that GWAS associations that may affect traits through non-eQTL related processes, such as exon splicing and protein abundance, may be missed.

1.3.3 Network-based approaches

GWASs for BMD have identified many genetic loci implicating disparate biological processes and mechanisms, suggesting a complex web of networks operating within and between various bone cell-types. Identifying these interactions is important as they can inform our understanding of "emergent properties" of bone that are not evident from the function of individual genes in isolation. This is analogous to identifying a car battery and alternator as elements involved in starting an engine. However, it would be impossible to understand their true function without knowing that they worked together in a car's electrical system. It is also likely that genetic variation is a major perturbation that shapes underlying biological networks. As a result, systems, rather than reductionist, approaches to bone genetics are critical to understand the role of genetics in systems-level function. Understanding bone molecular networks and how they are influenced by genetic variation is

also important in the context of discovering and evaluating potential anti-osteoporotic therapeutic targets ^{75–77}.

1.3.3.1 Biological networks

Networks are prevalent in all aspects of our lives. The internet, social media, and economic markets are all examples of networks that impact us daily. In biology, many types of networks exist including protein-protein interaction, transcriptions factor binding, metabolic, and gene regulatory networks. Mathematically, a network (or graph) is a set of nodes (elements) connected by edges, which represent relationships between nodes ⁷⁸. Edges can be directed or undirected and either weighted or unweighted. An undirected gene coexpression network represents the relationships in co-expression between genes without an indication of which node is upstream of the other, while a directed network models the information flow between nodes (e.g., increased expression of gene A causes increased expression of gene B). Weights can represent the strength of evidence for the edge or the strength of the relationship between nodes. Methods used to generate and analyze networks are indispensable to systems genetics, as they allow for a shift of focus from reductionist methods, like GWAS, to more holistic, systems-level approaches. Mostly due to the scarcity of bone-relevant data, and the relative paucity of investigators applying such approaches, the use of network biology in the bone field has lagged behind others. However, there are emerging use cases. For example, by combining BMD GWAS data with functional genomic analysis, a PU.1-dependent transcription factor network essential for osteoclast differentiation has been identified ⁷⁹.

1.3.3.2 Co-expression networks

The most popular types of biological networks used in systems genetics applications are based on co-expression. There are many methods for generating co-expression networks and one of the most widely used is weighted gene co-expression network analysis (WGCNA)⁸⁰. WGCNA organizes transcriptomic data into modules, or clusters, of coexpressed genes. It does this by analyzing co-expression (i.e., correlation in expression) across a set of perturbations, such as genetic background in mice or environmental exposures in a human population. Modules have been found to have a number of important features, such as containing functionally related genes that may be subject to co-regulation by similar factors ^{80,81}. As a result, one can think of co-expression network analysis as a way to organize biology in a relatively unbiased way, similar to the way that file folders are used to organize documents by topic.

There are two aspects of co-expression networks that make them particularly useful for systems genetics studies. First, unlike many other popular biological networks, co-expression networks retain tissue or cell-type specific information. While recent advances in proteomic technology have facilitated the study of protein-protein interactions *in vivo*, the vast majority of extant data is generated through *in vitro* methods, which may not accurately reflect physiological interactions ⁸². Second, unlike other biological networks, co-expression modules can be related to phenotypes from the individuals used to generate the transcriptomic profiles. For example, a WGCNA network was recently generated from blood cells in individuals with BMD measurements ⁸³. These data were then used to identify a module whose behavior (as summarized by its first principal component) was correlated with BMD. Once trait-correlated modules are identified they can be further analyzed to

identify key genes and relationships. For example, highly connected "hub" genes have been shown to drive modular associations with a trait ⁸⁴. A recent study generated a WGCNA network using bone transcriptomic data on 96 strains from the Hybrid Mouse Diversity Panel (HMDP)^{85,86}. An osteoblast-lineage specific module was identified (module 9) and shown to be highly correlated with femoral BMD in the same HMDP strains. The study showed that knockdown of the top two module 9 hub genes (Melanoma Antigen Family D1 (Maged1) and Par-6 Family Cell Polarity Regulator Gamma (Pard6g)) altered osteoblast proliferation, differentiation and mineralization in vitro and knockout of Maged1 decreased BMD in mice ^{84,86}. The authors mapped the first principal component of module 9 and demonstrated that the overall expression levels of module 9 genes were influenced by a local eQTL for Secreted Frizzled-related Protein 1 (Sfrp1), a key regulator of osteoblastogenesis⁸⁷. This demonstrates how co-expression network analysis in a genetics population can be used to understand the systems-level organization of genes. Similarly, another study generated a WGCNA network using gene expression data from female transiliac bone biopsies in humans. Through the integration of BMD GWAS data, this study identified a gene module and several candidate genes (Homer Protein Homolog 1 (HOMER1) and Spectrin Beta, Non-erythrocytic 1 (SPTBN1)), with putatively important roles in bone mass regulation⁸⁸.

Another use of co-expression networks is to inform GWAS. A number of studies have demonstrated that network information is a useful prioritization strategy for predicting causal genes for sets of GWAS associations⁸⁹. As an illustration, a recent study mapped genes located in 64 BMD GWAS associations onto the HMDP bone network described above ^{59,89}. This led to the identification of two modules that were enriched for genes implicated by GWAS. Using information on module genes with known roles in bone, it was predicted that novel module genes located in GWAS loci were causal and likely altered BMD via a role in osteoblasts. Two of the module genes, Microtubule Affinity Regulating Kinase 3 (*MARK3*) and Spectrin Beta, Non-erythrocytic 1 (*SPTBN1*), were experimentally confirmed to influence BMD when perturbed in mice. This study indicates that viewing GWAS data through the lens of a disease-relevant co-expression network can begin to highlight how key GWAS genes function together to regulate BMD.

1.3.3.3 Bayesian networks

Though initially described in the mid 1980's, Bayesian networks (BNs) have only recently begun to gain traction in biological research ⁹⁰. BNs are directed, acyclic graph representations of conditional dependencies between random variables 78. The directed, acyclic nature of the graphs is informative for reconstructing systems-level relationships between genes. For example, in a systems genetics context it is possible to apply a BN structure learning algorithm to a WGCNA module, as the dependence of gene expression on other genes can be observed in a hierarchical manner, which allows for an elucidation of the direction of the flow of molecular information. One scenario is where BN analysis methods are applied to trait-relevant WGCNA modules, in order to direct relationships between genes and identify key regulatory elements (Figure 1.3). This strategy was employed in a recent study, where an undirected co-expression network was constructed. Directional relationships between nodes were then established using Bayesian network analysis. This led to the identification of causal network structures relevant to late-onset Alzheimer's disease (LOAD) pathology as well as the identification of TYRO Protein Tyrosine Kinase Binding Protein (TYROBP) as a key regulator ⁹¹. In another study, a BN generated from coexpression modules was used to reveal regulatory driver genes affecting coronary artery

disease ⁹². Furthermore, this study used tissue-specific BNs to perform key driver analysis (KDA). Briefly, KDA leverages network topology and trait-or-disease-related gene sets to identify genes that are more connected with trait-or-disease-relevant genes than is expected by chance. These highly connected genes, termed "key drivers", are expected to be regulatory genes with strong evidence of having central roles in their associated networks. To our knowledge, BNs have not yet been applied in a systems genetics context in the bone field, and therefore provide an exciting avenue for future research.

One advantage of BNs is that they allow for the incorporation of prior knowledge, which allows for more informative modeling of gene relationships within modules. For example, network structure learning can be biased by "whitelisting" high-confidence edges (such as well-known gene-gene relationships or protein-protein interactions) a priori, or "blacklisting" improbable edges. Disparate data sources can be easily incorporated into BNs as well. For example, a BN from a WGCNA module can also include SNP nodes and trait nodes, in order to model information transfer from genetic element to gene expression and phenotypic outcomes ⁹³. Network-based approaches are not without limitations. One limitation involves the quality and type of the investigated bone phenotype. For example, BMD can be assayed in different anatomical locations by several different methods, which can lead to heterogeneity in the data that can obfuscate meaningful network relationships, or lead to network connections that are artificial and not mechanistically viable. Furthermore, a phenotype such as BMD is actually a composite of many different aspects of bone, which can also exacerbate the problems of interpretability. Therefore, careful selection of phenotypes should be performed *a priori*. Furthermore, biological networks often encompass multiple cell types, tissues and physiological microenvironments. In silico analyses based on data from *in vitro* sources, such as cultured osteoblasts, will not uncover many physiological

relationships that exist *in vivo*. This drawback is not unique to network analyses and pervades biological science, but should be carefully considered when designing experiments and drawing conclusions. Methodological drawbacks exist as well.



Figure 1.3 Network analysis can reveal systems-level information. In a typical network analysis workflow, gene expression profiles (**a**) (typically RNA-seq) can be analyzed from a network perspective. In this example, a co-expression network (**b**) is generated using WGCNA. In a WGCNA network, modules (colored clusters, **b**) consist of gene expression profiles connected by undirected edges, which signify the strength of the connection between two genes. A trait-correlated WGCNA module (within red circle, **b**) can be further dissected through Bayesian network analysis (**c**). Highly-connected hub genes (orange node, **d**) can signify functionally important genes. Bayesian networks differ from WGCNA networks in that they contain directed edges, and are acyclic (**c**). Furthermore, diverse biological information, such as SNP and trait data (blue and yellow boxes respectively, **e**) can be incorporated into Bayesian networks as prior information to improve network reconstruction. An advantage of Bayesian network analysis is the generation of more mechanistic hypotheses.

A significant drawback of using BNs to model biological relationships is in their acyclic nature. In biological processes, structures like feedback and feed-forward loops are prevalent. As BNs are acyclic, these network structures will be missed. Furthermore, depending on the algorithm used, it can be computationally impractical to learn the network structure of large sets of genes. These shortcomings make the aforementioned strategy of using BNs to dissect WGCNA modules an attractive one ⁹⁴.

1.4 The mouse as a model for bone-related traits

While the primary interest of biomedical research regarding bone is focused on understanding bone biology in humans, performing studies exclusively in humans is not feasible. For example, GWASs require very large sample sizes, and are generally not amenable to the interrogation of non-BMD traits in a well-powered fashion. Furthermore, ethical and practical considerations preclude most experiments in humans. To this end, the mouse has been an invaluable resource to the fields of bone biology and genetics.

The utility of the mouse as a model organism for the interrogation of bone biology lies in the high degree of similarity in both skeletal physiology and genetics between humans and mice, as well as in the relatively short reproductive cycle in mice, the relatively low cost of the mouse as an animal model, and the availability of the sequenced mouse genome ^{95–97}. These factors have enabled the use of the mouse to interrogate the genetic basis of many bone phenotypes, the identification of causal genes underlying bone traits, and the broad interrogation of the biology of bone. Classically, murine studies in bone have relied on the use of inbred mouse strains. After large differences in BMD were observed between various inbred mouse strains, researchers utilized mapping strategies in F₂ mice to identify bone-relevant quantitative trait loci (QTL) ⁹⁸. Briefly, inbred mouse strains with differing bone phenotypes (e.g., low vs. high BMD) are crossed to generate F₁ mice, which are then intercrossed to generate F₂ progeny ⁹⁹. These mice can then be used to identify broad QTL; for example, one study used F₂ mapping strategies to identify QTL for peak bone mass, by constructing intercrosses between SAMP6 (a murine model of senile osteoporosis with low peak bone mass) and SAMP2 ¹⁰⁰. While this strategy was informative in mapping QTL, these QTL were very large and were not amenable to the identification of causal genes ¹⁰¹.

In order to map narrower QTL regions that are more amenable to genetic analyses, several strategies were developed. One such strategy was to create inbred mouse panels for high-resolution association mapping, such as the Hybrid Mouse Diversity Panel (HMDP)⁸⁵. Using the HMDP, Farber et al. identified four significant associations affecting BMD, and identified Additional Sex Combs Like-2 (*Asx/2*) as a regulator of BMD and osteoclastogenesis¹⁰². More recently, newly available mouse reference panels and populations, such as the Collaborative Cross (CC) and the Diversity Outbred (DO) are enabling high-resolution mapping studies of many phenotypes, including bone microarchitectural phenotypes which are not amenable to interrogation in large human cohorts ¹⁰³. The CC is a recombinant inbred panel derived from eight inbred mouse strains, and the DO is an outbred mouse population derived from the same eight inbred mouse strains as the CC ^{104,105}. Recently, studies in the CC have successfully identified multiple QTL and candidate genes for several bone traits ^{103,106}.

Clearly, mouse models have been invaluable in furthering our understanding of the genetics of bone. One of the most important concepts that we've learned from murine studies of bone is that the adaptive nature of bone creates an additional layer of complexity in understanding osteoporosis. Studies have demonstrated that recombinant inbred mice with different genetic backgrounds build functional bones in different ways. For example, mice with genetically slender bones will compensate for this deficiency by increasing cortical thickness and mineralization, whereas mice with mineralization defects will increase bone size ^{107,108}. This genetically-based co-variation in traits serves as an example of a system adapting to perturbations. It also illustrates the importance of understanding not only how genetic variation impacts individual traits, but also the relationships between traits. A more encompassing approach to systems genetics has the potential to begin to understand how genetic variation contributes to these relationships and overall system function.

Finally, the use of mouse as a model for bone biology and genetics is made possible by the availability of several excellent resources, such as the International Mouse Phenotyping Consortium (IMPC) and the Origins of Bone and Cartilage Disease (OBCD) project, which systematically perturb mouse genes in order to screen the effects of genetic perturbations on myriad phenotypes, including several skeletal parameters ^{109,110}. Furthermore, resources such as the Mouse Genome Database (MGD) project provide integrated data such as integrated mouse strain phenotyping and genetic and genomic data¹¹¹.

1.5 Summary

In the past decade or so, advances in sequencing technologies have completely revolutionized biological science. In practically every biological field, a wealth of "-omics" data is being generated. However, our understanding of the underpinnings of biological processes and diseases is still far from complete. This is evident in the bone field, as many genetic associations with BMD have been described, yet we still know few of the responsible genes. These limitations reinforce the need for complementary strategies, such as systems genetics, to further advance our understanding of bone genetics.

One of the major limitations of genetic studies of bone is the primary focus on BMD. Although BMD is the single strongest predictor of osteoporotic fracture, there are many individuals with normal BMD who experience fracture ^{112,113}. The use of BMD has been necessitated by the difficulty, or impossibility, of measuring other aspects of bone fragility in humans. For example, biomechanical properties of bone strength, the single most important fracture-related trait, can only be measured in cadavers. Due to these limitations, a possible alternative is to use GWAS and systems genetics in mice and rats as a way of developing a more complete understanding of osteoporosis ^{102,114,115}.

In the field of systems genetics it is of the utmost importance to develop approaches for the effective understanding and utilization of available data. As biology is inherently complex, it is unreasonable to believe that a single, or few, types of genetic analyses will be sufficient to gain a thorough understanding of the genetics of complex bone traits. We argue that more realistic models of biological processes can be generated and analyzed by synthesizing and incorporating seemingly disparate data sources. For example, as mentioned above, integrating epigenetic data with genomic data can inform GWAS and identify candidate genes. Computational approaches such as TWAS and colocalization can utilize transcriptomic data to further inform the genetics of complex traits. Network-based approaches can inform GWAS and identify candidate genes while also providing insight into relationships between genes and traits. While barely scratching the surface, the systems genetics approaches described herein provide an avenue for such an endeavor. Of course, our understanding of many systems-level principles is still evolving. With increasingly accessible computational resources and by training researchers adept in the computational sciences, the transition to understanding bone biology and the impacts of genetic variation from a holistic, systems perspective will be within our reach.

In this work, we aim to address the two main limitations affecting the field of bone genetics; the strict focus of bone genetic studies on BMD, and the difficulties associated with the identification of causal genes underlying GWAS associations. We address these limitations in the following studies:

- (1) In Chapter 2, we use a novel mouse population, the Diversity Outbred (DO), in order to perform GWAS on 55 bone-related traits. We then use systems genetics approaches to identify Quiescin Sulfhydryl Oxidase 1 (*Qsox1*) as a novel gene affecting several bone traits, and perform experimental validation in a murine cohort. Furthermore, we use a network-based approach that utilizes expression data from the DO population, in order to prioritize putatively causal genes underlying human BMD GWAS associations.
- (2) In Chapter 3, we use a combined TWAS/eQTL colocalization approach, which leverages publicly available human gene expression and BMD GWAS data, in order to prioritize putatively causal genes underlying BMD GWAS associations.

Using this combined approach, we prioritize Protein Phosphatase 6 Regulatory Subunit 3 (*PPP6R3*) as a novel BMD gene, and perform functional validation in a murine cohort.

Overall, this work comprises a resource that aims to increase our understanding of the genetics underlying complex bone traits, and presents methodologies for the identification of putatively causal genes that underlie complex traits.

Chapter 2

Systems Genetics Analyses in Diversity Outbred Mice Inform BMD GWAS and Identify

Determinants of Bone Strength

Published in part in: Al-Barghouthi, B., Mesner, L., Calabrese, G., Brooks, D., Tommasini, S., Bouxsein, M., Horowitz, M., Rosen, C., Nguyen, K., Haddox, S., Farber, E., Onengut-Gumuscu, S., Pomp, D., and Farber, C. Systems genetics analyses in Diversity Outbred mice inform BMD GWAS and identify determinants of bone strength. *Nature Communications* 12, 3408 (2021). https://doi.org/10.1038/s41467-021-23649-0

2.1 Abstract

Genome-wide association studies (GWASs) for osteoporotic traits have identified over 1,100 associations; however, their impact has been limited by the difficulties of causal gene identification and a strict focus on bone mineral density (BMD). Here, we use Diversity Outbred (DO) mice to directly address these limitations by performing a systems genetics analysis of 55 complex skeletal phenotypes. We apply a network approach to cortical bone RNA-seq data to discover 66 genes likely to be causal for human BMD GWAS associations, including the genes *SERTAD4* and *GLT8D2*. We also perform GWAS in the DO for a wide-range of bone traits and identify *Qsax1* as a gene influencing cortical bone accrual and bone strength. In this work, we advance our understanding of the genetics of osteoporosis and highlight the ability of the mouse to inform human genetics.
2.2 Introduction

Osteoporosis is a condition of low bone strength and an increased risk of fracture ⁴. It is also one of the most prevalent diseases in the U.S., affecting over 10 million individuals⁷. Over the last decade, efforts to dissect the genetic basis of osteoporosis using genome-wide association studies (GWASs) of bone mineral density (BMD) have been tremendously successful, identifying over 1,100 independent associations ^{16,17,19}. These data have the potential to revolutionize our understanding of bone biology and the discovery of novel therapeutic targets ^{15,18}; however, progress to date has been limited.

One of the main limitations of human BMD GWAS is the difficulty in identifying causal genes. This is largely due to the fact that most associations implicate non-coding variation presumably influencing BMD by altering gene regulation ¹⁹. For other diseases, the use of molecular "-omics" data (e.g., transcriptomic, epigenomic, etc.) in conjunction with systems genetics approaches (e.g., identification of expression quantitative trait loci (eQTL) and network-based approaches) has successfully informed gene discovery ^{34,35}. However, few "-omics" datasets exist on bone or bone cells in large human cohorts (e.g., bone or bone cells were not part of the Geneotype-Tissue Expression (GTEx) project ⁵¹), limiting the use of systems genetics approaches to inform BMD GWAS ¹¹⁶.

A second limitation is that all large-scale GWASs have focused exclusively on $BMD^{16,17,19}$. BMD is a clinically relevant predictor of osteoporotic fracture; however, it explains only part of the variance in bone strength ^{117–120}. Imaging modalities and bone biopsies can be used to collect data on other bone traits such as trabecular microarchitecture and bone formation rates; however, it will be difficult to apply these techniques at scale (N=>100K). Additionally, many aspects of bone, including biomechanical properties,

cannot be measured *in vivo*. These limitations have hampered the dissection of the genetics of osteoporosis and highlight the need for resources and approaches that address the challenges faced by human studies.

The Diversity Outbred (DO) is a highly engineered mouse population derived from eight genetically diverse inbred founders (A/J, C57BL/6J, 129S1/SvImJ, NOD/ShiLtJ, NZO/HILtJ, CAST/EiJ, PWK/PhJ, and WSB/EiJ)¹⁰⁵. The DO has been randomly mated for over 30 generations and, as a result, it enables high-resolution genetic mapping and relatively efficient identification of causal genes ^{121,122}. As an outbred stock, the DO also more closely approximates the highly heterozygous genomes of a human population. These attributes, coupled with the ability to perform detailed and in-depth characterization of bone traits and generate molecular data on bone, position the DO as a platform to assist in addressing the limitations of human studies described above.

In this work, we present a resource for the systems genetics of bone strength consisting of information on 55 bone traits from over 600 DO mice, and RNA-seq data from marrow-depleted cortical bone in 192 DO mice. We demonstrate the utility of this resource in two ways. First, we apply a network approach to the bone transcriptomics data in the DO and identify 66 genes that are bone-associated nodes in Bayesian networks, and their human homologs are located in BMD GWAS loci and regulated by colocalizing eQTL in human tissues. Of the 66, 19 are not previously known to influence bone. The further investigation of two of the 19 novel genes, *SERTAD4* and *GLT8D2*, reveals that they are likely causal and influence BMD via a role in osteoblasts. Second, we perform GWASs in the DO for 55 complex traits associated with bone strength; identifying 28 QTL. By integrating QTL and bone eQTL data in the DO, we identify *Qsox1* as the gene responsible for a QTL

on Chromosome (Chr.) 1 influencing cortical bone accrual along the medial-lateral femoral axis and femoral strength. These data highlight the power of the DO mouse resource to complement and inform human genetic studies of osteoporosis.

2.3 Results

2.3.1 Development of a resource for the systems genetics of bone strength

An overview of the resource is presented in **Figure 2.1**. We measured 55 complex skeletal phenotypes in a cohort of DO mice (N=619; 314 males, 305 females; breeding generations 23-33) at 12 weeks of age. We also generated RNA-seq data from marrow-depleted femoral diaphyseal bone from a randomly chosen subset of the 619 phenotyped mice (N=192; 96/sex). All 619 mice were genotyped using the GigaMUGA¹²³ array (~110K SNPs) and these data were used to reconstruct the genome-wide haplotype structures of each mouse. As expected, the genomes of DO mice consisted of approximately 12.5% from each of the eight DO founders (**Figure 2.2A**). The collection of phenotypes included measures of bone morphology, microarchitecture, and biomechanics of the femur, along with tibial histomorphometry and marrow adiposity (**Supplementary Data 2.1 and 2.2**). Our data included quantification of femoral strength as well as many clinically relevant predictors of strength and fracture risk (e.g., trabecular and cortical microarchitecture). Traits in all categories (except tibial marrow adipose tissue (MAT)) were significantly (P_{adij}<0.05) correlated with femoral strength (**Supplementary Data 2.3**).





Additionally, all traits exhibited substantial variation across the DO cohort. For example, we observed a 30.8-fold variation (the highest measurement was 30.8 times greater than the lowest measurement) in trabecular bone volume fraction (BV/TV) of the distal femur and 5.6-fold variation in femoral strength (**Figure 2.2B**). After adjusting for covariates (age, DO generation, sex, and body weight) all traits had non-zero heritabilities (h²) (**Figure 2.2C**). Correlations between traits in the DO were consistent with expected relationships observed in previous mouse and human studies (**Supplementary Data 2.4**)^{124–}



Figure 2.2 Characterization of the experimental Diversity Outbred cohort. A) Allele frequency per chromosome, across the DO cohort. Intervals represent the eight DO founder strains: A/J (yellow), C57BL/6J (grey), 129S1/SvImJ (beige), NOD/ShiLtJ (dark blue), NZO/HILtJ (light blue), CAST/EiJ (green), PWK/PhJ (red), and WSB/EiJ (purple). B) Bone volume fraction and max load across the DO cohort. Insets are microCT images representing low and high bone volume fraction (BV/TV). C) Heritability of each bone trait. Phenotypes are colored by phenotypic category: morphology (purple), marrow adiposity (light blue), histomorphometry (green), microarchitecture (olive), and biomechanics (beige). Abbreviations for phenotypes are available in Supplementary Data 1.

In addition to standard RNA-seq quality control procedures (**Methods**), we also assessed RNA-seq quality by principal components analysis (PCA) and did not observe any major effect of sex, batch, and age in the first two principal components, which explained over 50% of the variance (**Supplemental Figure 2.1**). We did observe a separation of samples based on sex in the third PC, but it only explained 2.4% of the variance. Importantly, our PCA analysis did not identify any outliers in the bulk RNA-seq data. Furthermore, we performed differential expression analyses between sexes and between individuals with high versus low bone strength (**Supplementary Data 2.5 and 2.6**). As expected, the most significantly differentially expressed genes based on sex were located on the X chromosome. We identified 83 significantly (FDR<0.05) differentially-expressed transcripts in the analysis of low and high bone strength. Many were genes, such as *Absg*¹²⁸ and *Arg1*¹²⁹, which have previously been implicated in the regulation of bone traits.

2.3.2 Identification of bone-associated nodes

We wanted to address the challenge of identifying causal genes from BMD GWAS data, using the DO resource described above. To do so, we employed a network-based approach similar to one we have used in prior studies ^{59,130} (**Figure 2.3**). First, we partitioned genes into groups based on co-expression by applying weighted gene co-expression network analysis (WGCNA) to the DO cortical bone RNA-seq data ¹³¹. We generated three WGCNA networks; sex-combined, male, and female. The three networks contained a total of 124 modules (**Supplementary Data 2.7**).



Figure 2.3 Overview of the network approach used to identify genes potentially

responsible for BMD GWAS loci. Three WGCNA networks (124 total modules) were constructed from RNA-seq data on cortical bone in the DO (N=192). A Bayesian network was then learned for each module. We performed key driver analysis on each Bayesian network to identify BANs, by identifying nodes (genes) that were more connected to more known bone genes than was expected by chance. We colocalized GTEx human eQTL for each BAN with GWAS BMD SNPs to identify potentially causal genes at BMD GWAS loci. For the key driver analysis, the yellow node indicates the queried gene, red nodes indicate known bone genes, and grey nodes indicate non-bone genes. Abbreviations: DO – Diversity Outbred, WGCNA – weighted gene co-expression network analysis, BAN – bone-associated nodes, eQTL – expression quantitative trait loci, BMD – bone mineral density, GWAS – genome-wide association studies, SNP – single-nucleotide polymorphism. PPH4 – posterior probability of colocalization, hypothesis 4. A Gene Ontology (GO) analysis revealed that nearly all modules were enriched for genes involved in specific biological processes, including modules enriched for processes specific to bone cells (osteoblasts or osteoclasts) (**Supplementary Data 2.8**).

We next sought to infer causal interactions between genes in each module, and then use this information to identify genes likely involved in regulatory processes relevant to bone and the regulation of BMD. To do so, we generated Bayesian networks for each coexpression module, allowing us to model directed gene-gene relationships based on conditional independence. Bayesian networks allowed us to model causal links between coexpressed (and likely co-regulated) genes.

We hypothesized that key genes involved in bone regulatory processes would play central roles in bone networks and, thus, be more highly connected in the Bayesian networks. In order to test this hypothesis, we generated a list of genes implicated in processes known to impact bone or bone cells ("known bone gene" list (N=1,291); **Supplementary Data 2.9**; see **Methods**). The GWAS loci referenced in this study were enriched in human homologs of genes in the "known bone gene" list, relative to the set of protein-coding genes in the genome (OR=1.35, P=1.45⁻⁷). Across the three network sets (combined, male and female), we found that genes with putative roles in bone regulatory processes were more highly connected than all other genes (P=3.5 x 10^{-4} , P=1.7 x 10^{-2} , and P=2.9 x 10^{-5} for combined, male, and female network sets, respectively), indicating the structures of the Bayesian networks were not random with respect to connectivity.

To discover genes potentially responsible for GWAS associations, we identified bone-associated nodes (BANs). BANs were defined as genes connected in our Bayesian networks with more genes in the "known bone gene" list than would be expected by chance ^{92,132–134}. The analysis identified 1,370 genes with evidence ($P_{nominal} \le 0.05$) of being a BAN (i.e., sharing network connections with genes known to participate in a bone regulatory process) (Supplementary Data 2.10).

2.3.3 Using BANs to inform human BMD GWAS

We reasoned that the BAN list was enriched for causal BMD GWAS genes. In fact, of the 1,370 BANs, 1,173 had human homologs and 688 of those were within 1 Mbp of one of the 1,161 BMD GWAS lead SNPs identified in ¹⁶ and ¹⁹. This represents an enrichment of BANs within GWAS loci (\pm 1Mbp of GWAS SNP), relative to the number of protein-coding genes within GWAS loci (OR=1.26, P=9.49 x 10⁻⁵).

However, a gene being a BAN is likely not strong evidence, by itself, that a particular gene is causal for a BMD GWAS association. Therefore, to provide additional evidence connecting BMD-associated variants to the regulation of BANs, we identified local eQTL for each BAN homolog in 48 human non-bone samples using the Genotype-Tissue Expression (GTEx) project ^{51,135,136}. Our rationale for using GTEx was that while these data do not include information on bone tissues or bone cells, a high degree of local eQTL sharing has been observed between GTEx tissues ^{50,51}. This suggests that a colocalizing eQTL in a non-bone tissue may represent either a non-bone autonomous causal effect or may reflect the actions of a shared eQTL that is active in bone and shared across non-bone tissues. We then tested each eQTL for colocalization (i.e., probability that the eQTL and GWAS association share a common causal variant) with their respective BMD GWAS association ^{16,19}. Of the 688 BANs located in proximity of a BMD GWAS locus, 66 had colocalizing eQTL (PPH4≥0.75, **Supplementary Data 2.11**, see **Methods**) in at least one

GTEx tissue (**Supplementary Data 2.12**). Of these, 47 (71.2%) were putative regulators of bone traits (based on comparing to the known bone gene list (N=36) and a literature search for genes influencing bone cell function (N=11)), highlighting the ability of the approach to recover known biology. Based on overlap with the known bone gene list, this represents a highly significant enrichment of known bone genes in the list of BANs with colocalizing eQTL relative to the number of known bone genes in the list of GWAS-proximal BANs (OR=2.53, P=3.09 x 10⁻⁴). Our approach identified genes such as *SP7* (Osterix) ¹³⁷, *SOST* ^{138,139}, and *LRP5* ^{140–142}, which play central roles in osteoblast-mediated bone formation. Genes essential to osteoclast activity, such as *TNFSF11* (RANKL) ^{143–146}, *TNFRSF11A* (RANK) ^{147,148}, and *SLC4A2* ¹⁴⁹ were also identified. Nineteen (28.8%) genes were not previously implicated in the regulation of bone traits.

One of the advantages of the network approach is the ability to identify potentially causal genes and provide insight into how they may impact BMD based on their module memberships and network connections. For example, the cyan module in the female network (cyan_F) harbored many of the known BANs that influence BMD through a role in osteoclasts (the GO term "osteoclast differentiation" was highly enriched P=2.8 x 10⁻¹⁵ in the cyan_F module) (**Supplementary Data 2.8**). Three of the nineteen novel BANs with colocalizing eQTL (**Supplementary Data 2.12**), *ATP6V1A*, *PRKCH* and *AMZ1*, were members of the cyan module in the female network. Based on their cyan module memberships it is likely they play a role in osteoclasts. *ATP6V1A* is a subunit of the vacuolar ATPase V1 domain ¹⁵⁰. The vacuolar ATPase plays a central role in the ability of osteoclasts to acidify matrix and resorb bone, though *ATP6V1A* itself (which encodes an individual subunit) has not been directly connected to the regulation of BMD ¹⁵⁰. *PRKCH* encodes the eta isoform of protein kinase C and is highly expressed in osteoclasts ¹⁵¹. *AMZ1*

is a zinc metalloprotease and is relatively highly expressed in osteoclasts, and is highly expressed in macrophages, which are osteoclast precursors ¹⁵¹.

Next, we focused on two of the novel BANs with colocalizing eQTL, SERTAD4 (GTEx Adipose Subcutaneous; coloc PPH4=0.77; PPH4/PPH3=7.9) and GLT8D2 (GTEx Pituitary; coloc PPH4=0.88; PPH4/PPH3=13.4). Both genes were members of the royalblue module in the male network (royalblue_M). The function of SERTAD4 (SERTA domain-containing protein 4) is unclear, though proteins with SERTA domains have been linked to cell cycle progression and chromatin remodeling ¹⁵². GLT8D2 (glycosyltransferase 8 domain containing 2) is a glycosyltransferase linked to nonalcoholic fatty liver disease ¹⁵³. In the DO, the eigengene of the royalblue_M module was significantly correlated with several traits, including trabecular number (Tb.N; rho=-0.26; P= 9.5×10^{-3}) and separation (Tb.Sp; rho=0.27; P= 7.1×10^{-3}), among others (Supplementary Data 2.13). The royalblue M module was enriched for genes involved in processes relevant to osteoblasts such as "extracellular matrix" (P= 8.4×10^{-19}), "endochondral bone growth" (P= 5.7×10^{-4}), "ossification" (P= 8.9×10^{-4}) and "negative regulation of osteoblast differentiation" (P=0.04) (Supplementary Data 2.8). Additionally, Sertad4 and Glt8d2 were connected, in their local (3-step) Bayesian networks, to well-known regulators of osteoblast/osteocyte biology (such as Wnt16¹⁵⁴, Postn^{155,156}, and Col12a1¹⁵⁷ for Sertad4 and Pappa2¹⁵⁸, Pax1^{158,159}, and Tnn¹⁶⁰ for *Glt8d2*) (Figures 2.4A and 2.4B). *Sertad4* and *Glt8d2* were strongly expressed in calvarial osteoblasts with expression increasing (P<2.2 x 10^{-16} and P=6.4 x 10^{-10} , respectively) throughout the course of differentiation (Figure 2.4C). To further investigate their expression in osteoblasts, we generated single-cell RNA-seq (scRNA-seq) data on mouse bone marrow-derived stromal cells exposed to osteogenic differentiation media in vitro from our mouse cohort (N=5 mice (4 females, 1 male), 7,092 cells, Supplementary Data 2.14,

Supplemental Figure 2.2). Clusters of cell-types were grouped into mesenchymal progenitors, preadipocytes/adipocytes, osteoblasts, osteocytes, and non-osteogenic cells based on the expression of genes defining each cell-type (Supplementary Data 2.15). *Sertad4* was expressed across multiple cell-types, with its highest expression in a specific cluster (cluster 9) of mesenchymal progenitor cells and lower levels of expression in osteocytes (cluster 10) (Figure 2.4D, Supplemental Figure 2.3). *Glt8d2* was expressed in a relatively small number of cells in both progenitor and mature osteoblast populations (Figure 2.4D, Supplemental Figure 2.3).

Finally, we analyzed data from the International Mouse Phenotyping Consortium (IMPC) for Glt8d2 ¹⁶¹. After controlling for body weight, there was a significant (P=1.5 x 10⁻³) increase in BMD in male $Glt8d2^{-/-}$ and no effect (P=0.88) in female $Glt8d2^{-/-}$ mice (sex interaction P= 6.9 x 10⁻³) (Figure 2.4E). These data were consistent with the direction of effect predicted by the human GLT8D2 eQTL and eBMD GWAS locus where the effect allele of the lead eBMD SNP (rs2722176) was associated with increased GLT8D2 expression and decreased BMD. Together, these data suggest that *SERTAD4* and *GLT8D2* are causal for their respective BMD GWAS associations and they likely impact BMD through a role in modulating osteoblast-centric processes.



Figure 2.4 Identifying SERTAD4 and GLT8D2 as putative regulators of BMD. A) Local 3-step neighborhood around Sertad4. Known bone genes highlighted in green. Sertad4 highlighted in red. B) Local 3-step neighborhood around Glt8d2. Known bone genes highlighted in green. Glt8d2 highlighted in red. C) Expression of Sertad4 and Glt8d2 in calvarial osteoblasts. For each time point, N=3 independent biological replicates were examined. Error bars represent the standard error of the mean. TPM – transcripts per million. D) Single-cell RNA-seq expression data. Each point represents a cell (N=7,092 cells). The top panel shows UMAP clusters and their corresponding cell-type. The bottom two panels show the expression of Sertad4 and Glt8d2. The color scale indicates normalized gene expression value. E) Bone mineral density in Glt8d2 knockout mice from the IMPC. N=7 females and N=7 males for Glt8d2^{-/-} mice, N=1,466 females and N=1,477 males for Glt8d2^{+/+} mice. Boxplots indicate the median (middle line), the 25th and 75th percentiles (box) and the whiskers extend to 1.5 * IQR. Colors indicate genotype.

2.3.4 Identification of QTLs for strength-related traits in the DO

The other key limitation of human genetic studies of osteoporosis has been the strict focus on BMD, though many other aspects of bone influence its strength. To directly address this limitation using the DO, we performed GWAS for 55 complex skeletal traits. This analysis identified 28 genome-wide significant (permutation-derived P<0.05) QTLs for 20 traits mapping to 10 different loci (defined as QTL with peaks within a 1.5 Mbp interval) (**Table 2.1 and Supplemental Figure 2.4**). These data are presented interactively in a webbased tool (http://qtlviewer.uvadcos.io/). Of the 10 loci, four impacted a single trait (e.g., medial-lateral femoral width (ML) QTL on Chr2@145.4Mbp), while the other six impacted more than one trait (e.g., cortical bone morphology traits, cortical tissue mineral density (TMD), and cortical porosity (Ct.Por) QTL on Chr. 1@155Mbp). The 95% confidence intervals (CIs) for the 21 autosomal associations ranged from 615 Kbp to 5.4 Mbp with a median of 1.4 Mbp.

2.3.5 Overlap with human BMD GWAS

We anticipated that the genetic analysis of bone strength traits in DO mice would uncover novel biology not captured by human BMD GWAS. To evaluate this prediction, we identified overlaps between the 10 identified mouse loci and human BMD GWAS associations ^{3,19}. Of the 10 mouse loci, the human syntenic regions (**Supplementary Data 2.16**) for six (60%) contained at least one independent GWAS association (**Supplemental Figure 2.5**). We calculated the number expected by chance by randomly selecting 10 human regions (of the same size) 1000 times, followed by identifying overlaps. Six overlaps corresponded to the 57th percentile of the null distribution.

Locus	Trait	LOD	Chr.	Position (Mbp)	95% CI (Mbp)	# missense variants	Genes with colocalizing eQTL
1	TMD	23.9	1	155.1	154.8 - 155.6	7	Ier5
1	Ma.Ar	12.8	1	155.3	155.1 - 155.7	7	Ier5, Qsox1
1	Tt.Ar	11.5	1	155.2	155.1 - 156.2	7	Ier5, Qsox1
1	Ct.Por	11.4	1	155.4	155.1 - 156.4	7	Ier5, Qsox1
1	ML	10	1	155.4	155.1 - 155.7	7	Ier5, Qsox1
1	рМОІ	8.8	1	155.1	154.8 - 158.2	7	Ier5, Qsox1
1	Ct.Ar/Tt.Ar	8.5	1	155.3	154.3 - 155.7	7	Ier5, Qsox1
1	Imax	8.3	1	155.1	155.1 - 158.2	7	Ier5, Qsox1
2	ML	7.9	2	145.4	144.1 - 145.6	-	-
3	Ma.Ar	8.8	3	68.1	66.6 - 70	8	Mfsd1, Il12a, Gm17641, 1110032F04Rik
4	Ma.Ar	8	4	114.6	113 - 118.4	-	-
4	Tt.Ar	8.2	4	114.6	113.6 - 114.8	-	-
5	Ct.Ar/Tt.Ar	8.1	4	127.7	125.4 - 128.1	-	Csf3r, Gm12946, Clspn, Ncdn, Gm12941, Zmym6, Gm25600
6	BMD	7.8	8	103.5	102.7 - 104.4	-	-
7	TMD	14.6	10	23.5	23.1 - 24.6	-	-
7	W	13.6	10	24.3	23.5 - 24.6	-	-
7	Wpy	11.9	10	23.8	23.5 - 25.3	-	-
7	Dfx	10.7	10	23.7	23.3 - 24.6	-	-
7	DFmax	9.4	10	23.7	21.8 - 25.2	-	C920009B18Rik
8	Fmax	8.8	16	23.3	22.3 - 23.4	-	-
8	Ffx	8.2	16	23.1	22.6 - 23.4	-	-
9	Ct.Ar	13.5	Х	59.4	58.4 - 71.2	-	-
9	Imax	11	Х	59.5	58.4 - 69.6	-	-
9	рМОІ	10.4	Х	59.4	58.4 - 61.4	-	-
9	Imin	8.4	Х	59.5	57.3 - 61.2	-	Zic3
10	Ct.Th	9.9	Х	73.4	58.4 - 74.1	-	-
10	TbSp	8.6	Х	73.8	72.7 - 77.5	-	Pls3
10	Tb.N	7.9	Х	74	72.7 - 76.8	-	Fundc2, Cmc4, Pls3

Table 2.1 QTL identified for complex skeletal traits in the DO

2.3.6 Identification of potentially causal genes

For each locus, we defined the causal gene search space as the widest confidence interval given all QTL start and end positions ± 250 Kbp. We then used a previously described approach, merge analysis, to fine-map QTL and identify likely causal genes (**Figure 2.5**)¹⁶². Merge analysis was performed by imputing all known variants from the genome sequences of the eight founders onto haplotype reconstructions for each DO mouse, and then performing single variant association tests. We focused on variants in the top 15% of each merge analysis as those are most likely to be causal ¹⁶².

We next identified missense variants that were top merge analysis variants common to all QTL in a locus. We identified seven missense variants in locus 1, and eight missense variants in locus 3 (**Table 2.1**). Of the seven missense variants in locus 1, three (rs243472661, rs253446415, and rs33686629) were predicted to be deleterious by SIFT. They are all variants in the uncharacterized protein coding gene *BC034090*. In locus 3, three (rs250291032, rs215406048 and rs30914256) were predicted to be deleterious by SIFT (**Supplementary Data 2.17**). These variants were located in myeloid leukemia factor 1 (*Mlf1*), Iqcj and Schip1 fusion protein (*Iqschfp*), and Retinoic acid receptor responder 1 (*Rarres1*), respectively.

We next used the cortical bone RNA-seq data to map 10,399 local eQTL in our DO mouse cohort (**Supplementary Data 2.18**). Of these, 174 local eQTL regulated genes located within bone trait QTL. To identify colocalizing eQTL, we identified trait QTL/eQTL pairs whose top merge analysis variants overlapped. This analysis identified 18 genes with colocalizing eQTL in 6 QTL loci (**Table 2.1**).



Figure 2.5 Overview of our approach to QTL fine-mapping. A) Overview of merge analysis. LOD- logarithm of the odds, QTL – quantitative trait loci, DO – Diversity Outbred, SNP – single-nucleotide polymorphism, INDEL – insertion-deletion, SV – structural variant. **B**) Overview of merge analysis as performed for the identification of missense variants. **C**) Overview of merge analysis as performed for the identification of colocalizing trait QTL/ gene eQTL within a locus. The pink columns around the QTL in each association plot represent the QTL 95% confidence intervals. The yellow box in panel (**c**) represents the gene search space for a locus, defined as the region within \pm 250 Kbp around the outer boundaries of the 95% confidence intervals within a locus. eQTL – expression quantitative trait loci.

2.3.7 Characterization of a QTL on Chromosome 1 influencing bone morphology

Locus 1 (Chr. 1) influenced cortical bone morphology (medullary area (Ma.Ar), total cross sectional area (Tt.Ar), medial-lateral femoral width (ML), polar moment of inertia (pMOI), cortical bone area fraction (Ct.Ar/Tt.Ar), and maximum moment of inertia (I_{max})), tissue mineral density (TMD), and cortical porosity (Ct.Por) (**Figure 2.6A**). We focused on this locus due to its strong effect size and the identification of candidate genes (*Ier5, Qsox1,* and *BC034090*) (**Table 2.1**). Additionally, we had previously measured ML in an independent cohort of DO mice (N=577; 154 males/423 females) from earlier generations (generations G10 and G11) and a QTL scan of those data uncovered the presence of a similar QTL on Chr. 1 ¹⁶³ (**Supplemental Figure 2.6, Methods**). The identification of this locus across two different DO cohorts (which differed in generations, diets, and ages) provided robust replication justifying further analysis.

The traits mapping to this locus fell into two phenotypic groups, those influencing different aspects of cross-sectional size (e.g., ML and Tt.Ar) and TMD/cortical porosity. We suspected that locus 1 QTL underlying these two groups were distinct, and that QTL for traits within the same phenotypic group were linked. This hypothesis was further supported by the observation that correlations among the size traits were strong and cross-sectional size traits were not correlated with TMD or porosity **(Supplementary Data 2.4)**.

Therefore, we next tested if the locus affected all traits or was due to multiple linked QTL. The non-reference alleles of the top merge analysis variants for each QTL were private to WSB/EiJ. To test if these variants explained all QTL, we performed the same association scans for each trait, but included the genotype of the lead ML QTL variant (rs50769082; 155.46 Mbp; ML was used as a proxy for all the cortical morphology traits) as

an additive covariate. This led to the ablation of all QTL except for TMD which remained significant (**Figure 2.6B**).



Figure 2.6 *QTL* (*locus 1*) *on chromosome 1. A*) For each plot, the top panel shows allele effects for the DO founders for each of the 8 QTL (quantitative trait loci) across an interval on chromosome 1 (Mbp, colors correspond to the founder allele in the legend). Bottom panels show each respective QTL scan. The red horizontal lines represent LOD (logarithm of the odds) score thresholds (genome-wide $P \le 0.05$). *B*) QTL scans across the same interval as panel (A), after conditioning on rs50769082. C) QTL scans after conditioning on rs248974780. Phenotype abbreviations: TMD – tissue mineral density, ML – medial-lateral femoral width, pMOI – polar moment of intertia, Imax – maximum moment of inertia, Ct.Ar/Tt.Ar – bone area fraction, Tt.Ar – total area, Ma.Ar – medullary area, Ct.Por – cortical porosity.

We then repeated the analysis using the lead TMD QTL variant (rs248974780; 155.06 Mbp) as an additive covariate (**Figure 2.6C**). This led to the ablation of all QTLs. These results supported the presence of at least two loci both driven by WSB/EiJ alleles, one influencing cortical bone morphology and Ct.Por and the other influencing TMD, as well as possibly influencing cortical bone morphology and Ct.Por.

2.3.8 *Qsox1* is responsible for the effect of locus 1 on cortical bone morphology

Given the importance of bone morphology to strength, we sought to focus on identifying the gene(s) underlying locus 1 and impacting cortical bone morphology. We reevaluated candidate genes in light of the evidence for two distinct QTL. Immediate Early Response 5 (*Ier5*) and Quiescin Sulfhydryl Oxidase 1 (*Qsox1*) were identified as candidates based on the DO mouse eQTL analysis and *BC034090* as a candidate based on missense variants (**Table 2.1**). Interestingly, *Ier5* and *Qsox1* eQTL colocalized with all QTL, except the TMD QTL, where only *Ier5* colocalized, providing additional support for two distinct loci (**Table 2.1 and Figure 2.7A**). We cannot exclude the involvement of the missense variants in *BC034090*; however, without direct evidence that they impacted *BC034090* function, we put more emphasis on the eQTL. As a result, based on its colocalizing eQTL and known biological function (see below), we predicted that *Qsox1* was at least partially responsible for locus 1.

QSOX1 is the only known secreted catalyst of disulfide bond formation and a regulator of extracellular matrix integrity ¹⁶⁴. It has not been previously linked to skeletal development. We found that *Qsax1* was highly expressed in calvarial osteoblasts and its expression decreased (P=6.4 x 10^{-6}) during differentiation (**Figure 2.7B**).



Figure 2.7 Characterization of Qsox1. A) The top panel shows allele effects for the DO founders for Ier5 and Qsox1 expression an interval on chromosome 1 (Mbp, colors correspond to the founder allele in the legend). Yaxis units are best linear unbiased predictors (BLUPs). Bottom panels show each respective QTL scan. LOD (logarithm of the odds) score threshold for autosomal eQTL is 10.89 (alpha=0.05). B) Qsox1 expression in calvarial osteoblasts. For each time point, N=3 independent biological replicates were examined. Error bars represent the standard error of the mean. TPM – transcripts per million. C) Single-cell RNA-seq expression data. Each point represents a cell (N=7,092 cells). The top panel shows UMAP clusters and their corresponding celltype. The bottom panel shows the expression of Qsox1. The color scale indicates normalized gene expression value.

In scRNA-seq on bone marrow-derived stromal cells exposed to osteogenic differentiation media *in vitro*, we observed *Qsox1* expression in all osteogenic cells with its highest expression seen in a cluster of mesenchymal progenitors defined by genes involved in skeletal development such as *Grem2*, *Lmna*, and *Prrx2* (cluster 1) (**Supplementary Data 2.19 and Figure 2.7C**). Additionally, in the DO cortical bone RNA-seq data, *Qsox1* was highly co-expressed with many key regulators of skeletal development and osteoblast activity (e.g., *Runx2*; rho=0.48, P=<2.2 x 10^{-16} , *Lrp5*; rho=0.41, P=6.2 x 10^{-9}).

To directly test the role of *Qsox1*, we used CRISPR/Cas9 to generate *Qsox1* mutant mice. We generated five different mutant lines harboring unique mutations, including two 1bp frameshifts, a 171-bp in-frame deletion of the QSOX1 catalytic domain, and two large deletions (756 bp and 1347 bp) spanning most of the entire first exon of Osox1 (Figure 2.8A, Supplementary Data 2.20 and 2.21). All five mutations abolished QSOX1 activity in serum (Figure 2.8B). Given the uniform lack of QSOX1 activity, we combined phenotypic data from all lines to evaluate the effect of QSOX1 deficiency on bone. We hypothesized based on the genetic and DO mouse eQTL data, that QSOX1 deficiency would increase all traits mapping to locus 1, except TMD. Consistent with this prediction, ML was increased overall (P=1.8 $\times 10^{-9}$), and in male (P=5.6 $\times 10^{-7}$) and female (P=3.5 $\times 10^{-3}$) mice as a function of *Qsox1* mutant genotype (Figure 2.8C). Also consistent with the genetic data, we observed no difference in other gross morphological traits including anterior-posterior femoral width (AP) (P=0.31) (Figure 2.8D) and femoral length (FL) (P=0.64) (Figure 2.8E). We next focused on male $Qsox1^{+/+}$ and $Qsox1^{-/-}$ mice and used microCT to measure other bone parameters. We observed increased pMOI (P=0.02) (Figure 2.8F), Imax (P=0.009) (Figure 2.8G), and Ct.Ar/Tt.Ar (P=0.031) (Figure 2.8H). Total area (Tt.Ar) (Figure 2.8I) was increased, but the difference was only suggestive (P=0.08). Medullary area (Ma.Ar, P=0.93)

was not different (**Figure 2.8J**). We observed no change in TMD (P=0.40) (**Figure 2.8K**). We also observed no difference in cortical porosity (Ct.Por) (P=0.24) (**Figure 2.8L**).

Given the strength of locus 1 on bone morphology and its association with biomechanical strength, we were surprised the locus did not impact femoral strength. Typically, in four-point bending assays, the force is applied along the AP axis. We replicated this in femurs from $Qsox1^{+/+}$ and $Qsox1^{-/-}$ mice and saw no significant impact on strength (P=0.20) (**Figure 2.8M**). However, when we tested femurs by applying the force along the ML axis, we observed a significant increase in strength in $Qsox1^{-/-}$ femurs (P=1.0 x 10⁻³) (**Figure 2.8N**). Overall, these data demonstrate that absence of QSOX1 activity leads to increased cortical bone accrual specifically along the ML axis (**Figure 2.8O**).



Figure 2.8 Qsox1 is responsible for several chromosome 1 QTL. A) Representative image of the Qsox1 knockout mutations. B) QSOX1 activity assay in serum. Data is grouped by mouse genotype. Boxplots indicate the median (middle line), the 25th and 75th percentiles (box) and the whiskers extend to 1.5 * IQR. Colors indicate mutation type. C-E) Femoral morphology in Qsox1 mutant mice. F-L) microCT measurements of chromosome 1 QTL phenotypes in Qsox1 knockout mice. M) Bone strength (max load, F_{max}) in the AP orientation, measured via four-point bending. N) Bone strength (max load, F_{max}) in the ML orientation, measured via four-point bending. O) Representative microCT images of the effect of Qsox1 on bone size. In panels C-N, P-values above plots are ANOVA P-values for the genotype term, while P-values in the plots are contrast P-values, adjusted for multiple comparisons. The center points of the plots represent the least-squares mean, while the error bars represent the confidence intervals at a confidence level of 0.95. Abbreviations: ML – medial-lateral femoral width, AP – anterior-posterior femoral width, FL – femoral length, pMOI – polar moment of inertia, Imax – maximum moment of inertia, Ct.Ar/Tt.Ar – bone area fraction, Tt.Ar – total area, Ma.Ar – medullary area, TMD – tissue mineral density, Ct.Por – cortical porosity.

2.4 Discussion

Human GWASs for BMD have identified over 1,100 loci. However, progress in causal gene discovery has been slow and BMD explains only part of the variance in bone strength and the risk of fracture ¹⁸. The goal of this study was to demonstrate that systems genetics in DO mice can help address these limitations. Towards this goal, we used cortical bone RNA-seq data in the DO and a network-based approach to identify 66 genes likely causal for BMD GWAS loci. Nineteen of the 66 were novel. We provide further evidence supporting the causality of two of these genes, *SERTAD4* and *GLT8D2*. Furthermore, GWAS in the DO identified 28 QTLs for a wide-range of strength associated traits. From these data, *Qsox1* was identified as a genetic determinant of cortical bone mass and strength. These data highlight the power of systems genetics in the DO and demonstrate the utility of mouse genetics to inform human GWAS and bone biology.

To inform BMD GWAS, we generated Bayesian networks for cortical bone and used them to identify BANs. Our analysis was similar to key driver analyses ^{132–134} where the focus has often been on identifying genes with strong evidence (P_{adj} <0.05) of playing central roles in networks. In contrast, we used BAN analysis as a way to rank genes based on the likelihood ($P_{nominal} \leq 0.05$) that they are involved in a biological process important to bone (based on network connections to genes known to play a role in bone biology). We then identified genes most likely to be responsible for BMD GWAS associations by identifying BANs regulated by human eQTL that colocalize with BMD GWAS loci. Together, a gene being both a BAN in a GWAS locus and having a colocalizing eQTL is strong support of causality. This is supported by the observation that ~71% of the 66 BANs with colocalizing eQTL were putative regulators of bone traits, based on a literature review and overlap with the "known bone gene" list.

One advantage of our network approach was the ability to not only identify causal genes, but use network information to predict the cell-type through which these genes are likely acting. We demonstrate this idea by investigating the two novel BANs with colocalizing human eQTL from the royalblue_M module. The royalblue_M module was enriched in genes involved in bone formation and ossification, suggesting the module as a whole and its individual members were involved in osteoblast-driven processes. This prediction was supported by the role of genes in osteoblasts that were directly connected to Sertad4 and Glt8d2, the expression of the two genes in osteoblasts, and for Glt8d2, its regulation of BMD in vivo. Little is known regarding the specific biological processes that are likely impacted by Sertad4 and Glt8d2 in osteoblasts; however, it will be possible to utilize this information in future experiments designed to investigate their specific molecular functions. For example, Sertad4 was connected to Wnt16, Ror2, and Postn all of which play roles in various aspects of osteoblast/osteocyte function. Wnt signaling is a major driver of osteoblast-mediated bone formation and skeletal development ¹⁶⁵. Interestingly, Wnt16 and Ror2 play central roles in canonical (Wnt16) and non-canonical (Ror2 in the Wnt5a/Ror2 pathway) Wnt signaling ¹⁶⁶ and have been shown to physically interact in chondrocytes ¹⁶⁷. Postn has also been shown to influence Wnt signaling 167,168. These data suggest a possible role for Sertad4 in Wnt signaling.

Despite their clinical importance, we know little about the genetics of bone traits other than BMD. Here, we set out to address this knowledge gap. Using the DO, we identified 28 QTL for a wide-range of complex bone traits. The QTL were mapped at highresolution, most had 95% CIs < 1 Mbp¹²². This precision, coupled with merge and eQTL analyses in DO mice, allowed us to identify a small number of candidate genes for many loci. Overlap of existing human BMD GWAS association and mouse loci was no more than what would be expected by chance, suggesting that our approach has highlighted biological processes impacting bone that are independent of those with the largest effects on BMD. This new knowledge has the potential to lead to novel pathways which could be targeted therapeutically to increase bone strength. Future studies extending the work presented here will lead to the identification of additional genes and further our understanding of the genetics of a broad range of complex skeletal traits.

Using multiple approaches, we identified Qsox1 as responsible for at least part of the effect of the locus on Chr. 1 impacting bone morphology. We use the term "at least part" because it is clear that the Chr. 1 locus is complex. Using ML width as a proxy for all the bone morphology traits mapping to Chr. 1, the replacement of a single WSB/EiJ allele was associated with an increase in ML of 0.064 mm. Based on this, if Qsox1 was fully responsible for the Chr. 1 locus we would expect at least an ML increase of 0.128 mm in Qsox1 knockout mice; however, the observed difference was 0.064 mm (50% of the expected difference). This could be due to differences in the effect of Qsox1 deletion in the DO compared to the SJL x B6 background of the Qsox1 knockout or to additional QTL in the Chr. 1 locus. The latter is supported by our identification of at least two QTL in the region. Further work will be needed to fully dissect this locus.

Disulfide bonds are critical to the structure and function of numerous proteins ¹⁶⁹. Most disulfide bonds are formed in the endoplasmic reticulum ¹⁷⁰; however, the discovery of QSOX1 demonstrated that disulfide bonds in proteins can be formed extracellularly ¹⁶⁴. Ilani et al. ¹⁶⁴ demonstrated that fibroblasts deficient in QSOX1 had a decrease in the number of disulfide bonds in matrix proteins. Moreover, the matrix formed by these cells was defective in supporting cell-matrix adhesion and lacked incorporation of the alpha-4 isoform of laminin. QSOX1 has also been associated with perturbation of the extracellular matrix in the context of cancer and tumor invasiveness ^{171,172}. It is unclear at this point how QSOX1 influences cortical bone mass; however, it likely involves modulation of the extracellular matrix.

In summary, we have used a systems genetics analysis in DO mice to inform human GWAS and identify genetic determinants for a wide-range of complex skeletal traits. Through the use of multiple synergistic approaches, we have expanded our understanding of the genetics of BMD and osteoporosis. This work has the potential to serve as a framework for how to use the DO, and other mouse genetic reference populations, to complement and inform human genetic studies of complex disease.

2.5 Methods

Diversity Outbred mouse population and tissue harvesting:

A total of 619 (315 males, 304 females) Diversity Outbred (J:DO, JAX stock #0039376) mice, across 11 generations (gens. 23-33) were procured from The Jackson Laboratory at 4 weeks of age. DO mice were fed standard chow (Envigo Teklad LM-485 irradiated mouse/rat sterilizable diet. Product # 7912). The mice were maintained on a 12hour light/12-hour dark cycle, at a temperature range of 60°C-76°C, with a humidity range of 20%-70%. Mice were injected with calcein (30 mg/g body weight) both 7 days and 1 day prior to sacrifice. Mice were weighed and fasted overnight prior to sacrifice. Mice were sacrificed at approximately 12 weeks of age (median: 86 days, range: 76-94 days). Immediately prior to sacrifice, mice were anesthetized with isoflurane, nose-anus length was recorded and blood collected via submandibular bleeding. At sacrifice, femoral morphology (length and width) was measured with digital calipers (Mitoyuto American, Aurora, IL). Right femora were wrapped in PBS soaked gauze and stored in PBS at -20°C. Right tibiae were stored in 70% EtOH at room temperature. Left femora were flushed of bone marrow (which was snap frozen and stored in liquid nitrogen, see below – Single-cell RNA-seq of bone marrow stromal cells exposed to osteogenic differentiation media in vitro) and were immediately homogenized in Trizol. Homogenates were stored at -80°C. Left tibiae were stored in 10% neutral buffered formalin at 4°C. Tail clips were collected and stored at -80°C.

Measurement of trabecular and cortical microarchitecture:

Right femora were scanned using a 10 μ m isotropic voxel size on a desktop μ CT40 (Scanco Medical AG, Brüttisellen, Switzerland), following the Journal of Bone and Mineral Research guidelines for assessment of bone microstructure in rodents ¹⁷³. Trabecular bone

architecture was analyzed in the endocortical region of the distal metaphysis. Variables computed for trabecular bone regions include: bone volume, BV/TV, trabecular number, thickness, separation, connectivity density and the structure model index, a measure of the plate versus rod-like nature of trabecular architecture. For cortical bone at the femoral midshaft, total cross-sectional area, cortical bone area, medullary area, cortical thickness, cortical porosity and area moments of inertia about principal axes were computed.

Biomechanical testing:

The right femur from each mouse was loaded to failure in four-point bending in the anterior to posterior direction, such that the posterior quadrant is subjected to tensile loads. The widths of the lower and upper supports of the four-point bending apparatus are 7 mm and 3 mm, respectively. Tests were conducted with a deflection rate of 0.05 mm/s using a servohydraulic materials test system (Instron Corp., Norwood, MA). The load and mid-span deflection were acquired directly at a sampling frequency of 200 Hz. Load-deflection curves were analyzed for strength (maximum load), stiffness (the slope of the initial portion of the curve), post-yield deflection, and total work. Post-yield deflection at yield. Yield is defined as a 10% reduction of stiffness relative to the initial (tangent) stiffness. Work, which is a measure of toughness, is defined as the area under the load-deflection curve. Femora were tested at room temperature and kept moist with phosphate buffered saline during all tests.

Assessment of bone marrow adipose tissue (MAT):

Fixed right tibiae, dissected free of soft tissues, were decalcified in EDTA for 20 days, changing the EDTA every 3-4 days and stained for lipid using a 1:1 mixture of 2% aqueous osmium tetroxide (OsO₄) and 5% potassium dichromate. Decalcified bones were

imaged using μ CT performed in water with energy of 55 kVp, an integration time of 500 ms, and a maximum isometric voxel size of 10 μ m (the "high" resolution setting with a 20mm sample holder) using a μ CT35 (Scanco). To determine the position of the MAT within the medullary canal and to determine its change in volume, the bone was overlaid. MAT was recorded in 4 dimensions.

Histomorphometry:

Fixed right tibiae were sequentially dehydrated and infiltrated in graded steps with methyl methacrylate. Blocks were faced and 5 µm non-decalcified sections cut and stained with toludine blue to observe gross histology. This staining allows for the observation of osteoblast and osteoclast numbers, amount of unmineralized osteoid and the presence of mineralized bone. Histomorphometric parameters were analyzed on a computerized tablet using Osteomeasure software (Osteometrics, Atlanta, GA). Histomorphometric measurements were made on a fixed region just below the growth plate corresponding to the primary spongiosa.

Bulk RNA isolation, sequencing and quantification:

We isolated RNA from a randomly chosen subset (n=192, 96/sex) of the available mice at the time (mice number 1-417), constrained to have an equal number of male and female mice. Total RNA was isolated from marrow-depleted homogenates of the left femora, using the mirVana[™] miRNA Isolation Kit (Life Technologies, Carlsbad, CA). Total RNA-Seq libraries were constructed using Illumina TruSeq Stranded Total RNA HT sample prep kits. Samples were sequenced to an average of 39 million 2 x 75 bp paired-end reads (total RNA-seq) on an Illumina NextSeq500 sequencer in the University of Virginia Center for Public Health Genomics Genome Sciences Laboratory (GSL). A custom bioinformatics pipeline was used to quantify RNA-seq data. Briefly, RNA-seq FASTQ files were quality controlled using FASTQC (version 0.11.5) ¹⁷⁴ and MultiQC (version 1.0.dev0) ¹⁷⁵, aligned to the mm10 genome assembly with HISAT2 (version 2.0.5) ¹⁷⁶, and quantified with Stringtie (version 1.3.3) ¹⁷⁷. Read count information was then extracted with a Python script provided by the Stringtie website (prepDE.py). Finally, we filtered our gene set to include genes that had more than 6 reads, and more than 0.1 transcripts per million (TPM), in more than 38 samples (20% of all samples). This filtration resulted in 23,648 genes remaining from an initial set of 53,801 genes. (Note that most of these genes were defined by StringTie internally as genes, but indicate loci – contiguous regions on the genome where the exons of transcripts overlap). Sequencing data is available on GEO at accession code GSE152708[https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE152708].

Bulk RNA differential expression analyses:

RNA-seq data were subjected to a variance stabilizing transformation using the DESeq2 (version 1.20.0) R package¹⁷⁸, and the 500 most variable genes were used to calculate the principal components using the PCA function from the FactoMineR (version 2.4) R package¹⁷⁹. For visualization, age was binarized into "high" and "low", with "low" defined as age equal to, or less than, the median age at sacrifice (85 days) and "high" defined as age higher than 85 days. Differential expression was then performed using DESeq2, for both sex and bone strength (max load). For differential expression based on sex, we used a design formula of ~batch+age+sex. For bone strength, we binarized bone strength into "high" and "low" for each sex independently, using the median bone strength value for each sex (35.66 and 37.42 for males and females, respectively). Differential expression was performed using the following design formula: ~sex+batch+age+bone strength. Log2 fold

changes for both differential expression analyses were then shrunken using the lfcShrink function in DESeq2, using the adaptive t prior shrinkage estimator from the apeglm (version 1.4.2) R package¹⁸⁰.

Mouse genotyping:

DNA was collected from mouse tails from all 619 DO mice, using the PureLink Genomic DNA mini kit (Invitrogen). DNA was used for genotyping with the GigaMUGA array ¹²³ by Neogen Genomics (GeneSeek; Lincoln, NE). Genotyping reports were preprocessed for use with the qtl2 (version 0.20) R package ^{181 182}, and genotypes were encoded using directions and scripts from (kbroman.org/qtl2/pages/prep_do_data.html). Quality control was performed using the Argyle (version 0.2.2) R package ¹⁸³, where samples were filtered to contain no more than 5% no calls and 50% heterozygous calls. Samples that failed QC were re-genotyped. Furthermore, genotyping markers were filtered to contain only tier 1 and tier 2 markers. Markers that did not uniquely map to the genome were also removed. Finally, a qualitative threshold for the maximum number of no calls and a minimum number of homozygous calls was used to filter markers.

We calculated genotype and allele probabilities, as well as kinship matrices using the qtl2 R package. Genotype probabilities were calculated using a hidden Markov model with an assumed genotyping error probability of 0.002, using the Carter-Falconer map function. Genotype probabilities were then reduced to allele probabilities, and allele probabilities were used to calculate kinship matrices, using the "leave one chromosome out" (LOCO) parameter. Kinship matrices were also calculated using the "overall" parameter for heritability calculations. Further quality control was then performed ¹⁸⁴, which led to the removal of several hundred more markers that had greater than 5% genotyping errors, after which genotype and allele probabilities and kinship matrices were recalculated. After the aforementioned successive marker filtration, 109,427 markers remained, out of 143,259 initial genotyping markers. As another metric for quality control, we calculated the frequencies of the eight founder genotypes of the DO.

WGCNA network construction:

Gene counts, as obtained above, were pruned to remove genes that had fewer than 10 reads in more than 90% of samples. Genes not located on the autosomes or X chromosome were also removed. This led to the retention of 23,335 out of 23,648 genes. Variance-stabilizing transformation (DeSeq2 ¹⁷⁸) was applied, followed by RNA-seq batch correction using sex and age at sacrifice in days as covariates (sex was not included as a covariate in the sex-specific networks), using ComBat (sva (version 3.30.0) R package ¹⁸⁵). We then used the WGCNA (version 1.68) R package to generate signed co-expression networks with a soft thresholding power of 4 (power=5 for male networks) ^{80,186}. We used the blockwiseModules function to construct networks with a merge cut height of 0.15 and minimum module size of 30. WGCNA networks had 39, 45 and 40 modules for the sex-combined, female and male networks, respectively.

Bayesian network learning:

Bayesian networks for each WGCNA module were learned with the bnlearn (version 4.5) R package ¹⁸⁷. Specifically, expression data for genes within a WGCNA module were obtained as above (WGCNA network construction), and these data were used to learn the

structure of the underlying Bayesian network using the Max-Min Hill Climbing algorithm (function mmhc in bnlearn).

Construction of the "known bone gene" list:

We constructed a list of bone genes using Gene Ontology (GO) terms and the Mouse Genome Informatics (MGI) database ¹⁸⁸ ¹⁸⁹. Using AmiGO2, we downloaded GO terms for "osteo*", "bone" and "ossif*", using all three GO domains (cellular component, biological process and molecular function), without consideration of GO evidence codes ¹⁹⁰. The resulting GO terms were pruned to remove some terms that were not related to bone function or regulation. We then used the MGI Human and Mouse Homology data table to convert human genes to their mouse homologs. We also downloaded human and mouse genes which had the terms "osteoporosis", "bone mineral density", "osteoblast", "osteoclast", and "osteocyte", from MGI's Human – Mouse: Disease Connection (HMDC) database. Human genes were converted to their mouse counterparts as above. GO and MGI derived genes were merged and duplicates were removed. Finally, we removed genes that were not expressed in our dataset. That is to say, they were not considered in generating the WGCNA modules or Bayesian networks.

Bone Associated Node (BAN) analysis:

We used a custom script that utilized the igraph (version 1.2.4.1) R package to perform BAN analysis ¹⁹¹. Briefly, within a Bayesian network underlying a WGCNA module, we counted the number of neighbors for each gene, based on a neighborhood step size of 3. Neighborhood sizes also included the gene itself. BANs were defined as genes that were more highly connected to bone genes than would be expected by chance. We merged all genes from all Bayesian networks together in a matrix, and removed genes that were unconnected or only connected to 1 neighbor (neighborhood size ≤ 2). We then pruned all genes whose neighborhood size was greater than 1 standard deviation less than the mean neighborhood size across all modules. These pruning steps resulted in 13,009/17,264, 11,861/16,446 and 11,877/17,042 genes remaining for the full, male and female Bayesian networks, respectively.

Then, for each gene, we calculated if they were more connected to bone genes in our bone list (see construction of bone list above) than expected by chance using the hypergeometric distribution (phyper, R stats (version 3.5.1) package). The arguments were as follows: q: (number of genes in neighborhood that are also bone genes) – 1; m: total number of bone genes in our bone gene set; n: (number of genes in networks prior to pruning) – m; k: neighborhood size of the respective gene; lower.tail = false.

GWAS-eQTL colocalization:

We converted mouse genes with evidence of being a BAN ($P \le 0.05$) to their human homologs using the MGI homolog data table. If the human homolog was within 1Mbp of a GWAS association, we obtained all eQTL associations within ± 200 kb of the GWAS association in all 48 tissues of version 7 of the Geneotype-Tissue Expression project (GTEx). These eQTL variants were colocalized with the GWAS variants, using the coloc.abf function from the R coloc (version 3.2.1) package ⁵⁶. This returned posterior probabilities (PP) for five hypotheses:

- H0: No association with either trait.
- H1: Association with trait 1, not with trait 2.
- H2: Association with trait 2, not with trait 1.
- H3: Association with traits 1 and 2, two independent SNPs.
- H4: Association with traits 1 and 2, one shared SNP.

Genes were considered colocalizing if PPH4 ≥ 0.75 .

Gene ontology:

Gene ontology analysis for WGCNA modules was performed for each individual module using the topGO (version 2.32.0) package in R¹⁹². Enrichment tests were performed for the "Molecular Function", "Biological Process" and "Cellular Component" ontologies, using all genes in the network. Enrichment was performed using the "classic" algorithm with Fisher's exact test. P-values were not corrected for multiple testing.

Assessing the expression of Glt8d2 and Sertad4 in publicly available bone cell data:

We used bioGPS expression data from GEO with the accession code of GSE10246[https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE10246] to assay the expression of *Sertad4*, *Glt8d2*, and *Qsox1* in osteoblasts ¹⁵¹. We also downloaded the data from GEO with the accession code

GSE54461[https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE54461] to query expression in primary calvarial osteoblasts.

Analysis of BMD data on $Glt8d2^{-}$ mice from the IMPC:

The International Mouse Knockout Consortium ¹⁶¹ and the IMPC¹⁹³ have generated and phenotyped mice harboring null alleles for *Glt8d2* (*Glt8d2*^{tm1a(KOMP)Wtsi}, *Glt8d2*^{-/-}) (N=7 females and N=7 males). Phenotypes for the appropriate controls (C57BL/6) were also collected (N=1,466 females and N=1,477 males). A description of the battery of phenotypes collected on mutants can be found at

(https://www.mousephenotype.org/impress/PipelineInfo?id=4). The mice were 14 weeks of age at DEXA scanning and both sexes were included. We downloaded raw BMD, body weight and metadata for *Glt8d2* mutants from the IMPC webportal

[https://www.mousephenotype.org/data/charts?accession=MGI:1922032&allele_accession _id=MGI:4364018&pipeline_stable_id=MGP_001&procedure_stable_id=IMPC_DXA_00 1¶meter_stable_id=IMPC_DXA_004_001&zygosity=homozygote&phenotyping_cente r=WTSI]. These data were analyzed using the PhenStat (version 2.18.1) R package ¹⁹⁴. PhenStat was developed to analyze data generated by the IMPC in which a large number of wild-type controls are phenotyped across a wide-time range in batches and experimental mutant animals are tested in small groups interspersed among wild-type batches. We used the Mixed Model framework in PhenStat to analyze BMD data. The mixed model framework starts with a full model (with fixed effects of genotype, sex, genotype x sex and weight and batch as a random effect) and ends with final reduced model and genotype effect evaluation procedures ^{194,195}.

QTL mapping:

Phenotypes that notably deviated from normality were \log_{10} -transformed (the MAT phenotypes as well as PYD and W_{py} were transformed after a constant of 1 was added). Then, QTL mapping with a single-QTL model was performed via a linear mixed model using the scan1 function of the qtl2 R package. A kinship matrix as calculated by the "leave one chromosome out" method was included. Mapping covariates were sex, age at sacrifice in days, bodyweight, and DO mouse generation. Peaks were then identified with a minimum LOD score of 4 and a peak drop of 1.5 LODs. To identify significant QTL peaks, we permuted each phenotype scan 1000 times (using the scan1perm function of the qtl2 package) with the same mapping covariates as above, and calculated the significance threshold for each phenotype at a 5% significance level. Heritability for the phenotypes was calculated using the est_herit function of the qtl2 R package, using the same covariates as above, but with a kinship matrix that was calculated using the "overall" argument.

DO eQTL mapping:

Variance stabilizing transformation was applied to gene read counts from above using the DESeq2 R package, followed by quantile-based inverse Normal transformation¹⁹⁶. Then, hidden determinants of gene expression were calculated from these transformed counts, using Probabilistic Estimation of Expression Residuals (PEER (version 1.3))¹⁹⁷. 48 PEER factors were calculated using no intercept or covariates. Sex and the 48 PEER covariates were used as mapping covariates, and eQTL mapping was performed using the scan1 function, as above. To calculate a LOD score threshold, we randomly chose 50 genes and permuted them 1000 times, as above. Since all genes were transformed to conform to the same distribution, we found that using 50 was sufficient. Thresholds were set as the highest permuted LOD score each for autosomal chromosomes and the X-chromosome (10.89 and 11.55 LODs, respectively). Finally, we identified peaks as above, and defined eQTL as peaks that exceeded the LOD threshold and were no more than 1Mbp away from their respective transcript's start site, as defined by the Stringtie output.

Merge analysis:

We performed merge analysis, a previously published approach, using the SNPassociation methods in the qtl2 R package ^{162,182}. For each DO mouse QTL or eQTL peak, we imputed all variants within the 95% confidence interval of a peak, and tested each variant for association with the respective trait. This was performed using the scan1snps function of the qtl2 R package, with the same mapping covariates for QTL or eQTL, respectively. Then, we identified "top" variants by taking variants that were within 85% of the maximum SNP association's LOD score. For conditional analyses using a variant, we performed the same QTL scan as above, but included the genotype of the respective SNP as an additive mapping covariate, encoding it as a 0, 0.5 or 1, for homozygous alternative, heterozygous or homozygous reference, respectively.

BMD-GWAS overlap:

To identify BMD GWAS loci that overlapped with our DO mouse associations, we defined a mouse association locus as the widest confidence interval given all QTL start and end CI positions mapping to each locus. We then used the UCSC liftOver tool (https://genome.ucsc.edu/cgi-bin/hgLiftOver) ¹⁹⁸ (minimum ratio of bases that must remap = 0.1, minimum hit size in query = 100000) to convert the loci from mm10 to their syntenic hg19 positions. We then took all genome-wide significant SNPs ($P \le 5 \ge 10^{-8}$) from the Morris et al. GWAS for eBMD and the Estrada et al. GWAS for FNBMD and LSBMD, and identified variants that overlapped with the syntenic mouse loci (GenomicRanges (version 1.32.7) R package ¹⁹⁹).

SIFT annotations for merge analysis missense variants were queried using Ensembl's Variant Effect Predictor tool (<u>https://useast.ensembl.org/Tools/VEP</u>)²⁰⁰. All options were left as default.

Prior ML QTL mapping:

The cohorts used for the earlier QTL mapping of ML consisted of 577 Diversity Outbred mice from breeding generations G10 and G11¹⁶³. G10 cohort mice consisted of both males and females fed a defined synthetic diet (D10001, Research Diets, New Brunswick, NJ), and were euthanized and analyzed at 12–15 weeks of age. G11 cohort mice were all females fed a defined synthetic diet (D10001, Research Diets, New Brunswick, NJ) until 6 weeks of age, and were then subsequently fed either a high-fat, cholesterolcontaining (HFC) diet (20% fat, 1.25% cholesterol, and 0.5% cholic acid) or a low-fat, high protein diet (5% fat and 20.3% protein) (D12109C and D12083101, respectively, Research Diets, New Brunswick, NJ), and were euthanized and analyzed at 24–25 weeks of age. Mice were weighed and then euthanized by CO₂ asphyxiation followed by cervical dislocation. Carcasses were frozen at -80°C. Subsequently, the femur was dissected and length, AP width, and ML width were measured two independent times to 0.01 mm using digital calipers. Mice were genotyped using the MegaMUGA SNP array (GeneSeek; Lincoln, NE) designed with 77,800 SNP markers, and QTL mapping was performed as above, but with the inclusion of sex, diet, age and weight at sacrifice as additive covariates.

Generation of Qsox1 mutant mice:

Qsox1 knockout mice used in this study were generated using the CRISPR/Cas9 genome editing technique essentially as reported in Mesner, et al²⁰¹. Briefly, Cas9 enzyme that was injected into B6SJLF2 embryos (described below) was purchased from (PNA Bio) while the guide RNA (sgRNA) was designed and synthesized as follows: the 20 nucleotide (nt) sequence that would be used to generate the sgRNA was chosen using the CRISPR design tool developed by the Zhang lab (crispr.mit.edu). The chosen sequence and its genome map position is homologous to a region in Exon 1 that is ~ 225 bp 3' of the translation start site and ~20bp 5' of the Exon1/Intron1 boundary (Supplementary Data **2.22**). To generate the sgRNA that would be used for injections oligonucleotides of the chosen sequence, as well as the reverse complement (Supplementary Data 2.22, primers 1 and 2, respectively), were synthesized such that an additional 4 nts (CACC and AAAC) were added at the 5' ends of the oligonucleotides for cloning purposes. These oligonucleotides were annealed to each other by combining equal molar amounts, heating to 90°C for 5 min. and allowing the mixture to passively cool to room temperature. The annealed oligonucleotides were combined with BbsI digested pX330 plasmid vector (provided by the Zhang lab through Addgene; https://www.addgene.org/) and T4 DNA ligase (NEB) and subsequently used to transform Stbl3 competent bacteria (Thermo Fisher) following the manufacturer's' protocols. Plasmid DNAs from selected clones were sequenced from primer 3 (Supplementary Data 2.22) and DNA that demonstrated accurate sequence and position of the guide were used for all downstream applications. The DNA template used in the synthesis of the sgRNA was the product of a PCR using the verified plasmid DNA and primers 4 and 5 (Supplementary Data 2.22). The sgRNA was synthesized via in vitro transcription (IVT) by way of the MAXIscript T7 kit (Thermo Fisher) following the

manufacturer's protocol. sgRNAs were purified and concentrated using the RNeasy Plus Micro kit (Qiagen) following the manufacturer's protocol.

B6SJLF1 female mice (Jackson Laboratory) were super-ovulated and mated with B6SJLF1 males. The females were sacrificed and the fertilized eggs (B6SJLF2 embryos) were isolated from the oviducts. The fertilized eggs were co-injected with the purified Cas9 enzyme (50 ng/ μ l) and sgRNA (30 ng/ μ l) under a Leica inverted microscope equipped with Leitz micromanipulators (Leica Microsystems). Injected eggs were incubated overnight in KSOM-AA medium (Millipore Sigma). Two-cell stage embryos were implanted on the following day into the oviducts of pseudopregnant ICR female mice (Envigo). Pups were initially screened by PCR of tail DNA using primers 6 and 7 with subsequent sequencing of the resultant product from primer 8, when the PCR products suggested a relatively large deletion had occurred in at least one of the alleles (Supplementary Data 2.22). For those samples which indicated a small or no deletion had occurred, PCR of tail DNA using primers 9 and 10 was performed with subsequent sequencing of the resultant products from primer 11 (Supplementary Data 2.22). Finally, deletions were fully characterized by ligating, with T4 DNA ligase (NEB), the PCR products from either primer pairs 6/7 or 9/10 with the plasmid vector pCR 2.1 (Thermo Fisher) followed by transformation of One Shot Top 10 chemically competent cells (Thermo Fisher) following the manufacturers recommendations (Supplementary Data 2.22).

The resulting founder mice (see **Supplementary Data 2.20**) were mated to C57BL/6J mice (Jackson Laboratory), with CRISPR/Cas9-deletion heterozygous F1 offspring from the 1st and 2nd litters mated to generate the F2 offspring used in the study of bone related properties reported herein. In addition, mouse B (**Supplementary Data 2.20**) was subsequently mated to an SJL/J male (Jackson Laboratory), and the F2 offspring from

the heterozygous F1 crosses, as outlined above, were also used in this study. All F1 and F2 mice from all deletion 'strains' were genotyped using primer pairs 9/10, with the PCR products sequenced from primer 11 for mice possessing the 7+6 and 1bp deletions (**Supplementary Data 2.22**). An additional PCR using primers 6 and 7 was performed with tail DNA from mice carrying the 1347 bp and 756 bp deletions; the products from this 2nd PCR assisted in determining between heterozygous and homozygous deleted genotypes (**Supplementary Data 2.22**).

ML was measured for both femurs using calipers on a population of 12-week old F2 mice and ML was averaged between the two femurs. A linear model with genotype, mutation type, length, and weight was generated separately for males and females. For the sex-combined data, a sex term was also included in the model. ANOVAs were performed using the Anova function from the car (version 3.0.7) R package ²⁰². Lsmeans were calculated using the emmeans (version 1.4.1) R package ²⁰³. The same procedure was performed for the AP and FL sex-combined data.

We randomly selected 50 male F2 mice (25 wt + 25 mut) from the same population, and microarchitectural phenotypes were measured as above, but on left femurs. Bone strength was measured as above but in both the AP and ML orientations. A linear model with genotype, mutation type and weight was generated, and lsmeans were calculated using the emmeans R package ²⁰³.

Measuring Qsox1 activity in serum:

Serum was collected via submandibular bleeding from isoflurane anesthetized mice, prior to sacrifice and isolation of femurs for bone trait analysis. Blood samples were incubated at room temperature for 20-30 m followed by centrifugation at 2000 xg for 10 m at 4°C. The supernatants were transferred to fresh tubes and centrifuged again as described above. The 2nd supernatant of each sample was separated into 50-100 µl aliquots, snap frozen on dry ice and stored at -70°C. Only 'clear' serum samples were used for determining QSOX1 activity, because pink-red colored samples had slight-moderate activity, presumably due to sulfhydryl oxidase enzymes released from lysed red blood cells.

Sulfhydryl oxidase activity was determined as outlined in Israel et al., 2014 ²⁰⁴ with minor modifications. Briefly, serum samples were thawed on wet ice whereupon 5 μ l was used in a 200 μ l final reaction volume which consisted of 50 mM KPO₄, pH7.5, 1mM EDTA (both from Sigma), 10 μ M Amplex UltraRed (Thermo Fisher), 0.5% (v/v) Tween 80 (Surfact-Amps, low peroxide; Thermo Fisher), 50 nM Horseradish Peroxidase (Sigma), and initiated with the addition of dithiothreitol (Sigma) to 50 μ M initial concentration. The reactions were monitored with the 'high-sensitive dsDNA channel' of a Qubit Fluorimeter (Thermo Fisher) by measuring the fluorescence every 15-30s for 10m. The assay was calibrated by adding varying concentrations (0-3.2 μ M) of freshly diluted H₂O₂ (Sigma) to the reaction mixture minus serum. Enzyme activity was expressed in units of (pmol H₂O₂/min/ μ l serum) and typically calculated within the first several minutes of the reaction for wild-type and heterozygous mutant mice. Enzyme activity was calculated during the entire 10 minutes of the reaction for homozygous mutant genotypes.

Single-cell RNA-seq of bone marrow stromal cells exposed to osteogenic differentiation media in vitro:

Bone marrow isolation:

The left femur was isolated and cleaned thoroughly of all muscle tissue followed by removal of its distal epiphysis. The marrow was exuded by centrifugation at 2000 xg for 30

seconds into a sterile tube containing 35 µl freezing media (90% FBS, 10% DMSO). The marrow was then triturated 6 times on ice after addition of 150 µl ice cold freezing media and again after further addition of 1ml ice cold freezing media until no visible clumps remained prior to being placed into a Mr. Frosty Freezing Container (Thermo Scientific) and stored overnight at -80° C. Samples were transferred the following day to liquid nitrogen for long term storage.

Bone marrow culturing:

Previously frozen bone marrow samples from 5 DO mice (mouse IDs: 12, 45, 48, 50, and 84) were thawed at 37°C, resuspended into 5 ml bone marrow growth media (Alpha MEM, 10% FBS, 1% Pen/Strep, 0.01% Glutamax), pelleted in a Sorvall tabletop centrifuge at 212 xg for 5 minutes at room temperature and then subjected to red blood cell lysis by resuspending and triturating the resultant pellet into 5 ml 0.2% NaCl for 20 seconds, followed by addition and thorough mixing of 1.6% NaCl. Cells were pelleted again, resuspended into 1 ml bone marrow growth media, plated into one well per sample of a 48 well tissue culture plate and placed into a 37° C, 5% CO₂ incubator undisturbed for 3 days post-plating, at which time the media was aspirated, cells were washed with 1 ml DPBS once and bone marrow growth media was replaced at 300 µl volume. The process was repeated through day 5 post-plating. At day 6 post-plating, cells were washed in same manner; however, we performed a standard *in vitro* osteoblast differentiation protocol, by replacing bone marrow growth media with 300 µl osteogenic differentiation media (Alpha MEM, 10% FBS, 1% Penicillin Streptomycin, 0.01% Glutamax, 50 mg/ml Ascorbic Acid, 1M Bglycerophosphate, 100µM Dexamethasome). Cells undergoing differentiation were assessed for accumulated mineralization on days 4, 6, 8 and 10 of the differentiation process as

follows: IRDye 680 BoneTag Optical Probe (Li-Cor Biosciences, product #926-09374) was reconstituted according to the manufacturer's instructions. On days 3, 5, 7 and 9, 0.006 nmoles were added to each sample. Twenty-four hours later the cells were washed with 0.5 mls DPBS (Gibco, product #14190250) and media was replaced. The cells were then placed on the Odyssey CLx Imaging System (Li-Cor Biosciences) to measure mineralization density as reflected by IRDye 680 BoneTag Optical Probe incorporation. Final values for mineralization were computed by subtracting the average number of fluorescent units recorded in designated background wells from the number of fluorescent units recorded in the sample wells.

RNA isolation:

The isolation procedure outlined below was inspired by ²⁰⁵. Mineralized cultures were washed twice with Dulbecco's Phosphate Buffered Saline (DPBS). 0.5ml 60mM EDTA (pH 7.4, made in DPBS) was added for 15-minute room temperature (RT) incubation. EDTA solution was aspirated and replaced for a second 15-minute RT incubation. Cultures were then washed with 0.5ml Hank's Balanced Salt Solution (HBSS) and incubated with 0.5ml 8mg/ml collagenase in HBSS/4mM CaCl₂ for 10 minutes at 37° C with shaking. Cultures were triturated 10x and incubated for an additional 20 minutes and 37° C. Cultures were then transferred to a 1.5ml Eppendorf tube, and spun at 500 xg for 5 minutes at RT in a Sorvall tabletop centrifuge. Cultures were resuspended in 0.5ml 0.25% trypsin-EDTA (Gibco, Gaithersburg, MD) and incubated for 15 minutes at 37° C. Cultures were then triturated and incubated for an additional 15 minutes. 0.5ml of media were added, triturated and spun at 500 xg for 5 minutes at RT. Cultures were then resuspended in 0.5ml bone marrow differentiation media and cells were counted.

Library preparation, sequencing and analysis:

The samples were pooled and concentrated to 800 cells/ μ l in sterile PBS supplemented with 0.1% BSA. The single cell suspension was loaded into a 10x Chromium Controller (10X Genomics, Pleasanton, CA, USA), aiming to capture 8,000 cells, with the Single Cell 3' v2 reagent kit, according to the manufacturer's protocol. Following GEM capturing and lysis, cDNA was amplified (13 cycles) and the manufacturer's protocol was followed to generate the sequencing library. The library was sequenced on the Illumina NextSeq500 and the raw sequencing data was processed using CellRanger toolkit (version 2.0.1). The reads were mapped to the mm10 mouse reference genome assembly using STAR (version 2.5.1b) ²⁰⁶. Overall, 7,188 cells were sequenced, to a mean depth of 57,717 reads per cell. Sequencing data is available on GEO at accession code

GSE152806[https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE152806].

Analysis was performed using Seurat (version 3.1.4) ^{207,208}. Features detected in at least 3 cells where at least 200 features were detected were used. We then filtered out cells with less than 800 reads and more than 5800 reads, as well as cells with 10% or more mitochondrial reads. This resulted in 7,105 remaining cells. Expression measurements were multiplied by 10,000 and log normalized, and the 3000 most variable features were identified. The data were then scaled. Cells were then scored by cell cycle markers, and these scores, as well as the percentage of mitochondrial reads, were regressed out ²⁰⁹. Finally, clusters were found with a resolution of 1 and the UMAP was generated. An outlier cluster consisting of 13 cells was removed, resulting in 7,092 remaining cells. Cluster cell types were manually annotated after performing differential expression analyses of the expression of genes in each cluster relative to all other clusters (**Supplementary Data 2.15**), using the Seurat FindAllMarkers function, with the only.pos=TRUE argument.

2.6 Data availability

Raw genotyping data, calculated genotype and allele probabilities, and R/qtl2 cross files are available from Zenodo at DOI:10.5281/zenodo.4265417 [https://zenodo.org/record/4265417] ²¹⁰. Raw sequencing data is available from the NCBI Gene Expression Omnibus database with accession codes GSE152708 [https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE152708] and GSE152806 [https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE152806] . Mapped DO mouse QTL and eQTL can be viewed at our web-based tool [http://qtlviewer.uvadcos.io/].

eBMD GWAS summary statistics used for this study are available from GEFOS [http://www.gefos.org/?q=content/data-release-2018], as are the FN and LS BMD GWAS summary statistics [http://www.gefos.org/?q=content/data-release-2012].

We used bioGPS expression data from GEO with the accession code of GSE10246[https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE10246] to assay the expression of *Sertad4*, *Glt8d2*, and *Qsox1* in osteoblasts. We also downloaded the data from GEO with the accession code

GSE54461[https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE54461] to query expression in primary calvarial osteoblasts. *Glt8d2* knockout data was downloaded from the IMPC

[https://www.mousephenotype.org/data/charts?accession=MGI:1922032&allele_accession _id=MGI:4364018&pipeline_stable_id=MGP_001&procedure_stable_id=IMPC_DXA_00 1¶meter_stable_id=IMPC_DXA_004_001&zygosity=homozygote&phenotyping_cente r=WTSI]. Mouse-Human homologs were obtained from MGI [http://www.informatics.jax.org/downloads/reports/HOM_MouseHumanSequence.rpt]. We also obtained data from the MGI Human-Mouse:Disease Connection database [http://www.informatics.jax.org/diseasePortal]. Gene Ontologies were obtained from AmiGO2 [http://amigo.geneontology.org/amigo].

Finally, we obtained expression data from version 7 of the Genotype-Tissue Expression project [https://gtexportal.org/home/datasets].

2.7 Code availability

Analysis code is available on GitHub [https://github.com/baselmaher/DO_project]²¹¹.

2.8 Acknowledgements

Research reported in this publication was supported in part by the National Institute of Arthritis and Musculoskeletal and Skin Diseases of the National Institutes of Health under Award Numbers AR057759 to C.J.R., M.C.H., and C.R.F., and AR077992 to C.R.F. B.M.A-B was supported in part by a National Institutes of Health, Biomedical Data Sciences Training Grant (5T32LM012416). The authors acknowledge Wenhao Xu (University of Virginia) and the Genetically Engineered Mouse Models (GEMM) core and the University of Virginia Cancer Center Support Grant (CCSG) P30CA044579 from NCI for their support in generating *Qsox1* mutant mice. The authors also acknowledge the Yale School of Medicine Department of Orthopaedics and Rehabilitation's Histology and Histomorphometry Laboratory for all their work. We thank Matt Vincent (The Jackson Laboratory) and Gary Churchill (The Jackson Laboratory) for developing the QTL Viewer software and Neal Magee (University of Virginia) for hosting QTL Viewer on UVA servers. We thank the IMPC for accessibility to BMD data on *Glt8d2* knockout mice (www.mousephenotype.org). The data used for the analyses described in this manuscript were obtained from the IMPC Portal on 11/5/19. The Genotype-Tissue Expression (GTEx) Project was supported by the Common Fund of the Office of the Director of the National Institutes of Health, and by NCI, NHGRI, NHLBI, NIDA, NIMH, and NINDS. The data used for the analyses described in this manuscript were obtained from the GTEx Portal on 01/15/18.

Chapter 3

Transcriptome-wide Association Study and eQTL Colocalization Identify Potentially Causal

Genes Responsible for Bone Mineral Density GWAS Associations

Basel M. Al-Barghouthi, Will Rosenow, Kang-Ping Du, Jinho Heo, Robert Maynard, Larry Mesner, Gina Calabrese, Aaron Nakasone, Bhavya Senwar, Louis Gerstenfeld, Virginia Ferguson, Cheryl Ackert-Bicknell, Elise Morgan, David L. Brautigan and Charles R. Farber

(In preparation)

3.1 Abstract

Genome-wide association studies (GWASs) for bone mineral density (BMD) have identified over 1,100 associations to date. However, identifying causal genes implicated by such studies has been challenging. Recent advances in the development of transcriptome reference datasets and computational approaches such as transcriptome-wide association studies (TWASs) and expression quantitative trait loci (eQTL) colocalization have proven to be informative in identifying putatively causal genes underlying GWAS associations. Here, we used TWAS/eQTL colocalization in conjunction with transcriptomic data from the Genotype-Tissue Expression (GTEx) project to identify potentially causal genes for the largest BMD GWAS performed to date. Using this approach, we identified 512 genes as significant (Bonferroni ≤ 0.05) using both TWAS and eQTL colocalization. This set of genes was enriched for regulators of BMD and members of bone relevant biological processes. To investigate the significance of our findings, we selected *PPP6R3*, the gene with the strongest support from our analysis which was not previously implicated in the regulation of BMD, for further investigation. We observed that *Ppp6r3* deletion in mice decreased BMD. In this work, we provide an updated resource of putatively causal BMD genes and demonstrate that PPP6R3 is a putatively causal BMD GWAS gene. These data increase our understanding of the genetics of BMD and provide further evidence for the utility of combined TWAS/colocalization approaches in untangling the genetics of complex traits.

3.2 Introduction

Osteoporosis is a disease characterized by low bone mineral density (BMD), decreased bone strength, and an increased risk of fracture that affects over 10 million individuals in the U.S. ^{4,7}. BMD is the single strongest predictor of fracture and a highly heritable quantitative trait ^{212 10 8}. Over the last decade, genome-wide association studies (GWASs) have identified over 1,100 independent associations for BMD ^{16,17,19}. However, despite the success of GWAS, few of the underlying causal genes have been identified ^{18,213}.

One of the main difficulties in GWAS gene discovery is that the vast majority (>90%) of associations are driven by non-coding variation ^{214,215}. Over the last decade, approaches such as transcriptome-wide association studies (TWASs) and expression quantitative trait locus (eQTL) colocalization, have been developed which leverage transcriptomic data in order to inform gene discovery by connecting non-coding disease associated variants to changes in transcript levels ^{56,57,66,68,73}. These approaches have proven successful for a wide array of diseases and disease-associated quantitative traits ^{73,216,217}. However, the osteoporosis field has lagged behind such efforts, due to the limited number of large-scale bone-related transcriptomic datasets.

In a TWAS, genetic predictors of gene expression (e.g., local eQTL - sets of genetic variants that influence the expression of a gene in close proximity ²¹⁸) identified in a reference population (e.g., the Genotype-Tissue Expression (GTEx) project ⁴⁹) are used to impute gene expression in a GWAS cohort. Components of gene expression due to genetic variation are then associated with a disease or disease-associated quantitative trait. Genes identified by TWAS are often located in GWAS associations, suggesting that the genetic regulation of their expression is the mechanism underlying such associations. Several tools

(e.g., FUSION, PrediXcan, MultiXcan^{66,67,69}) have been developed to perform TWASs. Most of these tools use GWAS summary statistics, making TWAS widely applicable to large GWAS datasets. In contrast, eQTL colocalization is a statistical approach that determines if there is a shared genetic basis for two associations (e.g., a local eQTL and BMD GWAS locus). Recently, it has been demonstrated that prioritizing genes using both TWAS and eQTL colocalization provides a way to identify genes with the strongest support for causality^{68,73}.

The GTEx project has generated RNA-seq data on over 50 tissues across hundreds of individuals ⁵⁰. Even though data on the tissues/cell-types likely to be most relevant to BMD (bone or bone cells) were not included, the project demonstrated that many expression quantitative trait loci (eQTL) were shared across tissues ^{50,51}. Additionally, it is well known that effects in a wide-array of non-bone cell-types and tissues can impact bone and BMD ^{219,220}. As a result, we sought to use the GTEx resource in conjunction with TWAS and eQTL colocalization to leverage non-bone gene expression data to identify putatively causal genes underlying BMD GWAS.

Here, we performed TWAS and eQTL colocalization using the GTEx resource and the largest BMD GWAS performed to date to identify potentially causal genes ¹⁹. Using this approach we identified 512 genes significantly associated via TWAS with a significant colocalizing eQTL. To investigate the significance of our findings we selected Protein Phosphatase 6 Regulatory Subunit 3 (*PPP6R3*), the gene with the strongest support not previously implicated in the regulation of BMD, for further investigation. We demonstrate using mutant mice that *Ppp6r3* is a regulator of lumbar spine BMD. These results highlight the power of leveraging GTEx data, even in the absence of data from the most relevant tissue/cell-types, to increase our understanding of the genetic architecture of BMD.

3.3 Results

3.3.1 TWAS and eQTL colocalization identify potentially causal BMD GWAS genes

To identify potentially causal genes responsible for BMD GWAS associations, we combined TWAS and eQTL colocalization using GTEx data (**Figure 3.1A**). We began by performing a TWAS using reference gene expression predictions from GTEx (Version 8; 49 tissues) and the largest GWAS performed to date for heel estimated BMD (eBMD) (>1,100 independent associations)^{19,50}. The analysis was performed using S-MultiXcan, which allowed us to leverage information across all 49 GTEx tissues ⁶⁹. Our analysis focused on protein-coding genes (excluded non-coding genes). A total of 2,156 protein-coding genes were significantly (Bonferroni-adjusted P-value ≤ 0.05) associated with eBMD

(Supplementary Data 3.1).

Next, we identified colocalizing eQTL from each of the 49 tissues in GTEx using fastENLOC ^{57,73}. We identified 1,182 colocalizing protein-coding genes with a regional colocalization probability (RCP) of 0.1 or greater (**Supplementary Data 3.2**). In total, 512 protein-coding genes were significant in both the TWAS and eQTL colocalization analyses (**Table 3.1 and Supplementary Data 3.3**). Among the identified genes were many with well-known roles in the regulation of BMD, such as *RUNX2* (**Figure 3.1B**), *IGF1*, and *LRP6*, as well as novel genes such as *RERE* (**Figure 3.1C**). Overall, the identified genes had significantly colocalizing eQTL in all 49 GTEx tissues, with eQTL from cultured fibroblasts

(132 genes), subcutaneous adipose tissue (117 genes), tibial artery (115 genes) and tibial nerve (114 genes) exhibiting the highest number of significant colocalizations

(Supplementary Data 3.4). TWAS predictors were only generated for genes on autosomes and of the 1,103 independent associations identified by Morris, et al.¹⁹, 1,097 were autosomal. For each of these, we defined a locus as the region consisting of \pm 1 Mbp around each association. Of the 1,097 loci, almost half (542; 49%) of the loci contained at least one of the 512 prioritized genes. Most loci overlapped one gene (mean = 1.7, median = 1); however, 184 loci overlapped multiple genes, including a locus on Chromosome (Chr) 20 (lead SNP rs6142137) which contained 9 prioritized genes. (Figure 3.1D).

3.3.2 Characterization of genes identified by TWAS/eQTL colocalization

To evaluate the ability of the combined TWAS/colocalization approach to identify genes previously implicated in the regulation of BMD, other bone traits, or the activity of bone cells, we queried the presence of "known bone genes" within the list of the 512 prioritized protein-coding genes. To do so, we created a database-curated set of genes previously implicated in the regulation of bone processes (henceforth referred to as our "known bone genes" list, N=1,399, **Supplementary Data 3.5**). Of the 512 genes identified above, 66 (12.9%) were known bone genes, representing a significant enrichment (odds ratio = 1.72; P = 1.0×10^{-4}) over what would be expected by chance (**Supplementary Data 3.6**).

We also performed a Gene Ontology enrichment analysis of the 512 prioritized genes. We observed enrichments in several bone-relevant ontologies, such as "regulation of ossification" (P=2.6 x 10^{-5}), "skeletal system development" (P=2.8 x 10^{-5}), and "regulation of osteoblast differentiation" (P=8.7 x 10^{-5}) (Figure 3.2A, Supplementary Data 3.7).



Figure 3.1 *TWAS* and *eQTL* colocalization identify potentially causal BMD GWAS genes. *A)* Overview of the analysis. The human image was obtained from BioRender.com. *TWAS*/colocalization plot for genes in the locus around RUNX2 (B) and RERE (C). The $-log_{10}$ Bonferroni-adjusted P-values from the *TWAS* analysis (top panel) and the maximum RCPs from the colocalization analyses (bottom panel). Genes alternate in color for visual clarity. Triangles represent RUNX2 (B) and RERE (C). D) Distribution of prioritized genes within eBMD GWAS loci.

Gene	Tissue with greatest RCP	Max. RCP	TWAS P-value (Bonferroni)
SPTBN1	Cells_Cultured_Fibroblasts	0.9469	<5 x 10 ⁻³²⁴
CCDC170	Spleen	0.6582	<5 x 10 ⁻³²⁴
FAM3C	Artery_Tibial	0.4917	<5 x 10 ⁻³²⁴
SEPT5	Skin_Sun_Exposed	0.4868	2.26 x 10 ⁻²⁸⁶
FGFRL1	Cells_Cultured_Fibroblasts	0.1611	5.31 x 10 ⁻²⁷²
GREM2	Cells_Cultured_Fibroblasts	0.9998	4.31 x 10-257
GPATCH1	Whole_Blood	0.3564	3.44 x 10 ⁻²²⁶
RHPN2	Pituitary	0.2181	8.71 x 10 ⁻²²¹
BMP4	Brain_Cortex	0.5468	5.49 x 10 ⁻¹⁶⁹
RUNX2	Esophagus_Gastroesophageal_Junction	0.2372	1.99 x 10 ⁻¹⁴⁶

Table 3.1 Top 10 protein-coding genes significant by colocalization (RCP ≥ 0.1) and TWAS, sorted by TWAS P-value.

The International Mouse Phenotype Consortium (IMPC) has recently measured whole body BMD on hundreds of mouse knockouts ^{110,193}. We searched the IMPC database for any of the 512 genes identified by TWAS and eQTL colocalization. Of the 512, 142 (27.7%) had been tested by the IMPC and 64 (12.5% of the 512 prioritized genes, 45% of the 142 IMPC-tested genes) had a nominally significant (P \leq 0.05) alteration of whole body BMD in knockout/knockdown mice, compared to controls. Of the 64, 49 (76.5%) were not members of the "known bone gene" list.

An example of one of the 64 novel genes is *GPATCH1*, located within a GWAS association on human chromosome 19q13.11. Of all the genes in the region, *GPATCH1* had the strongest TWAS association (P=3.44 x 10⁻²²⁶) (**Figure 3.2B**) and the strongest eQTL colocalization (whole blood, RCP=0.36) (**Figures 3.2B-D**).



Figure 3.2 TWAS and eQTL colocalization identify GPATCH1 as a novel potentially causal BMD GWAS gene. A) The top 40 terms from a Gene Ontology analysis of the 512 potentially causal BMD genes identified by our analysis. Terms with clear relevance to bone are highlighted in red. Only terms from the "Biological Process" (BP) subontology are listed. B) TWAS/colocalization plot for genes in the locus around GPATCH1 (\pm 1.5 Mbp). The $-log_{10}$ Bonferroni-adjusted P-values from the TWAS analysis (top panel) and the maximum RCPs from the colocalization analyses (bottom panel). Genes alternate in color for visual clarity. Triangles represent GPATCH1. C) Mirrorplot showing eBMD GWAS locus (top panel) and colocalizing GPATCH1 eQTL in whole blood (bottom panel). SNPs are colored by their LD with rs11881367 (purple), the most significant GWAS SNP in the locus. D) Scatterplot of $-log_{10}$ P-values for GPATCH1 eQTL versus eBMD GWAS SNPs. SNPs are colored by their LD with rs11881367 (purple). E) Bone mineral density (BMD) in Gpatch1 knockdown mice. N=7 females and N=4 males for Gpatch1^{+/-} mice, N=880 females and N=906 males for Gpatch1^{+/+} mice. Boxplots indicate the median (middle line), the 25th and 75th percentiles (box) and the whiskers extend to 1.5 * IQR.

The eQTL and BMD GWAS allele effects for the top SNPs were in the same direction, suggesting that decreasing the expression of *GPATCH1* would lead to decreased BMD. BMD data from the IMPC showed that female mice heterozygous for a *Gpatch1* null allele had decreased BMD ($P=2.17 \times 10^{-8}$) (**Figure 3.2E**). Together, these data suggest that many of the genes identified by the combined TWAS/colocalization approach are likely causal BMD GWAS genes.

3.3.3 *PPP6R3* is a candidate causal gene for a GWAS association on Chr. 11

To identify novel candidate genes for functional validation, we focused on genes with the strongest evidence of being causal. To do so, we increased the colocalization RCP threshold to 0.5, and then sorted genes based on TWAS Bonferroni-adjusted P-values. Furthermore, we constrained the list of candidates for functional validation to genes which were not members of the "known bone gene" list or genes with a nominal (P \leq 0.05) alteration in whole-body BMD as determined by the IMPC. This yielded 137 putatively causal BMD genes (**Table 3.2, Supplementary Data 3.8**).

Though it was not on the "known bone gene" list, the first gene ranked by TWAS Pvalue, *SPTBN1*, has been demonstrated to play a role in the regulation of BMD ⁵⁹. The second, *PPP6R3*, has not been previously implicated in the regulation of BMD. *PPP6R3* is located on human Chr. 11 within 1 Mbp of seven independent eBMD GWAS SNPs identified by Morris et al. ¹⁹ (subsequently referred to as "eBMD lead SNPs") (**Figure 3.3A**). Of all the protein-coding genes (N=29) in the ~1.8 Mbp region surrounding *PPP6R3*, its expression was the most significantly associated with eBMD by TWAS (Bonferroni = 5.7 x 10^{-93}) (**Figure 3.3B**).

Gene	Tissue with greatest RCP	Max. RCP	TWAS P-value (Bonferroni)
SPTBN1	Cells_Cultured_fibroblasts	0.9469	<5 x 10-324
PPP6R3	Thyroid	0.5291	5.7 x 10 ⁻⁹³
BARX1	Colon_Transverse	0.7764	6.36 x 10-63
MEOX2	Brain_Nucleus_accumbens_basal_ganglia	0.6286	3.21 x 10 ⁻⁵³
RERE	Adipose_Subcutaneous	0.6431	6.95 x 10 ⁻⁴⁶
SIPA1	Nerve_Tibial	0.9981	4.26 x 10-41
CAPZB	Testis	0.6716	3.64 x 10-33
B4GALNT3	Artery_Aorta	0.9241	2.67 x 10 ⁻³³
TRPC4AP	Breast_Mammary_Tissue	0.5577	8.62 x 10-31
AXL	Minor_Salivary_Gland	0.6205	9.74 x 10-31

Table 3.2 Top 10 novel protein-coding genes significant by colocalization (RCP ≥ 0.5) and TWAS, sorted by TWAS P-value.

Furthermore, *PPP6R3* was the only gene in the region with eQTL (in four GTEx tissues, thyroid, ovary, brain_putamen_basal_ganglia, and stomach with RCPs = 0.53, 0.50, 0.28 and 0.14, respectively) that colocalized with at least one of the eBMD associations (**Figure 3.3B**). Based on these data, we chose to further investigate *PPP6R3* as a potentially causal BMD gene.

We first determined which of the seven associations colocalized with the *PPP6R3* eQTL (**Figure 3.3C**). The most significant *PPP6R3* eQTL SNP in thyroid tissue (the tissue with the highest RCP) was rs10047483 (Chr. 11, 68.464237 Mbp) (*PPP6R3* eQTL P = 6.99 x 10^{-8} , eBMD GWAS P = 1.2×10^{-100}) located in intron 1 of *PPP6R3*. The most significant eBMD lead SNP in the locus was rs11228240 (Chr. 11, 68.450822 Mbp, eBMD GWAS P = 6.6×10^{-101} , *PPP6R3* eQTL P = 3.7×10^{-6}), located upstream of *PPP6R3*. Consistent with

the colocalization analysis, these two variants are in high LD ($r^2=0.941$) and rs10047483 does not exhibit strong LD ($r^2 < 0.104$) with any of the other six eBMD lead SNPs in the locus. The eQTL and BMD GWAS allele effects for rs10047483 were opposing, suggesting that a decrease in the expression of *PPP6R3* would lead to an increase in BMD.

A recent fracture GWAS identified 14 significant associations, one of which was located in the *PPP6R3* region (rs35989399, Chr. 11, 68.622433 Mbp)¹⁹. We analyzed the fracture GWAS in the same manner as we did above for eBMD. We found that *PPP6R3* expression when analyzed by TWAS was significant for fracture (TWAS Bonferroni-pval = $6.0 \ge 10^{-3}$) and the same *PPP6R3* eQTL colocalized with the fracture association (RCP = 0.49 in ovary, RCP = 0.36 in thyroid) (**Figure 3.3D**). Together, these data highlight *PPP6R3* as a strong candidate for one of the seven eBMD/fracture associations in this region.

3.3.4 *PPP6R3* is a regulator of femoral geometry, BMD, and vertebral microarchitecture

To assess the effects of *PPP6R3* expression on bone phenotypes, we characterized mice harboring a gene trap allele (*Ppp6r3*^{tm1a(KOMP)Wtsi}) (**Figure 3.4A**). We intercrossed mice heterozygous for the mutant allele to generate mice of all three genotypes (wild-type (WT), heterozygous (HET), and mutant (MUT)). The absence of PPP6R3 protein in MUT mice was confirmed through Western blotting (**Figure 3.4B**).

The BMD analyses presented above used heel eBMD GWAS data. We used these data because they represent the largest, most well-powered BMD GWAS to date ¹⁶. However, to determine whether perturbation of *Ppp6r3* would be expected to impact

femoral or lumbar spine BMD in a similar manner, we turned to a smaller GWAS to look at both of these traits. In a GWAS by Estrada et al. ¹⁶, a total of 56 loci were identified for femoral neck (FNBMD) and lumbar spine (LSBMD) BMD. One of the 56 loci corresponded to the same SNPs associated with the *PPP6R3* eQTL. The locus was significant for LSBMD; however, it did not reach genome-wide significance for FNBMD (**Supplemental Figure 3.1**).

We evaluated BMD at both the femur and the lumbar spine in *Ppp6r3*^{em1a(KOMP)Wtsi} mice, with the expectation, based on the above data, that perturbation of *Ppp6r3* would have a stronger impact on BMD at the lumbar spine. At approximately 9 weeks of age, we measured areal BMD (aBMD) at the femur and lumbar spine using dual X-ray absorptiometry (DXA). First, we observed no change in body weight at 9 weeks that might impact bone phenotypes (**Supplemental Figure 3.2A**). As the above analysis predicted, we observed a significant effect of *Ppp6r3* genotype on aBMD at the lumbar spine (WT vs. MUT P=0.01, **Figure 3.4C**), but not the femur (WT vs. MUT P=0.26, **Figure 3.4D**). It should also be noted, however, that we observed significantly decreased femoral width, but not length, in *Ppp6r3* mutant mice (anterior-posterior (AP) femoral width, WT vs. MUT P=0.02; medial-lateral (ML) femoral width, WT vs. MUT P=2.2 x 10⁻⁶, **Supplemental Figures 3.2B-D**).



Figure 3.3 PPP6R3 is a top-10 novel eBMD gene. A) eBMD GWAS SNPs around the PPP6R3 locus (± 1 Mbp). The y-axis represents $-log_{10}$ eBMD GWAS P-values. Highlighted SNPs (black) are the seven lead eBMD GWAS SNPs in the locus. B) TWAS/colocalization plot for genes in the locus around PPP6R3 (± 1 Mbp). The $-log_{10}$ Bonferroniadjusted P-values from the TWAS analysis (top panel) and the maximum RCPs from the colocalization analyses (bottom panel). Genes alternate in color for visual clarity. Triangles represent PPP6R3. Mirrorplot of the eBMD locus (C) and PPP6R3 eQTL in thyroid, and fracture locus and PPP6R3 eQTL in thyroid (D). The panels on the right are scatterplots of $-log_{10}$ P-values for PPP6R3 eQTL and eBMD GWAS SNPs (C) and the PPP6R3 eQTL and fracture GWAS SNPs (D). SNPs are colored by their LD with rs10047483 (purple), the most significant PPP6R3 eQTL in the locus. Not all genes are shown.

Due to the significant effect of *Ppp6r3* genotype on lumbar spine aBMD, we further characterized the effects of *Ppp6r3* genotype on microarchitectural phenotypes via microcomputed tomography (μ CT). We observed significant (P \leq 0.05) decreases in trabecular bone volume fraction (BV/TV, WT vs. MUT P=0.015, **Figure 3.4E-F**) and volumetric BMD (vBMD, WT vs. MUT P=0.015, **Figure 3.4G**) of the lumbar spine as a function of *Ppp6r3* genotype, but found no significant changes in tissue mineral density (TMD, **Supplemental Figure 3.2E**), trabecular separation (TbSp), trabecular thickness (TbTh) or trabecular number (TbN) (**Figures 3.4 H-J**).



Figure 3.4 Ppp6r3 functional validation shows an effect of genotype on bone mass. A) Schematic of the Ppp6r3 gene-trap allele (Ppp6r3 tm1a(KOMP)Wtsi). Image obtained from the IMPC. B) Western blot of the Ppp6r3 experimental mice. Top left panel shows that PPP6R1 protein (control) levels are not affected by the Ppp6r3 gene-trap allele. Top right panel shows the effect of the gene-trap allele on PPP6R3 protein levels. The two bands are ostensibly due to different PPP6R3 isoforms. Bottom panel shows that PP6C protein (control) levels are not affected by the Ppp6r3 gene-trap allele. Least-squares means for spinal (C) and femoral (D) areal BMD (aBMD) DXA in Ppp6r3 experimental mice. Scale is shown on the bottom left. F-J) Least-squares means for μ CT measurements in the lumbar spines of Ppp6r3 WT, HET, and MUT mice. Contrast P-values, adjusted for multiple comparisons are presented. *P \leq 0.05. Abbreviations: BV/TV - bone volume fraction, vBMD – volumetric bone mineral density, TbSp – trabecular separation, TbTh – trabecular thickness, TbN – trabecular number.

3.4 Discussion

BMD GWASs have identified over 1,100 associations to date. However, identifying causal genes remains a challenge. To aid researchers in further dissecting the genetics of complex traits, reference transcriptomic datasets and computational methods have been developed for the prioritization and identification of causal genes underlying GWAS associations. In this work, our goal was to utilize these data and tools to prioritize putatively causal genes underlying BMD GWAS associations. Specifically, we used the GTEx eQTL reference dataset in 49 tissues to perform TWAS and eQTL colocalization on the largest BMD GWAS. Using this approach, we identified 512 putatively causal protein-coding genes that were significant in both the TWAS and colocalization approaches.

Our approach was inspired by a recent study that used the GTEx resource and a TWAS/eQTL colocalization approach similar to the one we employed. Pividori et al. ⁷³ recently combined TWAS and eQTL colocalization to GTEx and GWAS data on 4,091 traits, including BMD, from the UK Biobank data. A total of 76 protein-coding genes were identified and of the 76, we identified 55 (72.4%) of the same genes in our implementation. There are several reasons for this discrepancy in the number of prioritized genes. First, both studies used a GWAS based on the UK BioBank ²²¹; however, there were significant differences in sample size. The PhenomeXcan project utilized GWAS data based on the analysis of ~207,000 individuals, whereas we used GWAS data based on the analysis of ~426,000 individuals ^{19,73}. Second, the two GWAS studies utilized different association models. Finally, due to the breadth of the PhenomeXcan project, they had a higher multiple-testing burden than we did, which led to different Bonferroni-adjusted P-value thresholds ($P<5.49 \times 10^{-10} \text{ vs. } P\leq 2.38 \times 10^{-6}$).

One of many novel genes identified in our study was *PPP6R3*, which was also identified in the PhenomeXcan project ⁷³. *PPP6R3* is a regulatory subunit of protein phosphatase 6 and has been implicated in several cancers^{222,223}. In humans, the PPP6R3 protein shows ubiquitous expression across tissues, and may have an important role in maintaining immune self-tolerance ²²³. It is unclear how *PPP6R3* may be influencing BMD. However, Protein Phosphatase 6 has been shown to oppose activation of the nuclear factor kappa-light-chain enhancer of activated B cells (NF-xB) pathway in lymphocytes ²²⁴. Since the NF-xB signaling pathway is highly involved in osteoclastogenesis and bone resorption, it is possible that *PPP6R3* may be involved in the regulation of this pathway in osteoclasts²²⁵.

The *PPP6R3* locus demonstrated a high level of complexity, containing seven independent GWAS associations, at least one of which was also associated with fracture. Interestingly, just upstream of *PPP6R3* is *LRP5*, a WNT signaling co-receptor ²²⁶. *LRP5* is a well-known regulator of BMD and gain and loss of function mutations lead to high bone mass syndrome and osteoporosis pseudoglioma, respectively ^{140,142,227,228}. *LRP5* expression was not significantly associated with eBMD by TWAS (Bonferroni P= 1), nor did it have a colocalizing eQTL in GTEx tissues (most significant RCP=1.6 x 10⁻² in pancreas). However, another eBMD lead SNP in the region, rs4988321, is a missense mutation in *LRP5* (Val667Met) that has been associated with BMD in multiple studies ^{229–231}. While this variant represents an association that is independent of the rs10047483 association ($r^2 = 0.104$), it further highlights the complexity of this locus both in terms of the number of associations as well as target genes.

To determine the effect of *Ppp6r3* expression on bone, we characterized bone phenotypes in mice harboring a gene-trap allele (*Ppp6r3*^{tm1a(KOMP)Wtsi}). Consistent with the

observation that the *PPP6R3* eQTL SNPs were significantly associated with lumbar spine, but not femoral neck BMD, we observed that *Ppp6r3* deletion had a significant effect on lumbar spine BMD, but we did not observe an overall effect on femoral BMD. Using μ CT, we further characterized the effect of *Ppp6r3* deletion on lumbar spine microarchitecture. We observed significant decreases in trabecular bone volume fraction (BV/TV) and volumetric BMD of the lumbar spine as a function of *PPP6R3* genotype. While we did not observe significant effects of *Ppp6r3* deletion on trabecular thickness or number, the direction of effects for those phenotypes suggests that the observed decrease in bone volume fraction and BMD may be explained by the cumulative effects of *Ppp6r3* deletion on trabecular thickness and number.

Our hypothesis regarding the directions of effect of *Ppp6r3* expression on BMD based on the eQTL and eBMD/lumbar spine BMD GWAS were opposite to what we observed. There are several reasons that may explain this. First, our hypothesis was based on expression data in non-bone tissues and cell-types. Recent studies have shown that the direction of eQTL effects can differ between different cells and tissues within humans ^{232,233}. Second, our hypothesis was based on human data, while our functional experiments were performed in mice. Third, we globally deleted *Ppp6r3* in mice, as opposed to ablating it in a bone-specific knockout. Future studies of the *PPP6R3* eQTL in bone cells as well as the generation of conditional *Ppp6r3* knockouts will allow us to unravel the precise role of this association and *PPP6R3* in the regulation of bone mass.

As we and others have shown, the use of both TWAS and eQTL colocalization can prioritize putatively causal genes underlying GWAS associations. Here, we have shown the utility of this approach even in the absence of eQTL data from the most phenotype-relevant

tissue. However, it is important to highlight the limitations of our analysis. While studies have shown that many eQTL are shared among tissues, the lack of eQTL data in bone and bone cells means that bone-specific eQTL were missed. In addition, the use of multiple nonbone tissues may have inflated the number of false positives based on coincidence of strong TWAS and eQTL colocalization signals that have no biological impact on bone. Furthermore, the lack of bone transcriptomic data may also explain the observed disparity between our hypothesized and observed direction-of-effect for PPP6R3. It is also important to note that due to the reliance of this approach on eQTL data, genes that affect BMD via non-expression related mechanisms were not captured. Another limitation of our approach arises from the definition of loci based on linkage disequilibrium (LD). We used a set of previously-defined approximately independent LD blocks, derived from a cohort of European individuals, in our fastENLOC analysis ²³⁴. The inexact nature of these data may lead to spurious colocalizations due to mismatches in LD structure between the reference LD blocks and the GWAS/eQTL populations. Additionally, because the GWAS and eQTL data have mismatching LD structures, due to their being derived from cohorts with different ancestries, our analyses, particularly the colocalization analyses, may suffer from reduced power⁶¹. This also raises the related issue of the reduced generalizability of our results in non-European individuals, which brings further attention to the necessity of performing GWASs and providing reference data in diverse and underrepresented populations. Finally, another issue arises when considering correlations in expression, and predicted expression, between genes in a locus, which may lead to spurious associations in TWAS analyses ⁷⁰.

In summary, we applied a combined TWAS/colocalization approach using GTEx and identified 512 putatively causal BMD genes. We further investigated *PPP6R3* and demonstrated that it is a regulator of lumbar spine BMD. We believe this work provides a

valuable resource for the bone genetics community and may serve as a framework for prioritizing genes underlying GWAS associations using publicly available tools and data for a wide range of diseases.

3.5 Methods

fastENLOC colocalization:

For each of the eBMD and fracture GWASs, we performed colocalization using fastENLOC, by following the tutorial and guidelines available at https://github.com/xqwen/fastenloc.

Briefly, for each GWAS, we converted variant coordinates to the hg38 human genome assembly, using the UCSC liftOver tool (minimum ratio of bases that must remap = 1) [https://genome.ucsc.edu/cgi-bin/hgLiftOver]. We calculated z-scores by dividing GWAS betas by standard errors. We then defined loci based on European linkage disequilibrium (LD) blocks, as defined based on the results of Berisa and Pickrell, 2015²³⁴.

Z-scores were then converted to posterior inclusion probabilities (PIPs) using torus ²³⁵. Finally, these data were colocalized with fastENLOC for all 49 GTEx V8 tissues, with the "total_variants" flag set to 14,000,000. Colocalization was performed using pre-computed GTEx multi-tissue annotations, obtained from <u>https://github.com/xqwen/fastenloc</u>. Finally, to identify protein-coding genes in the results, we utilized Ensembl's "hsapiens_gene_ensembl" dataset using biomaRt (version 2.45.8).

S-MultiXcan:

We conducted a transcriptome-wide association study by integrating genome-wide SNP-level association summary statistics from an estimated bone mineral density GWAS
(Morris et al. 2018) with GTEx version 8 gene expression QTL data from 49 tissue types. We used the S-MultiXcan (Barbeira et al. 2019) approach for this analysis, to correlate gene expression across tissues to increase power and identify candidate susceptibility genes. Default parameters were used, with the exception of the "--cutoff_condition_number" parameter, which was set to 30. Bonferroni-correction of P-values was performed on the resultant gene set (22,337 genes), using R's p.adjust function. This was followed by the removal of non-protein-coding genes. The analysis was also performed in the same manner using summary statistics from a fracture GWAS (CITE). Finally, to identify protein-coding genes in the results, we utilized Ensembl's "hsapiens_gene_ensembl" dataset using biomaRt^{236,237}.

Creation of the "known bone gene" list:

We generated a "known bone gene" set as follows: First, we downloaded Gene Ontology IDs for the following terms: "osteo*", "bone", and "ossif*" from AmiGO2 (version 2.5.13)¹⁹⁰. After removal of non-bone related terms, we extracted all mouse and human genes related to the GO terms, using biomaRt. From this list, we retained proteincoding genes.

We also used the "Human-Mouse: Disease Connection" database available at the Mouse Genome Informatics website, to download human and mouse genes annotated with the terms "osteoporosis", "bone mineral density", "osteoblast", "osteoclast" and "osteocyte". We used biomaRt to identify the gene biotypes, and retained protein-coding genes. We then used the MGI human-mouse homology table [http://www.informatics.jax.org/downloads/reports/HOM_MouseHumanSequence.rpt] to convert all mouse genes to their human homologs. Finally, we removed genes that weren't interrogated in both the colocalization and the TWAS analyses.

Gene Ontology enrichment analyses:

Gene ontology analysis was performed for the set of protein-coding genes passing the colocalization threshold RCP \geq 0.1 and S-MultiXcan Bonferroni P-value \leq 0.05, using the "topGO" package (version 2.40.0) in R²³⁸. Enrichment tests were performed for the "Molecular Function", "Biological Process" and "Cellular Component" ontologies, using all protein-coding genes that were subjected to colocalization and multiXcan analysis as background. Enrichment was performed using the "classic" algorithm with Fisher's exact test. P-values were not adjusted for multiple testing.

Linkage disequilibrium calculations:

Linkage disequilibrium between variants was calculated using the LDlinkR (version 1.0.2) R package, using the "EUR" population ²³⁹.

PPP6R3 knockout mouse generation:

The study was carried out in strict accordance with NIH's Guide for the Care and Use of Laboratory Animals. Additionally, the University of Virginia Institutional Animal Care and Use Committee approved all animal procedures. *Ppp6r3* gene trap mice were generated using targeted embryonic stem cell clones heterozygous for the *Ppp6r3*^{tm1a(KOMP)Wtsi} gene trap allele obtained from the International Knockout Mouse Project (KOMP; [https://www.komp.org]). KOMP ES clones were karyotyped and injected using a

XYClone Laser (Hamilton Thorne, Beverly, MA) into B6N-Tyr^{c-Brd}/BrdCrCrl (Charles River, Wilmington, MA) 8-cell stage embryos to create chimeric mice. Resultant chimeras were bred to B6N-Tyr^{c-Brd}/BrdCrCrl mice to obtain germline transmission of the *Ppp6r3* gene trap allele. From a breeding pair of two heterozygous mice, we generated our experimental population through HET x HET matings. Breeder mice were fed a breeder chow diet (Envigo Teklad S-2335 mouse breeder sterilizable diet, irradiated. Product # 7904), and experimental mice were fed a standard chow diet (Envigo Teklad LM-485 irradiated mouse/rat sterilizable diet. Product #7912).

Genotyping of PPP6R3 mice:

DNA for genotyping was extracted from tail clips as follows: tail clips were incubated overnight at 55° C in a solution of 200uL digestion/lysis buffer (Viagen Direct PCR (tail), Los Angeles, CA) and 1mg/mL proteinase K (Viagen, Los Angeles, CA). After overnight incubation, tails were heated at 85° C for 45 minutes, and solutions were subsequently stored at 4° C.

For genotyping, PCR reactions were set up as follows. For each reaction, 1 μ L of DNA was mixed with 24 μ L of a master mix consisting of 19.5 μ L nuclease-free H₂O, 2.5 μ L 10x PCR reaction buffer (Invitrogen, Waltham, MA), 0.75 μ L of mgCl2 (Invitrogen, Waltham, MA), 0.5 μ L of 10mMol Quad dNTPs (Roche Diagnostics GmbH, Mannheim, Germany) 0.25 μ L of Platinum Taq DNA polymerase (Invitrogen, Waltham, MA), and 0.25 μ L of each primer, diluted to 20 μ Mol. PCR primers were obtained from Integrated DNA Technologies, Coralville, IA.

Forward primer: 5'- CAC CTG GGT TGG TTA CAT CC -3'

Reverse primer: 5'- GAC CCT GCC TTA AAA CCA AA -3'

The following PCR settings were used:

- Initialization: 94° C, 120s
- Denaturation: 94° C, 30s (37 cycles)
- Annealing: 54° C, 30s (37 cycles)
- Elongation: 72° C, 35s (37 cycles)
- Final elongation: 72° C, 300s

PCR products were run on a 2% agarose gel for 150 minutes at 60 volts, to distinguish between wild-type, heterozygous and mutant *Ppp6r3* mice.

PPP6R3 Western blotting:

Mouse spleens 20-40 mg in weight were suspended in 1% NP40 buffer (50 mM Tris (pH 8) 100 mM NaCl, 1 % NP40, 1 mM EGTA, 1 mM EDTA, Protease inhibitor cocktail (04-693-116-001, Roche), 1 mM PMSF, 50 mM NaF, 0.2 mM sodium vanadate). The tissue was homogenized by RNase-free disposable pestles (ThermoFisher #12-141-364) and incubated for 10 min on ice. After brief sonication, the sample was centrifuged for 10 min at 13,000 x rpm at 4C. The protein concentration in the extract was measured by Bradford assay. 100ug of sample protein was boiled 5 min in SDS sample buffer, loaded in each lane, resolved by gradient SDS– PAGE (Bio-Rad #456-1085) and immunoblotted as described in Guergnon et al ²⁴⁰. Primary antibodies were diluted 1:1000. (SAPS1 Ab: ThermoFisher #PA5-44275, SAPS3 Ab: ThermoFisher #PA5-58405, PP6C Ab: Sigma #HPA050940)

PPP6R3 functional validation:

Experimental mice were sacrificed at approximately 9 weeks of age (mean age = 61 days). At sacrifice, the right femurs were isolated, and femoral morphology (length and widths in AP and ML orientations) was measured with digital calipers (Mitoyuto American, Aurora, IL). Femurs were then wrapped in PBS-soaked gauze and stored at -20° C, until analysis. Lumbar vertebrae L3-L5 were also dissected at sacrifice and were wrapped in PBS-soaked gauze and frozen at -20° C.

Dual X-ray absorptiometry:

Individual right femurs and the lumbar spine (L5 vertebrae) were isolated from surrounding soft tissues and frozen at -20° C in PBS. Dual X-ray absorptiometry (DXA) was performed on the femurs and lumbar vertebrae using the Lunar Piximus II (GE Healthcare) as described previously by Beamer et al ²⁴¹. In short, 10 isolated bones were placed in the detector field at a time and the samples were analyzed one by one, such that the region of interest (ROI) was set for one specimen at a time for data collection. The ROI for the femurs was on the entire isolated femur. For the spine, was on the entire isolated L5. Care was taken to ensure that the sample orientation was identical for all samples.

Micro-computed tomography and image analysis:

All μ CT analyses were carried out at the μ CT Imaging Core Facility at Boston University using a Scanco Medical μ CT 40 instrument (Brütisellen, Switzerland). The power, current, and integration time used for all scans were 70 kVp, 113 μ A, and 200 msec respectively. The L5 vertebrae were scanned at a resolution of 12 microns/voxel. Two

volumes of interest (VOIs) were selected for analysis: 1) the entire portion of the L5 vertebra extending from 60 microns caudal to the cranial growth plate in the vertebral body to 60 microns cranial to the caudal growth plate; and 2) only the trabecular centrum contained in the first VOI. Semi-automated-edge detection (Scanco Medical) was used to define the boundary between the trabecular centrum and cortical shell to produce the second VOI. Gaussian filtering (sigma = 0.8, support = 1) was used for partial background noise suppression. A scan of a potassium hydroxyapatite phantom allowed conversion of grayvalues to mineral density. For segmentation of bone tissue, the threshold was set at a 16bit gray value of 7143 (521 mgHA/ccm), and this global threshold was applied to all of the samples. For each VOI, the following were calculated: total volume (TV), bone volume (BV), bone volume fraction (BV/TV), bone mineral density (BMD), and tissue mineral density (TMD). BMD was defined as the average density of all voxels in the VOI, whereas TMD was defined as the average density of all voxels in the VOI above the threshold 173 . For the second VOI, the following additional parameters were calculated: trabecular thickness (Tb.Th), trabecular separation (Tb.Sp), trabecular number (Tb.N), connectivity density (Conn.D), and structure model index (SMI)¹⁷³.

Statistical analyses:

To calculate the enrichment of bone genes in prioritized genes, we performed Fisher's exact test, using R's "fisher.test" function, with the alternative hypothesis set as "greater".

For the statistical analysis of the phenotyping results, we calculated least-squares means (lsmeans) using the "emmeans" R package (version 1.5.2.1) ²⁴². Input for the lsmeans function was a linear model including terms for genotype, weight and age in days. For sex-

combined data, we also added a term for sex. For all phenotypes, we also included a term for weight. Finally, for DXA phenotypes, we included a term for "CenterRectX" and "CenterRectY".

We used Tukey's HSD test to test for significant differences in lsmeans, for each pair of genotype levels. Tukey's HSD also controls the family-wise error rate.

Analyses involving data from the International Mouse Phenotyping Consortium:

For the IMPC data, we obtained data using their "statistical-result" SOLR database, using the "solrium" R package (version 1.1.4) ²⁴³. We obtained experimental results using the "Bone*Mineral*Density" parameter. We then pruned the resulting data to only include "Successful" analyses, and removed experiments that included the skull. To generate the *Gpatch1* boxplot, we obtained raw data using from IMPC's "statistical-raw-data" SOLR database for *Gpatch1*, and analyzed the data in the same manner as IMPC, using the "OpenStats" R package (version 1.0.2), using the method="MM" and MM_BodyWeightIncluded = TRUE arguments ²⁴⁴. Finally, mouse genes were converted to their human syntenic counterparts using Ensembl's "hsapiens_gene_ensembl" and "mmusculus_gene_ensembl" datasets through biomaRt.

PhenomeXcan data analysis:

We obtained all significant PhenomeXcan gene-trait associations from their paper [https://advances.sciencemag.org/content/6/37/eaba2083], and used data for the "3148_raw-Heel_bone_mineral_density_BMD" phenotype ⁷³. Furthermore, we constrained our search to only include genes that were annotated by the authors as "protein_coding".

LSBMD/FNBMD GWAS analysis:

We obtained sex-combined LSBMD and FNBMD GWAS summary statistics from GEFOS [http://www.gefos.org/?q=content/data-release-2012], and then used a custom script that utilized the biomaRt R package to convert variants to their GRCh38 coordinates.

3.6 Data availability

eBMD and fracture GWAS summary statistics were obtained from GEFOS, as were the LSBMD and FNBMD GWAS summary statistics. GTEx eQTL data were obtained from the GTEx web portal. Data from the PhenomeXcan project were obtained from Pividori et al ⁷³. Statistical data from the IMPC were obtained using an R interface to their SOLR database. *Ppp6r3* experimental data are provided on our GitHub [https://github.com/baselmaher/BMD_TWAS_colocalization]. Mouse-Human homologs were obtained from MGI [http://www.informatics.jax.org/downloads/reports/HOM_MouseHumanSequence.rpt]. We also obtained data from the MGI Human-Mouse:Disease Connection database [http://www.informatics.jax.org/diseasePortal]. Gene Ontologies were obtained from AmiGO2 [http://amigo.geneontology.org/amigo].

3.7 Code availability

Analysis code and the raw data for our *Ppp6r3* functional validation analyses are available on GitHub [https://github.com/basel-maher/BMD_TWAS_colocalization].

3.8 Acknowledgements

Research reported in this publication was supported in part by the National Institute of Arthritis and Musculoskeletal and Skin Diseases of the National Institutes of Health under Award Number AR071657 to C.R.F., L.C.G and E.F.M., and by the National Center for Research Resources of the National Institutes of Health under Award Number S10RR021072 to E.F.M. B.M.A-B was supported in part by a National Institutes of Health, Biomedical Data Sciences Training Grant (5T32LM012416). The authors acknowledge Wenhao Xu (University of Virginia) and the Genetically Engineered Mouse Models (GEMM) core for their technical assistance in generating the *Ppp6r3* gene-trap mice. We thank the IMPC for accessibility to BMD data in knockout mice (www.mousephenotype.org). The data used for the analyses described in this manuscript were obtained from the IMPC SOLR database on 3/8/21. The Genotype-Tissue Expression (GTEx) Project was supported by the Common Fund of the Office of the Director of the National Institutes of Health, and by NCI, NHGRI, NHLBI, NIDA, NIMH, and NINDS. The data used for the analyses described in this manuscript were obtained from the GTEx Portal on 6/30/20. Chapter 4

Concluding Remarks and Future Directions

4.1 Summary and conclusions

In the field of bone genetics, GWASs have been very successful in identifying genomic loci associated with bone-related traits, with over 1,100 associated loci identified to date ^{16,19}. While these findings have contributed to our understanding of the genetics underlying bone biology, our progress as a field has been hindered by two main shortcomings of bone-related GWASs: 1) GWASs in the bone field have overwhelmingly been focused on BMD as a trait, and 2) researchers have encountered difficulties in identifying the causal genes underlying GWAS associations. In this dissertation, we tackled these shortcomings, using both mouse and human data, by utilizing systems genetics approaches.

4.1.1 Association mapping of bone-related traits

By performing a GWAS for 55 bone-related traits in Chapter 2, we began to address the issue of the strict focus of bone-related GWAS on BMD. We utilized a powerful new mouse resource, the Diversity Outbred, in order to perform a GWAS for 55 bone-related traits. Due to the strict focus of bone-related GWASs on BMD in humans, our analysis begins to address a significant knowledge gap in our understanding of the genetics of bone. From the GWAS that we performed, we identified 28 QTL for 20 different bone-related phenotypes, which captured novel biology not yet implicated by human GWAS. Using transcriptomic data that we generated from our DO mouse cohort, in conjunction with genomic data from the fully sequenced DO mouse founder mice, we were able to identify 18 putatively causal genes within six of the QTL loci. Through the use of fine mapping and the characterization of a mouse knockout, we were able to demonstrate that *Qsox1* was at least partially responsible for one of the identified QTL affecting cortical bone morphology.

Our results not only provide a resource that serves to increase our understanding of the genetics of bone, but also illustrate the particular utility of the DO in identifying highresolution QTL and prioritizing some of the genes underlying these loci. One of the most exciting examples of the utility of the DO is demonstrated in our dissection of the locus on Chr. 1. Since our GWAS encompassed many traits, the Chr. 1 locus contained QTL for traits in two different phenotypic categories: cross-sectional size and TMD/cortical porosity. We were able to leverage the genetic diversity of the DO population in order to tease apart the locus and conclude that it comprised of at least two loci affecting different aspects of bone. These results demonstrate the unique power of the DO (and the utility of performing multitrait GWASs on the same population) in understanding complex trait genetics.

4.1.2 Systems genetics analyses in mouse inform human BMD GWAS

In conceiving the work presented in this dissertation, we hypothesized that systems genetics analyses in the mouse can inform human BMD GWAS. By using network-based approaches, we leveraged mouse transcriptomic data to identify putatively causal genes underlying human GWAS associations, thus tackling another shortcoming of bone-related GWAS. Using transcriptomic data from DO mouse cortical bone, we generated coexpression networks, and then learned Bayesian networks underlying each co-expression network. We then used a method inspired by key driver analysis in order to identify 688 genes (BANs) that were significantly connected to known bone genes in the Bayesian networks. In order to identify genes most likely to be responsible for BMD GWAS associations, we then identified BANs likely to be regulated by human eQTL, using eQTL colocalization analysis. This led to the identification of 66 genes that were both associated in their expression with known bone genes, and also had eQTL that colocalized with BMD GWAS loci. Forty-seven (~71%) of these genes were novel, showing that our approach was both able to recall known biology and identify novel putatively causal genes underlying GWAS loci.

To further illustrate the utility of our network-based approach, we then used the results of our co-expression network analysis, in addition to single-cell RNA sequencing, to generate hypotheses regarding the specific bone traits that some of these genes (e.g., *SERTAD4* and *GLT8D2*) may affect, as well as the bone cell-types through which their effects may be exerted. The results of our network-based analyses reinforce the utility of the mouse as a model for informing human genetics via systems genetics methods, and provide an avenue for further work aimed at understanding the relationships between genes and specific traits, and the elucidation of the trait-relevant cell-types.

Overall, the study performed in Chapter 2 demonstrates the power of the DO in performing high-precision genetic mapping analyses in the context of bone and the ability of systems genetic analyses in mouse to inform human GWAS. Thus, this work provides a powerful resource for the bone genetics field, as well as methodologies for informing human GWAS associations that are applicable across myriad complex traits.

4.1.3 A combined TWAS/eQTL colocalization approach informs BMD GWAS

Recently, large-scale projects (e.g., GTEx) have produced a wealth of publicly available data that are amenable to systems genetics analyses. Great methodological and analytical strides have also been made, by many groups, to make use of these data. In Chapter 3, we sought to utilize publicly available data, using recently developed systems genetics methods, to further address the difficulty in identifying causal genes underlying BMD GWAS loci.

By utilizing the GTEx eQTL reference data, we performed a combined TWAS/eQTL colocalization approach in order to identify 512 putatively causal genes underlying GWAS loci. These 512 genes were enriched in known bone genes, lending support to our results. Furthermore, we presented an example of a novel candidate GWAS gene, *GPATCH1*, which showed an effect on BMD in publicly available mouse knockdown data. We then narrowed down the list of 512 genes to identify 137 novel, higher-confidence candidate genes. We characterized the locus surrounding the gene with strongest support as a candidate gene, *PPP6R3*, and performed experimental validation of its effect on BMD in mice. We found that *Ppp6r3* deletion decreased BMD.

The work presented in Chapter 3 serves to provide support for the utility of systems genetics approaches, namely combined TWAS/eQTL colocalization, in identifying and prioritizing genes underlying GWAS associations. Specifically, the results of the study serve as a resource for further elucidating the genetic determinants of BMD. Perhaps more interestingly, however, our study provides support for the notion that systems genetics approaches can inform complex trait genetics even in the absence of data from trait-relevant cells and tissues.

4.2 Limitations

4.2.1 Bone-specific -omics data

As we have previously noted, while our approaches have been successful in utilizing eQTL data from non-bone tissues, we have likely missed many bone-specific genes due to the lack of tissue-specific eQTL data in bone, and may have suffered from an increased false-positive rate due to spurious non-bone related associations. While the overwhelming majority of eQTLs in LD with a disease mutation are not tissue-specific, genes with tissue-specific eQTL stend to be more highly enriched in disease-associated genes, suggesting that approaches utilizing eQTL data, such as TWAS and eQTL colocalization, may benefit from the use of tissue-specific eQTL study in osteoclasts identified a significant proportion of osteoclast-specific eQTL, when compared to GTEx eQTL data²⁴⁶. While our group and others are currently generating small-scale bone-specific eQTL data, more concerted efforts from large consortia, such as GTEx, are required for the progression of the bone genetics field ⁵³⁻⁵⁵.

4.2.2 Single-cell -omics data

Furthermore, due to the cellular heterogeneity exhibited in tissues comprising reference datasets, such as GTEx, and bulk "-omics" data in general, analyses of these data only reveal findings from an average reading across constituent cell populations, and may miss context-dependent eQTLs that are often relevant to complex phenotypes ^{247,248}. While such data are highly informative, the fields of gene discovery and disease therapeutics can be

advanced by the use of cell-type specific data ²⁴⁹. For example, by using epigenomic annotations associated with regulatory region activity, Claussnitzer et al. identified preadipocytes as the cell-types of action underlying associations of the FTO locus with obesity. Using eQTL data in preadipocytes, they then identified two genes, IRX3 and IRX5, as causal, with strong eQTL at the preadipocyte level. However, when the analysis was repeated in whole-adipose tissue, they noted an absence of an eQTL signal, suggesting that the genetic locus likely functions in a cell-specific manner²⁵⁰. Bone is no different; as we have shown above in our scRNA-seq analysis (Chapter 2), bone tissue is comprised of transcriptionally heterogeneous cell-types at different developmental stages. Ongoing studies have also demonstrated cellular heterogeneity in cultured osteoblasts, and have even identified new transcriptionally distinct bone cell-types ^{251,252}. Therefore, we believe that increasing the availability of data from different bone-cell types at different developmental stages will surely increase our understanding of bone genetics and will, ultimately, enhance bone disease-related therapeutics. Currently, promising avenues for the generation of these data include projects aiming to provide reference maps of all human cells, such as the Human Cell Atlas²⁵³, as well as recent developments in computational and experimental methods that aim to generate data on the cellular level, such as single-cell genomics, singlecell assay for transposase-accessible chromatin using sequencing (sc-ATAC-seq), and celltype deconvolution of bulk transcriptomics data ^{254–256}. While in their infancy, methods are also being developed to leverage scRNA-seq data to inform GWAS. In general, these methods aim to infer trait-relevant tissues and cell-types using GWAS summary statistics^{257,258}.

4.2.3 Beyond steady-state transcriptomics data

In our approaches, we have primarily focused on the use of steady-state, *cis*-eQTL data to inform GWAS and identify putatively causal genes. However, only a small proportion of expression variation can be attributed to *cis*-eQTL²⁴⁵. It is likely that larger studies involving very large sample sizes may, providing the continued development of related computational approaches, be able to identify *trans*-eQTL that explain a large proportion of the variance in expression. However, since the effect sizes of *trans*-eQTL are very small, the contribution of such studies to the identification of causal genes is still to be seen. It is likely more beneficial to focus on applying the systems genetics approaches we have utilized to other regulatory data. For example, the GTEx project has recently made available splicing QTL data (although not in bone), which are amenable to systems genetics analyses ⁵⁰. Data concerning non-eQTL regulatory mechanisms, such as chromatin accessibility QTL, protein abundance QTL, and methylation QTL, as well as approaches that utilize these data, such as cistrome-wide association studies (CWAS, an extension of TWAS), will likely contribute significantly towards the identification of bone-related genes and drug discovery ²⁴⁸. While some bone-relevant regulatory data is available, such as the data provided by the ENCODE project from primary osteoblasts²⁵⁹, it is important to generate a wide variety of regulatory data, form large cohorts, encompassing myriad bone cell-types at different developmental stages.

4.2.4 Non-European ancestries

In our studies, we utilized GWAS and eQTL data generated in majority Europeanancestry populations. In order to increase both the power to discover novel biology underlying bone traits and the extensibility of our findings to non-European populations, it is crucial to perform large-scale studies in non-European populations. To date, a few small bone-related association studies have been performed in non-European populations. While small in size, the emerging trend of performing such association studies in non-European populations is encouraging ^{260–266}.

4.2.5 Investigating the genetics of other aspects of bone

Finally, a major challenge in the field of bone genomics is the focus on BMD as a phenotype. Although BMD is the most clinically-relevant predictor of osteoporosis, and is easily measured in large cohorts, studying non-BMD bone traits (such as bone strength and microarchitectural properties) will be extremely beneficial in increasing our understanding of bone development and in treating diseases of bone. Generating non-BMD data in large populations will be challenging, but technological progress is beginning to allow such studies²².

4.3 Future directions

While the work presented in this dissertation provides a wealth of information regarding the genetics of bone traits, more work can be done to further our understanding of the genetics of bone. Much of this work can be performed as a continuation of our results. For example, our data provide some support for *Ier5* as a causal gene underlying some of the Chr. 1 locus traits, particularly TMD. Performing a knockout experiment, in mice, to test this hypothesis will serve to further unravel the genetics of the Chr. 1 locus. Another

extension of our results would be to further characterize the effects of *PPP6R3* on bone, by generating a bone-specific knockout mouse, which may shed light on the discordance of our expected direction-of-effect with our observed results, and may also further elucidate the mechanisms by which *PPP6R3* affects bone traits.

Another avenue for future work involves the refinement of our results and methodologies. Performing colocalization analyses using splicing QTL from GTEx, for example, may serve to further inform BMD GWAS by providing more support for the prioritized genes identified by our network-based approaches in Chapter 2, and our results from Chapter 3. Another refinement that would be very interesting involves the integration of prior data, in the form of whitelists (network edges that must exist, for example edges between transcription factors and the genes they regulate), in order to learn (possibly) more accurate network structures. Additionally, the structures of our Bayesian networks can be refined by bootstrapping or cross-validation. We were not able to perform these analyses due to cost and time constraints, and our results demonstrate the utility of our models as they stand, but performing these analyses may lead to more informative networks and, therefore, higher-confidence putatively causal genes.

Finally, some future directions are contingent on the generation of more data. For example, extending our multi-trait GWAS in Chapter 2 with more DO mice may result in more, as well as narrower, QTL for bone traits. Furthermore, the increased number of mice (of newer generations) in the GWAS cohort will allow us to further dissect complex loci, such the Chr. 1 locus. We are currently also generating RNA-seq data from cultured osteoblasts and osteoclasts in the DO. Performing the same analyses that we performed on cortical bone RNA-seq data will provide a wealth of information regarding cell-type specific contributions to bone traits and BMD GWAS associations. Network-based analyses of these data can also be performed to further understand the coupling mechanisms between osteoblasts and osteoclasts. The most useful data to be generated for systems genetics analyses of bone, however, are bone-specific transcriptomic data in large cohorts. Once these data are generated, we will be better powered to identify genes underlying BMD GWAS, using the systems genetics approaches discussed herein.

In conclusion, while BMD GWASs have provided a wealth of information to the field of bone genetics, the focus on BMD as a trait and the inherent difficulty in identifying causal genes underlying GWAS associations have been major obstacles towards our understanding of the genetic underpinnings of bone development, function, and disease. In this work, we have contributed a data resource encompassing GWAS associations of 55 bone-related traits in mouse, and have demonstrated and applied analytical methodologies, utilizing systems genetics approaches, to identify putatively causal genes underlying BMD GWAS in humans. However, this work barely scratches the surface of our understanding of the genetics of bone. In order to advance our understanding of the genetics of bone, concerted, collaborative efforts across myriad disciplines must be undertaken to provide more relevant and diverse data, and to generate more sophisticated analytical, experimental and computational approaches for the effective utilization of these data. Given the current rapid and highly-collaborative progressions observed in the scientific community at-large, we are highly optimistic in the future advancements that will be achieved in the field of bone genetics. Appendix A

Supplementary Data

All supplementary data are available at:

https://zenodo.org/record/4876547

Supplementary Data 2.1 Phenotypes collected in the DO.

Supplementary Data 2.2 Raw measurements of the 55 mapped phenotypes.

Supplementary Data 2.3 Correlations of mapped phenotypes with bone strength.

Supplementary Data 2.4 Pairwise correlations of all 55 mapped phenotypes.

Supplementary Data 2.5 Results of bulk RNA-seq differential expression between sexes.

Supplementary Data 2.6 Results of bulk RNA-seq differential expression between individuals with high vs. low bone strength.

Supplementary Data 2.7 Module membership.

Supplementary Data 2.8 Modular Gene Ontology terms with a P-value ≤ 0.05 .

Supplementary Data 2.9 List of known bone genes.

Supplementary Data 2.10 Bayesian networks.

Supplementary Data 2.11 Significantly colocalizing genes.

Supplementary Data 2.12 Homologous human BANs with colocalizing eQTL.

Supplementary Data 2.13 Spearman correlations between WGCNA modules and phenotypes.

Supplementary Data 2.14 Mineralization readouts for bone marrow stromal cells exposed to osteogenic differentiation media *in vitro*.

Supplementary Data 2.15 Marker genes for scRNA-seq clusters.

Supplementary Data 2.16 Syntenic DO loci.

Supplementary Data 2.17 Variant effect predictor (VEP) output for 15 missense SNPs.

Supplementary Data 2.18 Mapped eQTL.

Supplementary Data 2.19 scRNA-seq cluster 1 genes.

Supplementary Data 2.20 Founder mice for CRISPR/Cas9 analysis.

Supplementary Data 2.21 Predicted amino acid sequences for Qsox1 deletion mutants.

Supplementary Data 2.22 Oligonucleotide sequences for CRISPR/Cas9 analysis.

Supplementary Data 3.1 2,156 protein-coding genes that are significant by TWAS.

Supplementary Data 3.2 1,182 protein-coding genes that are significant in the colocalization analysis.

Supplementary Data 3.3 The 512 protein-coding genes that are significant by both TWAS and colocalization.

Supplementary Data 3.4 Number of the 512 significantly colocalizing genes per GTEx tissue.

Supplementary Data 3.5 The known bone gene list.

Supplementary Data 3.6 The 66 genes that are significant by both TWAS and colocalization.

Supplementary Data 3.7 Gene Ontology enrichments for the 512 protein-coding genes that are significant by both TWAS and colocalization.

Supplementary Data 3.8 137 novel putatively causal BMD genes.

Appendix B

Supplemental Figures



Supplemental Figure 2.1 *Principal Component Analysis of bulk RNA-seq data. A)* Scree plot showing the percentage of explained variance for the first 10 principal components. **B**) Individuals in PC1 and PC2 space, colored by sex. **C**) Individuals in PC3 and PC4 space, colored by sex. **D**) Individuals in PC1 and PC2 space, colored by batch. **E**) Individuals in PC1 and PC2 space, colored by batch. **F**) Individuals in PC1 and PC2 space, colored by age (binarized, see Methods). **G**) Individuals in PC3 and PC4 space, colored by age (binarized, see Methods).



Supplemental Figure 2.2 Mineralization of bone marrow-derived stromal cells exposed to osteogenic differentiation media in vitro. During differentiation, cells from each

individual DO mouse were assessed for accumulated mineralization by IRDye 680 BoneTag Optical Probe incorporation. The final values for mineralization shown here were computed by subtracting the average number of fluorescent units recorded in designated background wells from the number of fluorescent units recorded in the sample wells. In the cultures from DO mouse #50, there was a much higher percentage of marrow adipogenic lineage precursor cells and a small number of osteoblasts. Consistent with this observation, mouse #50 also demonstrated high levels of marrow adiposity. This is likely the basis of the poor in vitro mineralization observed for the cultures from this mouse.



Supplemental Figure 2.3 Sex-specific UMAP visulization of single cell RNA-seq expression data on bone marrow stromal cells cultured in osteogenic differentiation media in vitro. Each point represents a cell. A) Sertad4 expression in cells from a male DO mouse. B) Glt8d2 expression in cells from female DO mice (N=4). C) Sertad4 expression in cells from a male DO mouse. D) Glt8d2 expression in cells from female DO mice (N=4). In all panels, the color scales indicate normalized gene expression values.



Supplemental Figure 2.4 Significant QTL associations. Twenty-eight mapped QTL exceeding permutation-based LOD score thresholds (alpha=0.05)



BMD GWAS associations in locus 1

Supplemental Figure 2.5 Overlap between BMD GWAS SNPs and QTL loci. A-J)Each panel corresponds to a QTL locus's syntenic human region. Panels a-j represent each of the 10 loci sequentially. Red circles represent BMD GWAS SNPs in the locus. The horizontal lines represent the genome-wide significance threshold ($P = 5 \times 10^{-8}$). Not all genes are shown.

Α



Chromosome 20 position (bp)



132

Chromosome 3 position (bp)



Chromosome 1 position (bp)



Е

Chromosome 1 position (bp)



Chromosome 16 position (bp)





Chromosome 3 position (bp)


Chromosome X position (bp)



Chromosome X position (bp)



Supplemental Figure 2.6 ML mapping in a replication cohort. The top panel shows allele effects for the DO founders for ML in an interval on chromosome 1 (Mbp). Y-axis units are best linear unbiased predictors (BLUPs). The bottom panel shows the QTL scan.

ML QTL replication



Supplemental Figure 3.1 LSBMD and FNBMD GWAS SNPs in the PPP6R3 locus. A) GWAS SNPs for LSBMD in the PPP6R3 locus. SNPs in red are significant PPP6R3 eQTL in thyroid tissue. The dashed line represents the genome-wide significance level (P-value= 5×10 -8). B) GWAS SNPs for FNBMD in the PPP6R3 locus. SNPs in red are significant PPP6R3 eQTL in thyroid tissue. The dashed line represents the genome-wide significance level (P-value= 5×10 -8). C) Mirroplot of LSBMD SNPs and PPP6R3 eQTL. SNPs are colored by their LD with rs10047483 (purple), the most significant PPP6R3 eQTL in thyroid. In the LSBMD panel, the most significant SNP is highlighted in purple, as rs10047483 was not assayed. D) Mirroplot of FNBMD SNPs and PPP6R3 eQTL. SNPs are colored by their LD with rs10047483 (purple), the most significant PPP6R3 eQTL in thyroid. In the FNBMD panel, the most significant SNP is highlighted in purple, as rs10047483 was not assayed.



Supplemental Figure 3.2 *Ppp6r3 functional validation. A*) Weight. B) Anterior-posterior (AP) femoral width. C) Medial-lateral (ML) femoral width. D) Femoral length (FL). E) Tissue mineral density (TMD), as measured by µCT. In all panels, least-square means are plotted. P-values are contrast P-values, adjusted for multiple comparisons. Asterisks represent significance (P<=0.05).

References

- 1. Grabowski, P. Physiology of bone. Endocr. Dev. 28, 33–55 (2015).
- Karsenty, G. & Oury, F. Biology without walls: the novel endocrinology of bone. *Annu. Rev. Physiol.* 74, 87–105 (2012).
- Riddle, R. C. & Clemens, T. L. Bone cell bioenergetics and skeletal energy homeostasis. *Physiol. Rev.* 97, 667–698 (2017).
- Black, D. M. & Rosen, C. J. Clinical Practice. Postmenopausal Osteoporosis. N. Engl. J. Med. 374, 254–262 (2016).
- Cummings, S. R. & Melton, L. J. Epidemiology and outcomes of osteoporotic fractures. Lancet 359, 1761–1767 (2002).
- Harvey, N., Dennison, E. & Cooper, C. Osteoporosis: impact on health and economics. Nat. Rev. Rheumatol. 6, 99–105 (2010).
- Burge, R. *et al.* Incidence and economic burden of osteoporosis-related fractures in the United States, 2005-2025. *J. Bone Miner. Res.* 22, 465–475 (2007).
- 8. Peacock, M. Genetics of Osteoporosis. *Endocr. Rev.* 23, 303–326 (2002).
- Morrison, N. A. *et al.* Prediction of bone density from vitamin D receptor alleles. *Nature* 367, 284–287 (1994).
- Ralston, S. H. & Uitterlinden, A. G. Genetics of osteoporosis. *Endocr. Rev.* 31, 629–662 (2010).
- Richards, J. B. *et al.* Collaborative meta-analysis: associations of 150 candidate genes with osteoporosis and osteoporotic fracture. *Ann. Intern. Med.* 151, 528–537 (2009).
- 12. Ioannidis, J. P. *et al.* Meta-analysis of genome-wide scans provides evidence for sex- and site-specific regulation of bone mass. *J. Bone Miner. Res.* 22, 173–183 (2007).

- Kiel, D. P. *et al.* Genome-wide association with bone mass and geometry in the Framingham Heart Study. *BMC Med. Genet.* 8 Suppl 1, S14 (2007).
- Bush, W. S. & Moore, J. H. Chapter 11: Genome-wide association studies. *PLoS Comput. Biol.* 8, e1002822 (2012).
- Richards, J. B., Zheng, H.-F. & Spector, T. D. Genetics of osteoporosis from genomewide association studies: advances and challenges. *Nat. Rev. Genet.* 13, 576–588 (2012).
- Estrada, K. *et al.* Genome-wide meta-analysis identifies 56 bone mineral density loci and reveals 14 loci associated with risk of fracture. *Nat. Genet.* 44, 491–501 (2012).
- Kemp, J. P. *et al.* Identification of 153 new loci associated with heel bone mineral density and functional involvement of GPC6 in osteoporosis. *Nat. Genet.* 49, 1468–1475 (2017).
- Sabik, O. L. & Farber, C. R. Using GWAS to identify novel therapeutic targets for osteoporosis. *Transl. Res.* 181, 15–26 (2017).
- Morris, J. A. *et al.* An atlas of genetic influences on osteoporosis in humans and mice. *Nat. Genet.* 51, 258–266 (2019).
- Cummings, S. R. *et al.* Appendicular bone density and age predict hip fracture in women. The Study of Osteoporotic Fractures Research Group. *JAMA* 263, 665–668 (1990).
- Osterhoff, G. *et al.* Bone mechanical properties and changes with osteoporosis. *Injury* 47
 Suppl 2, S11-20 (2016).
- Nielson, C. M. *et al.* Novel genetic variants associated with increased vertebral volumetric BMD, reduced vertebral fracture risk, and increased expression of SLC1A3 and EPHB2. *J. Bone Miner. Res.* **31**, 2085–2097 (2016).

- 23. Hsu, Y.-H. *et al.* An integration of genome-wide association study and gene expression profiling to prioritize the discovery of novel susceptibility Loci for osteoporosis-related traits. *PLoS Genet.* **6**, e1000977 (2010).
- 24. Prins, B. P. *et al.* Genome-wide analysis of health-related biomarkers in the UK Household Longitudinal Study reveals novel associations. *Sci. Rep.* **7**, (2017).
- Zhao, L.-J. *et al.* Genome-wide association study for femoral neck bone geometry. *J.* Bone Miner. Res. 25, 320–329 (2010).
- Timpson, N. J., Greenwood, C. M. T., Soranzo, N., Lawson, D. J. & Richards, J. B. Genetic architecture: the shape of the genetic contribution to human traits and disease. *Nat. Rev. Genet.* 19, 110–124 (2018).
- Boyle, E. A., Li, Y. I. & Pritchard, J. K. An expanded view of complex traits: From polygenic to omnigenic. *Cell* 169, 1177–1186 (2017).
- Liu, X., Li, Y. I. & Pritchard, J. K. Trans effects on gene expression can drive omnigenic inheritance. *Cell* 177, 1022-1034.e6 (2019).
- Manolio, T. A. *et al.* Finding the missing heritability of complex diseases. *Nature* 461, 747–753 (2009).
- Styrkarsdottir, U. *et al.* Nonsense mutation in the LGR4 gene is associated with several human diseases and other traits. *Nature* 497, 517–520 (2013).
- 31. Styrkarsdottir, U. *et al.* Two rare mutations in the COL1A2 gene associate with low bone mineral density and fractures in Iceland. *J. Bone Miner. Res.* **31**, 173–179 (2016).
- 32. Zheng, H. *et al.* Whole-genome sequencing identifies EN1 as a determinant of bone density and fracture. *Nature* **526**, 112–117 (2015).
- Rivadeneira, F. *et al.* Twenty bone-mineral-density loci identified by large-scale metaanalysis of genome-wide association studies. *Nat. Genet.* 41, 1199–1206 (2009).

- Nadeau, J. H. & Dudley, A. M. Genetics. Systems genetics. Science 331, 1015–1016 (2011).
- Civelek, M. & Lusis, A. J. Systems genetics approaches to understand complex traits. Nat. Rev. Genet. 15, 34–48 (2014).
- 36. Tak, Y. G. & Farnham, P. J. Making sense of GWAS: using epigenomics and genome engineering to understand the functional relevance of SNPs in non-coding regions of the human genome. *Epigenetics Chromatin* 8, 57 (2015).
- 37. Onengut-Gumuscu, S. *et al.* Fine mapping of type 1 diabetes susceptibility loci and evidence for colocalization of causal variants with lymphoid gene enhancers. *Nat. Genet.*47, 381–386 (2015).
- Farh, K. K.-H. *et al.* Genetic and epigenetic fine mapping of causal autoimmune disease variants. *Nature* 518, 337–343 (2015).
- Castaldi, P. J. *et al.* Genome-wide association identifies regulatory Loci associated with distinct local histogram emphysema patterns. *Am. J. Respir. Crit. Care Med.* **190**, 399–409 (2014).
- Hazelett, D. J. *et al.* Comprehensive functional annotation of 77 prostate cancer risk loci. *PLoS Genet.* **10**, e1004102 (2014).
- 41. Chen, X.-F. *et al.* An osteoporosis risk SNP at 1p36.12 acts as an allele-specific enhancer to modulate LINC00339 expression via long-range loop formation. *Am. J. Hum. Genet.* 102, 776–793 (2018).
- 42. The ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).

- Ito, Y. *et al.* Cdc42 regulates bone modeling and remodeling in mice by modulating RANKL/M-CSF signaling and osteoclast polarization. *J. Clin. Invest.* 120, 1981–1993 (2010).
- 44. Teitelbaum, S. L. & Ross, F. P. Genetic regulation of osteoclast development and function. *Nat. Rev. Genet.* **4**, 638–649 (2003).
- Zhu, D.-L. *et al.* Multiple functional variants at 13q14 risk locus for osteoporosis regulate RANKL expression through long-range super-enhancer. *J. Bone Miner. Res.* 33, 1335–1346 (2018).
- Cookson, W., Liang, L., Abecasis, G., Moffatt, M. & Lathrop, M. Mapping complex disease traits with global gene expression. *Nat. Rev. Genet.* 10, 184–194 (2009).
- Farber, C. R. & Lusis, A. J. Integrating global gene expression analysis and genetics.
 Adv. Genet. 60, 571–601 (2008).
- Rockman, M. V. & Kruglyak, L. Genetics of global gene expression. *Nat. Rev. Genet.* 7, 862–872 (2006).
- GTEx Consortium. The Genotype-Tissue Expression (GTEx) project. Nat. Genet. 45, 580–585 (2013).
- 50. GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science* **369**, 1318–1330 (2020).
- GTEx Consortium *et al.* Genetic effects on gene expression across human tissues. Nature 550, 204–213 (2017).
- 52. Bonewald, L. F. The amazing osteocyte. J. Bone Miner. Res. 26, 229–238 (2011).
- Reppe, S. *et al.* Eight genes are highly associated with BMD variation in postmenopausal Caucasian women. *Bone* 46, 604–612 (2010).

- Grundberg, E. *et al.* Population genomics in a disease targeted primary cell model. *Genome Res.* 19, 1942–1952 (2009).
- 55. Mullin, B. H. *et al.* Expression quantitative trait locus study of bone mineral density GWAS variants in human osteoclasts. *J. Bone Miner. Res.* **33**, 1044–1051 (2018).
- Giambartolomei, C. *et al.* Bayesian Test for Colocalisation between Pairs of Genetic Association Studies Using Summary Statistics. *PLoS Genetics* vol. 10 e1004383 (2014).
- Wen, X., Pique-Regi, R. & Luca, F. Integrating molecular QTL data into genome-wide genetic association analysis: Probabilistic assessment of enrichment and colocalization. *PLoS Genet.* 13, e1006646 (2017).
- Hormozdiari, F. et al. Colocalization of GWAS and eQTL signals detects target genes.
 Am. J. Hum. Genet. 99, 1245–1260 (2016).
- Calabrese, G. M. *et al.* Integrating GWAS and Co-expression Network Data Identifies Bone Mineral Density Genes SPTBN1 and MARK3 and an Osteoblast Functional Module. *Cell Syst* 4, 46-59.e4 (2017).
- Zhang, Q. *et al.* Genomic variants within chromosome 14q32.32 regulate bone mass through MARK3 signaling in osteoblasts. *J. Clin. Invest.* 131, (2021).
- 61. Hukku, A. *et al.* Probabilistic colocalization of genetic variants from complex and molecular traits: promise and limitations. *Am. J. Hum. Genet.* **108**, 25–35 (2021).
- 62. Cano-Gamez, E. & Trynka, G. From GWAS to function: Using functional genomics to identify the mechanisms underlying complex diseases. *Front. Genet.* **11**, 424 (2020).
- 63. Aten, J. E., Fuller, T. F., Lusis, A. J. & Horvath, S. Using genetic markers to orient the edges in quantitative trait networks: the NEO software. *BMC Syst. Biol.* **2**, 34 (2008).
- Hemani, G., Bowden, J. & Davey Smith, G. Evaluating the potential role of pleiotropy in Mendelian randomization studies. *Hum. Mol. Genet.* 27, R195–R208 (2018).

- 65. Zhu, Z. *et al.* Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nat. Genet.* **48**, 481–487 (2016).
- 66. Gusev, A. *et al.* Integrative approaches for large-scale transcriptome-wide association studies. *Nat. Genet.* **48**, 245–252 (2016).
- 67. Gamazon, E. R. *et al.* A gene-based association method for mapping traits using reference transcriptome data. *Nat. Genet.* **47**, 1091–1098 (2015).
- Barbeira, A. N. *et al.* Exploring the phenotypic consequences of tissue specific gene expression variation inferred from GWAS summary statistics. *Nat. Commun.* 9, 1825 (2018).
- 69. Barbeira, A. N. *et al.* Integrating predicted transcriptome from multiple tissues improves association detection. *PLoS Genet.* **15**, e1007889 (2019).
- Wainberg, M. *et al.* Opportunities and challenges for transcriptome-wide association studies. *Nat. Genet.* 51, 592–599 (2019).
- Ma, M. *et al.* Integrating transcriptome-wide association study and mRNA expression profiling identifies novel genes associated with bone mineral density. *Osteoporos. Int.* 30, 1521–1528 (2019).
- 72. Grundberg, E. *et al.* Mapping cis- and trans-regulatory effects across multiple tissues in twins. *Nat. Genet.* **44**, 1084–1089 (2012).
- 73. Pividori, M. *et al.* PhenomeXcan: Mapping the genome to the phenome through the transcriptome. *Sci. Adv.* **6**, eaba2083 (2020).
- 74. Yin, P. *et al.* Integrating genome-wide association and transcriptome predicted model identify novel target genes with osteoporosis. *bioRxiv* (2019) doi:10.1101/771543.
- 75. Nelson, M. R. *et al.* The support of human genetic evidence for approved drug indications. *Nat. Genet.* **47**, 856–860 (2015).

- Plenge, R. M., Scolnick, E. M. & Altshuler, D. Validating therapeutic targets through human genetics. *Nat. Rev. Drug Discov.* 12, 581–594 (2013).
- 77. Barabási, A.-L., Gulbahce, N. & Loscalzo, J. Network medicine: a network-based approach to human disease. *Nat. Rev. Genet.* **12**, 56–68 (2011).
- Needham, C. J., Bradford, J. R., Bulpitt, A. J. & Westhead, D. R. A primer on learning in Bayesian networks for computational biology. *PLoS Comput. Biol.* 3, e129 (2007).
- 79. Carey, H. A. *et al.* Enhancer variants reveal a conserved transcription factor network governed by PU.1 during osteoclast differentiation. *Bone Res.* **6**, (2018).
- Langfelder, P. & Horvath, S. WGCNA: an R package for weighted correlation network analysis. *BMC Bioinformatics* 9, 559 (2008).
- Huang, R., Wallqvist, A. & Covell, D. G. Comprehensive analysis of pathway or functionally related gene expression in the National Cancer Institute's anticancer screen. *Genomics* 87, 315–328 (2006).
- Kaake, R. M. *et al.* A new in vivo cross-linking mass spectrometry platform to define protein-protein interactions in living cells. *Mol. Cell. Proteomics* 13, 3533–3543 (2014).
- Farber, C. R. Identification of a gene module associated with BMD through the integration of network analysis and genome-wide association data. *J. Bone Miner. Res.* 25, 2359–2367 (2010).
- 84. Langfelder, P., Mischel, P. S. & Horvath, S. When is hub gene selection better than standard meta-analysis? *PLoS One* **8**, e61505 (2013).
- Bennett, B. J. *et al.* A high-resolution association mapping panel for the dissection of complex traits in mice. *Genome Res.* 20, 281–290 (2010).
- Calabrese, G. *et al.* Systems genetic analysis of osteoblast-lineage cells. *PLoS Genet.* 8, e1003150 (2012).

- 87. Yao, W. *et al.* Overexpression of secreted frizzled-related protein 1 inhibits bone formation and attenuates parathyroid hormone bone anabolic effects. *J. Bone Miner. Res.* 25, 190–199 (2010).
- Chen, Y.-C. *et al.* Integrative analysis of genomics and transcriptome data to identify potential functional genes of BMDs in females. *J. Bone Miner. Res.* **31**, 1041–1049 (2016).
- 89. Leiserson, M. D. M., Eldridge, J. V., Ramachandran, S. & Raphael, B. J. Network analysis of GWAS data. *Curr. Opin. Genet. Dev.* **23**, 602–610 (2013).
- Pearl, J. Bayesian networks: A model of self-activated memory for evidential reasoning. in *Proceedings of the 7th Conference of the Cognitive Science Society* 329–334 (1985).
- Zhang, B. *et al.* Integrated systems approach identifies genetic nodes and networks in late-onset Alzheimer's disease. *Cell* 153, 707–720 (2013).
- 92. Mäkinen, V.-P. *et al.* Integrative genomics reveals novel molecular pathways and gene networks for coronary artery disease. *PLoS Genet.* **10**, e1004502 (2014).
- 93. Su, C., Andrew, A., Karagas, M. R. & Borsuk, M. E. Using Bayesian networks to discover relations between genes, environment, and disease. *BioData Min.* 6, (2013).
- Zhang, B., Tran, L., Emilsson, V. & Zhu, J. Characterization of Genetic Networks Associated with Alzheimer's Disease. *Methods Mol. Biol.* 1303, 459–477 (2016).
- 95. Elefteriou, F. & Yang, X. Genetic mouse models for bone studies--strengths and limitations. *Bone* **49**, 1242–1254 (2011).
- Jilka, R. L. The relevance of mouse models for investigating age-related bone loss in humans. J. Gerontol. A Biol. Sci. Med. Sci. 68, 1209–1217 (2013).
- Youlten, S. E. & Baldock, P. A. Using mouse genetics to understand human skeletal disease. *Bone* 126, 27–36 (2019).

- Beamer, W. G., Donahue, L. R., Rosen, C. J. & Baylink, D. J. Genetic variability in adult bone density among inbred strains of mice. *Bone* 18, 397–403 (1996).
- 99. Beamer, W. G., Donahue, L. R. & Rosen, C. J. Genetics and bone. Using the mouse to understand man. J. Musculoskelet. Neuronal Interact. 2, 225–231 (2002).
- 100. Shimizu, M. *et al.* Identification of peak bone mass QTL in a spontaneously osteoporotic mouse strain. *Mamm. Genome* **10**, 81–87 (1999).
- 101. Solberg Woods, L. C. QTL mapping in outbred populations: successes and challenges. *Physiol. Genomics* 46, 81–90 (2014).
- 102. Farber, C. R. *et al.* Mouse genome-wide association and systems genetics identify Asxl2 as a regulator of bone mineral density and osteoclastogenesis. *PLoS Genet.* 7, e1002038 (2011).
- 103. Levy, R., Mott, R. F., Iraqi, F. A. & Gabet, Y. Collaborative cross mice in a genetic association study reveal new candidate genes for bone microarchitecture. *BMC Genomics* 16, 1013 (2015).
- 104. Churchill, G. A. *et al.* The Collaborative Cross, a community resource for the genetic analysis of complex traits. *Nat. Genet.* **36**, 1133–1137 (2004).
- 105. Churchill, G. A., Gatti, D. M., Munger, S. C. & Svenson, K. L. The Diversity Outbred mouse population. *Mamm. Genome* 23, 713–718 (2012).
- 106. Levy, R. *et al.* A genome-wide association study in mice reveals a role for Rhbdf2 in skeletal homeostasis. *Sci. Rep.* **10**, 3286 (2020).
- 107. Jepsen, K. J. *et al.* Phenotypic integration of skeletal traits during growth buffers genetic variants affecting the slenderness of femora in inbred mouse strains. *Mamm. Genome* 20, 21–33 (2009).

- 108. Jepsen, K. J. Systems analysis of bone. Wiley Interdiscip. Rev. Syst. Biol. Med. 1, 73–88 (2009).
- 109. Freudenthal, B., Logan, J., Croucher, P. I., Williams, G. R. & Bassett, J. H. D. Rapid phenotyping of knockout mice to identify genetic determinants of bone strength. *J. Endocrinol.* 231, R31–R46 (2016).
- 110. Swan, A. L. *et al.* Mouse mutant phenotyping at scale reveals novel genes controlling bone mineral density. *PLoS Genet.* **16**, e1009190 (2020).
- 111. Blake, J. A. *et al.* Mouse Genome Database (MGD): Knowledgebase for mouse-human comparative biology. *Nucleic Acids Res.* **49**, D981–D987 (2021).
- 112. Kanis, J. A. et al. Assessment of fracture risk. Osteoporos. Int. 16, 581-589 (2005).
- 113. Stone, K. L. *et al.* BMD at multiple sites and risk of fracture of multiple types: long-term results from the Study of Osteoporotic Fractures. *J. Bone Miner. Res.* 18, 1947–1954 (2003).
- 114. Flint, J. & Eskin, E. Genome-wide association studies in mice. *Nat. Rev. Genet.* 13, 807–817 (2012).
- 115. Rat Genome Sequencing and Mapping Consortium *et al.* Combined sequence-based and genetic mapping analysis of complex traits in outbred rats. *Nat. Genet.* **45**, 767–775 (2013).
- 116. Al-Barghouthi, B. M. & Farber, C. R. Dissecting the Genetics of Osteoporosis using Systems Approaches. *Trends Genet.* **35**, 55–67 (2019).
- 117. Dufresne, T. E., Chmielewski, P. A., Manhart, M. D., Johnson, T. D. & Borah, B. Risedronate preserves bone architecture in early postmenopausal women in 1 year as measured by three-dimensional microcomputed tomography. *Calcif. Tissue Int.* **73**, 423– 432 (2003).

- 118. Cummings, S. R. *et al.* Improvement in spine bone density and reduction in risk of vertebral fractures during treatment with antiresorptive drugs. *Am. J. Med.* **112**, 281–289 (2002).
- 119. Lochmüller, E.-M. *et al.* Correlation of Femoral and Lumbar DXA and Calcaneal Ultrasound, Measured In Situ with Intact Soft Tissues, with the In Vitro Failure Loads of the Proximal Femur. *Osteoporosis International* vol. 8 591–598 (1998).
- 120. Melton, L. J., III, Chrischilles, E. A., Cooper, C., Lane, A. W. & Riggs, B. L. How Many Women Have Osteoporosis? *J. Bone Miner. Res.* **20**, 886–892 (2005).
- 121. Logan, R. W., Robledo, R. F. & Recla, J. M. High-precision genetic mapping of behavioral traits in the diversity outbred mouse population. *Genes Brain Behav.* (2013).
- 122. Svenson, K. L. *et al.* High-resolution genetic mapping using the Mouse Diversity outbred population. *Genetics* **190**, 437–447 (2012).
- 123. Morgan, A. P. *et al.* The Mouse Universal Genotyping Array: From Substrains to Subspecies. *G3* 6, 263–279 (2015).
- 124. Karasik, D. *et al.* Heritability and Genetic Correlations for Bone Microarchitecture: The Framingham Study Families. *J. Bone Miner. Res.* **32**, 106–114 (2017).
- 125. Ng, A. H. M., Wang, S. X., Turner, C. H., Beamer, W. G. & Grynpas, M. D. Bone quality and bone strength in BXH recombinant inbred mice. *Calcif. Tissue Int.* 81, 215– 223 (2007).
- 126. Turner, C. H. *et al.* Variation in bone biomechanical properties, microstructure, and density in BXH recombinant inbred mice. *J. Bone Miner. Res.* **16**, 206–213 (2001).
- 127. Schlecht, S. H. & Jepsen, K. J. Functional integration of skeletal traits: an intraskeletal assessment of bone size, mineralization, and volume covariance. *Bone* 56, 127–138 (2013).

- 128. Szweras, M. *et al.* alpha 2-HS glycoprotein/fetuin, a transforming growth factorbeta/bone morphogenetic protein antagonist, regulates postnatal bone growth and remodeling. *J. Biol. Chem.* 277, 19991–19997 (2002).
- 129. Yeon, J.-T., Choi, S.-W. & Kim, S. H. Arginase 1 is a negative regulator of osteoclast differentiation. *Amino Acids* **48**, 559–565 (2016).
- 130. Sabik, O. L., Calabrese, G. M., Taleghani, E., Ackert-Bicknell, C. L. & Farber, C. R. Identification of a core module for bone mineral density through the integration of a co-expression network and GWAS data. *Cell Rep.* **32**, 108145 (2020).
- 131. Zhang, B. & Horvath, S. A general framework for weighted gene co-expression network analysis. *Stat. Appl. Genet. Mol. Biol.* **4**, Article17 (2005).
- 132. Watson, C. T. *et al.* Integrative transcriptomic analysis reveals key drivers of acute peanut allergic reactions. *Nat. Commun.* **8**, 1943 (2017).
- 133. Huan, T. *et al.* Integrative network analysis reveals molecular mechanisms of blood pressure regulation. *Mol. Syst. Biol.* **11**, 799 (2015).
- 134. Wang, I.-M. *et al.* Systems analysis of eleven rodent disease models reveals an inflammatome signature and key drivers. *Mol. Syst. Biol.* **8**, 594 (2012).
- 135. Aguet, F. *et al.* The GTEx Consortium atlas of genetic regulatory effects across human tissues. *bioRxiv* 787903 (2019) doi:10.1101/787903.
- 136. Aguet, F. *et al.* Local genetic effects on gene expression across 44 human tissues. *bioRxiv*074450 (2016) doi:10.1101/074450.
- 137. Nakashima, K. *et al.* The novel zinc finger-containing transcription factor osterix is required for osteoblast differentiation and bone formation. *Cell* **108**, 17–29 (2002).
- 138. Balemans, W. *et al.* Increased bone density in sclerosteosis is due to the deficiency of a novel secreted protein (SOST). *Hum. Mol. Genet.* **10**, 537–543 (2001).

- 139. Brunkow, M. E. *et al.* Bone dysplasia sclerosteosis results from loss of the SOST gene product, a novel cystine knot-containing protein. *Am. J. Hum. Genet.* **68**, 577–589 (2001).
- 140. Gong, Y. *et al.* LDL receptor-related protein 5 (LRP5) affects bone accrual and eye development. *Cell* **107**, 513–523 (2001).
- 141. Little, R. D. *et al.* A mutation in the LDL receptor-related protein 5 gene results in the autosomal dominant high-bone-mass trait. *Am. J. Hum. Genet.* **70**, 11–19 (2002).
- 142. Boyden, L. M. et al. High bone density due to a mutation in LDL-receptor-related protein 5. N. Engl. J. Med. 346, 1513–1521 (2002).
- 143. Kong, Y. Y. *et al.* OPGL is a key regulator of osteoclastogenesis, lymphocyte development and lymph-node organogenesis. *Nature* **397**, 315–323 (1999).
- 144. Wong, B. R. et al. TRANCE (Tumor Necrosis Factor [TNF]-related Activation-induced Cytokine), a New TNF Family Member Predominantly Expressed in T cells, Is a Dendritic Cell–specific Survival Factor. Journal of Experimental Medicine vol. 186 2075– 2080 (1997).
- 145. Yasuda, H. *et al.* Osteoclast differentiation factor is a ligand for osteoprotegerin/osteoclastogenesis-inhibitory factor and is identical to TRANCE/RANKL. *Proc. Natl. Acad. Sci. U. S. A.* **95**, 3597–3602 (1998).
- 146. Lacey, D. L. *et al.* Osteoprotegerin ligand is a cytokine that regulates osteoclast differentiation and activation. *Cell* **93**, 165–176 (1998).
- 147. Anderson, D. M. *et al.* A homologue of the TNF receptor and its ligand enhance T-cell growth and dendritic-cell function. *Nature* **390**, 175–179 (1997).
- 148. Wong, B. R. *et al.* The TRAF family of signal transducers mediates NF-kappaB activation by the TRANCE receptor. *J. Biol. Chem.* **273**, 28355–28359 (1998).

- 149. Wu, J., Glimcher, L. H. & Aliprantis, A. O. HCO3-/Cl- anion exchanger SLC4A2 is required for proper osteoclast differentiation and function. *Proc. Natl. Acad. Sci. U. S. A.* 105, 16934–16939 (2008).
- 150. Duan, X., Yang, S., Zhang, L. & Yang, T. V-ATPases and osteoclasts: ambiguous future of V-ATPases inhibitors in osteoporosis. *Theranostics* **8**, 5379–5399 (2018).
- 151. Lattin, J. E. *et al.* Expression analysis of G Protein-Coupled Receptors in mouse macrophages. *Immunome Res.* **4**, 5 (2008).
- 152. Bennetts, J. S. *et al.* Evolutionary conservation and murine embryonic expression of the gene encoding the SERTA domain-containing protein CDCA4 (HEPP). *Gene* 374, 153– 165 (2006).
- 153. Zhan, Y. *et al.* Mechanism of the effect of glycosyltransferase GLT8D2 on fatty liver. *Lipids Health Dis.* **14**, 43 (2015).
- 154. Movérare-Skrtic, S. *et al.* Osteoblast-derived WNT16 represses osteoclastogenesis and prevents cortical bone fragility fractures. *Nat. Med.* **20**, 1279–1288 (2014).
- 155. Takeshita, S., Kikuno, R., Tezuka, K. & Amann, E. Osteoblast-specific factor 2: cloning of a putative bone adhesion protein with homology with the insect protein fasciclin I. *Biochem.* J 294 (Pt 1), 271–278 (1993).
- 156. Horiuchi, K. *et al.* Identification and characterization of a novel protein, periostin, with restricted expression to periosteum and periodontal ligament and increased expression by transforming growth factor beta. *J. Bone Miner. Res.* **14**, 1239–1249 (1999).
- 157. Izu, Y., Ezura, Y., Koch, M., Birk, D. E. & Noda, M. Collagens VI and XII form complexes mediating osteoblast interactions during osteogenesis. *Cell Tissue Res.* 364, 623–635 (2016).

- 158. Amiri, N. & Christians, J. K. PAPP-A2 expression by osteoblasts is required for normal postnatal growth in mice. *Growth Horm. IGF Res.* **25**, 274–280 (2015).
- 159. Wilm, B., Dahl, E., Peters, H., Balling, R. & Imai, K. Targeted disruption of Pax1 defines its null phenotype and proves haploinsufficiency. *Proc. Natl. Acad. Sci. U. S. A.* 95, 8692–8697 (1998).
- 160. Kimura, H., Akiyama, H., Nakamura, T. & de Crombrugghe, B. Tenascin-W inhibits proliferation and differentiation of preosteoblasts during endochondral bone formation. *Biochem. Biophys. Res. Commun.* 356, 935–941 (2007).
- 161. Koscielny, G. *et al.* The International Mouse Phenotyping Consortium Web Portal, a unified point of access for knockout mice and related phenotyping data. *Nucleic Acids Res.* 42, D802-9 (2014).
- 162. Yalcin, B., Flint, J. & Mott, R. Using Progenitor Strain Information to Identify Quantitative Trait Nucleotides in Outbred Mice. *Genetics* vol. 171 673–681 (2005).
- 163. Shorter, J. R. *et al.* Quantitative trait mapping in Diversity Outbred mice identifies two genomic regions associated with heart size. *Mamm. Genome* **29**, 80–89 (2018).
- 164. Ilani, T. *et al.* A secreted disulfide catalyst controls extracellular matrix composition and function. *Science* **341**, 74–76 (2013).
- 165. Huybrechts, Y., Mortier, G., Boudin, E. & Van Hul, W. WNT Signaling and Bone: Lessons From Skeletal Dysplasias and Disorders. *Front. Endocrinol.* **11**, 165 (2020).
- 166. Teufel, S. & Hartmann, C. Wnt-signaling in skeletal development. *Curr. Top. Dev. Biol.*133, 235–279 (2019).
- 167. Tong, W. et al. Wnt16 attenuates osteoarthritis progression through a PCP/JNKmTORC1-PTHrP cascade. Ann. Rheum. Dis. 78, 551–561 (2019).

- 168. Bonnet, N., Garnero, P. & Ferrari, S. Periostin action in bone. Mol. Cell. Endocrinol. 432, 75–82 (2016).
- 169. Rajpal, G. & Arvan, P. Chapter 236 Disulfide Bond Formation. in Handbook of Biologically Active Peptides (Second Edition) (ed. Kastin, A. J.) 1721–1729 (Academic Press, 2013).
- 170. Bulleid, N. J. & Ellgaard, L. Multiple ways to make disulfides. *Trends in Biochemical Sciences* vol. 36 485–492 (2011).
- 171. Feldman, T. *et al.* Inhibition of fibroblast secreted QSOX1 perturbs extracellular matrix in the tumor microenvironment and decreases tumor growth and metastasis in murine cancer models. *Oncotarget* **11**, 386–398 (2020).
- 172. Hanavan, P. D. *et al.* Ebselen inhibits QSOX1 enzymatic activity and suppresses invasion of pancreatic and renal cancer cell lines. *Oncotarget* **6**, 18418–18428 (2015).
- 173. Bouxsein, M. L. *et al.* Guidelines for assessment of bone microstructure in rodents using micro--computed tomography. *J. Bone Miner.* Res. 25, 1468–1486 (2010).
- 174. Andrews, S. & Others. FastQC: a quality control tool for high throughput sequence data. (2010).
- 175. Ewels, P., Magnusson, M., Lundin, S. & Käller, M. MultiQC: summarize analysis results for multiple tools and samples in a single report. *Bioinformatics* **32**, 3047–3048 (2016).
- 176. Kim, D., Paggi, J. M., Park, C., Bennett, C. & Salzberg, S. L. Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* 37, 907– 915 (2019).
- 177. Pertea, M. *et al.* StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat. Biotechnol.* **33**, 290–295 (2015).

- 178. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
- 179. Lê, S., Josse, J. & Husson, F. FactoMineR: AnRPackage for multivariate analysis. *J. Stat. Softw.* **25**, (2008).
- 180. Zhu, A., Ibrahim, J. G. & Love, M. I. Heavy-tailed prior distributions for sequence count data: removing the noise and preserving large differences. *Bioinformatics* 35, 2084– 2092 (2019).
- 181. Team, R. C. & Others. R: A language and environment for statistical computing. (2013).
- 182. Broman, K. W. et al. R/qtl2: Software for Mapping Quantitative Trait Loci with High-Dimensional Data and Multiparent Populations. Genetics 211, 495–502 (2019).
- 183. Morgan, A. P. argyle: An R Package for Analysis of Illumina Genotyping Arrays. G3 6, 281–286 (2015).
- 184. Broman, K. W., Gatti, D. M., Svenson, K. L., Sen, S. & Churchill, G. A. Cleaning Genotype Data from Diversity Outbred Mice. G3 9, 1571–1579 (2019).
- 185. Leek, J. T., Johnson, W. E., Parker, H. S., Jaffe, A. E. & Storey, J. D. The sva package for removing batch effects and other unwanted variation in high-throughput experiments. *Bioinformatics* 28, 882–883 (2012).
- 186. Langfelder, P. & Horvath, S. Fast R Functions for Robust Correlations and Hierarchical Clustering. J. Stat. Softw. 46, (2012).
- 187. Scutari, M. Learning Bayesian Networks with thebnlearnRPackage. *Journal of Statistical Software* vol. 35 (2010).
- 188. Ashburner, M. et al. Gene Ontology: tool for the unification of biology. Nature Genetics vol. 25 25–29 (2000).

- 189. Blake, J. A. *et al.* The Mouse Genome Database (MGD): premier model organism resource for mammalian genomics and genetics. *Nucleic Acids Res.* **39**, D842-8 (2011).
- 190. Carbon, S. *et al.* AmiGO: online access to ontology and annotation data. *Bioinformatics*25, 288–289 (2009).
- 191. Csardi, G., Nepusz, T. & Others. The igraph software package for complex network research. *InterJournal, complex systems* **1695**, 1–9 (2006).
- 192. Alexa, A. Rahnenfuhrer J. topGO: enrichment analysis for Gene Ontology. 2010. R *package version* **2**, 45 (2017).
- 193. Dickinson, M. E. *et al.* High-throughput discovery of novel developmental phenotypes. *Nature* **537**, 508–514 (2016).
- 194. Kurbatova, N., Karp, N., Mason, J. & Haselimashhadi, H. PhenStat: statistical analysis of phenotypic data. *R package version* **2**, (2015).
- 195. West, B. T., Welch, K. B. & Galecki, A. T. Linear Mixed Models: A Practical Guide Using Statistical Software, Second Edition. (CRC Press, 2014).
- 196. Yang, J. et al. FTO genotype is associated with phenotypic variability of body mass index. Nature **490**, 267–272 (2012).
- 197. Stegle, O., Parts, L., Piipari, M., Winn, J. & Durbin, R. Using probabilistic estimation of expression residuals (PEER) to obtain increased power and interpretability of gene expression analyses. *Nat. Protoc.* **7**, 500–507 (2012).
- 198. Hinrichs, A. S. *et al.* The UCSC Genome Browser Database: update 2006. *Nucleic Acids Res.* **34**, D590-8 (2006).
- 199. Lawrence, M. et al. Software for Computing and Annotating Genomic Ranges. PLoS Computational Biology vol. 9 e1003118 (2013).
- 200. McLaren, W. et al. The Ensembl Variant Effect Predictor. Genome Biol. 17, 122 (2016).

- 201. Mesner, L. D. *et al.* Mouse genome-wide association and systems genetics identifies Lhfp as a regulator of bone mass. *PLoS Genet.* **15**, e1008123 (2019).
- 202. Fox, J. & Weisberg, S. An R companion to applied regression. (SAGE Publications, 2018).
- 203. Russell, L. emmeans: Estimated Marginal Means, aka Least-Squares Means. R package version 1.4. (2019).
- 204. Israel, B. A., Jiang, L., Gannon, S. A. & Thorpe, C. Disulfide bond generation in mammalian blood serum: detection and purification of quiescin-sulfhydryl oxidase. *Free Radic. Biol. Med.* 69, 129–135 (2014).
- 205. Hanna, H., Mir, L. M. & Andre, F. M. In vitro osteoblastic differentiation of mesenchymal stem cells generates cell layers with distinct properties. *Stem Cell Res. Ther.*9, 203 (2018).
- 206. Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. Bioinformatics 29, 15–21 (2013).
- 207. Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat. Biotechnol.* 36, 411–420 (2018).
- 208. Stuart, T. et al. Comprehensive Integration of Single-Cell Data. Cell 177, 1888-1902.e21 (2019).
- 209. Tirosh, I. *et al.* Dissecting the multicellular ecosystem of metastatic melanoma by singlecell RNA-seq. *Science* **352**, 189–196 (2016).
- 210. Al-Barghouthi, B. *et al.* Systems genetics analyses in Diversity Outbred mice inform human bone mineral density GWAS and identify Qsox1 as a novel determinant of bone strength. (2020) doi:10.5281/ZENODO.4265417.

- 211. Al-Barghouthi, B. basel-maher/DO_project: (Zenodo, 2021). doi:10.5281/ZENODO.4666326.
- 212. Miller, P. D., Zapalowski, C., Kulak, C. A. & Bilezikian, J. P. Bone densitometry: the best way to detect osteoporosis and to monitor therapy. *J. Clin. Endocrinol. Metab.* 84, 1867–1871 (1999).
- 213. Rocha-Braz, M. G. M. & Ferraz-de-Souza, B. Genetics of osteoporosis: searching for candidate genes for bone fragility. *Arch. Endocrinol. Metab.* **60**, 391–401 (2016).
- 214. Giral, H., Landmesser, U. & Kratzer, A. Into the wild: GWAS exploration of noncoding RNAs. *Front. Cardiovasc. Med.* **5**, 181 (2018).
- 215. Edwards, S. L., Beesley, J., French, J. D. & Dunning, A. M. Beyond GWASs: illuminating the dark road from association to function. *Am. J. Hum. Genet.* **93**, 779–797 (2013).
- 216. Bhattacharya, A. *et al.* A framework for transcriptome-wide association studies in breast cancer in diverse study populations. *Genome Biol.* **21**, 42 (2020).
- 217. Thom, C. S. & Voight, B. F. Genetic colocalization atlas points to common regulatory sites and genes for hematopoietic traits and hematopoietic contributions to disease phenotypes. *BMC Med. Genomics* **13**, 89 (2020).
- 218. Nica, A. C. & Dermitzakis, E. T. Expression quantitative trait loci: present and future. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **368**, 20120362 (2013).
- 219. Fitzpatrick, L. A. Secondary causes of osteoporosis. Mayo Clin. Proc. 77, 453-468 (2002).
- 220. Mirza, F. & Canalis, E. Management of endocrine disease: Secondary osteoporosis: pathophysiology and management. *Eur. J. Endocrinol.* **173**, R131-51 (2015).
- 221. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).

- 222. Stefansson, B. & Brautigan, D. L. Protein phosphatase 6 subunit with conserved Sit4associated protein domain targets IkappaBepsilon. J. Biol. Chem. 281, 22624–22634 (2006).
- 223. Cristiano, L. PPP6R3 (protein phosphatase 6 regulatory subunit 3). Atlas Genet. Cytogenet. Oncol. Haematol. (2020) doi:10.4267/2042/70657.
- 224. Ziembik, M. A., Bender, T. P., Larner, J. M. & Brautigan, D. L. Functions of protein phosphatase-6 in NF-*μ*B signaling and in lymphocytes. *Biochem. Soc. Trans.* 45, 693–701 (2017).
- 225. Abu-Amer, Y. NF-*x*B signaling and bone resorption. Osteoporos. Int. 24, 2377–2386 (2013).
- 226. Mao, J. *et al.* Low-density lipoprotein receptor-related protein-5 binds to Axin and regulates the canonical Wnt signaling pathway. *Mol. Cell* **7**, 801–809 (2001).
- 227. Mizuguchi, T. *et al.* LRP5, low-density-lipoprotein-receptor-related protein 5, is a determinant for bone mineral density. *J. Hum. Genet.* **49**, 80–86 (2004).
- 228. Marques-Pinheiro, A. *et al.* Novel LRP5 gene mutation in a patient with osteoporosispseudoglioma syndrome. *Joint Bone Spine* **77**, 151–153 (2010).
- 229. van Meurs, J. B. J. *et al.* Large-scale analysis of association between LRP5 and LRP6 variants and osteoporosis. *JAMA* **299**, 1277–1290 (2008).
- 230. Brixen, K. et al. Polymorphisms in the low-density lipoprotein receptor-related protein 5 (LRP5) gene are associated with peak bone mass in non-sedentary men: results from the Odense androgen study. *Calcif. Tissue Int.* 81, 421–429 (2007).
- 231. Giroux, S., Elfassihi, L., Cardinal, G., Laflamme, N. & Rousseau, F. LRP5 coding polymorphisms influence the variation of peak bone mass in a normal population of French-Canadian women. *Bone* 40, 1299–1307 (2007).

- 232. Peters, J. E. *et al.* Insight into genotype-phenotype associations through eQTL mapping in multiple cell types in health and immune-mediated disease. *PLoS Genet.* **12**, e1005908 (2016).
- 233. Mizuno, A. & Okada, Y. Biological characterization of expression quantitative trait loci (eQTLs) showing tissue-specific opposite directional effects. *Eur. J. Hum. Genet.* 27, 1745–1756 (2019).
- 234. Berisa, T. & Pickrell, J. K. Approximately independent linkage disequilibrium blocks in human populations. *Bioinformatics* **32**, 283–285 (2016).
- 235. Wen, X. Effective QTL Discovery Incorporating Genomic Annotations. *bioRxiv* (2015) doi:10.1101/032003.
- 236. 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).
- 237. Durinck, S. *et al.* BioMart and Bioconductor: a powerful link between biological databases and microarray data analysis. *Bioinformatics* **21**, 3439–3440 (2005).
- 238. Alexa, A. & Rahnenfuhrer, J. topGO: Enrichment Analysis for Gene Ontology. R package version 2.40.0. (2020).
- 239. Myers, T. A., Chanock, S. J. & Machiela, M. J. LDlinkR: An R package for rapidly calculating linkage disequilibrium statistics in diverse populations. *Front. Genet.* **11**, 157 (2020).
- 240. Guergnon, J., Derewenda, U., Edelson, J. R. & Brautigan, D. L. Mapping of protein phosphatase-6 association with its SAPS domain regulatory subunit using a model of helical repeats. *BMC Biochem.* **10**, 24 (2009).

- 241. Beamer, W. G. *et al.* BMD regulation on mouse distal chromosome 1, candidate genes, and response to ovariectomy or dietary fat. *J. Bone Miner. Res.* **26**, 88–99 (2011).
- 242. Lenth, R. Emmeans: Estimated Marginal Means, aka Least-Squares Means. R package version 1.5.2.1. (2020).
- 243. Chamberlain, S. solrium: General Purpose R Interface to "Solr". R package version1.1.4. (2019).
- 244. Mashhadi, H. H. OpenStats: A Robust and Scalable Software Package for Reproducible Analysis of High-Throughput genotype-phenotype association. R package version 1.0.2. (2020).
- 245. Umans, B. D., Battle, A. & Gilad, Y. Where are the disease-associated eQTLs? *Trends Genet.* **37**, 109–124 (2021).
- 246. Mullin, B. H. *et al.* Characterisation of genetic regulatory effects for osteoporosis risk variants in human osteoclasts. *Genome Biol.* **21**, 80 (2020).
- 247. Griffiths, J. A., Scialdone, A. & Marioni, J. C. Using single-cell genomics to understand developmental processes and cell fate decisions. *Mol. Syst. Biol.* **14**, e8046 (2018).
- 248. Baca, S. C., Singler, C., Zacharia, S., Seo, J. H. & Morova, T. Genetic determinants of chromatin reveal prostate cancer risk mediated by context-dependent gene regulation. *bioRxiv* (2021).
- 249. Donovan, M. K. R., D'Antonio-Chronowska, A., D'Antonio, M. & Frazer, K. A. Cellular deconvolution of GTEx tissues powers discovery of disease and cell-type associated regulatory variants. *Nat. Commun.* **11**, 955 (2020).
- Claussnitzer, M. *et al.* FTO obesity variant circuitry and adipocyte browning in humans.
 N. *Engl. J. Med.* 373, 895–907 (2015).

- 251. Gong, Y. *et al.* A systematic dissection of human primary osteoblastsin vivoat single-cell resolution. *bioRxiv* (2020) doi:10.1101/2020.05.12.091975.
- 252. McDonald, M. M. *et al.* Osteoclasts recycle via osteomorphs during RANKL-stimulated bone resorption. *Cell* **184**, 1940 (2021).
- 253. Lindeboom, R. G. H., Regev, A. & Teichmann, S. A. Towards a Human Cell Atlas: Taking notes from the past. *Trends Genet.* (2021) doi:10.1016/j.tig.2021.03.007.
- 254. Hwang, B., Lee, J. H. & Bang, D. Single-cell RNA sequencing technologies and bioinformatics pipelines. *Exp. Mol. Med.* **50**, 1–14 (2018).
- 255. Buenrostro, J. D. *et al.* Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature* **523**, 486–490 (2015).
- 256. Avila Cobos, F., Alquicira-Hernandez, J., Powell, J. E., Mestdagh, P. & De Preter, K. Benchmarking of cell type deconvolution pipelines for transcriptomics data. *Nat. Commun.* 11, 5650 (2020).
- 257. Calderon, D. *et al.* Inferring relevant cell types for complex traits by using single-cell gene expression. *Am. J. Hum. Genet.* **101**, 686–699 (2017).
- 258. Watanabe, K., Umićević Mirkov, M., de Leeuw, C. A., van den Heuvel, M. P. & Posthuma, D. Genetic mapping of cell type specificity for complex traits. *Nat. Commun.* 10, 3222 (2019).
- 259. Sloan, C. A. et al. ENCODE data at the ENCODE portal. Nucleic Acids Res. 44, D726-32 (2016).
- 260. Cho, Y. S. *et al.* A large-scale genome-wide association study of Asian populations uncovers genetic factors influencing eight quantitative traits. *Nat. Genet.* 41, 527–534 (2009).

- 261. Kung, A. W. C. *et al.* Association of JAG1 with bone mineral density and osteoporotic fractures: a genome-wide association study and follow-up replication studies. *Am. J. Hum. Genet.* 86, 229–239 (2010).
- 262. Kou, I. *et al.* Common variants in a novel gene, FONG on chromosome 2q33.1 confer risk of osteoporosis in Japanese. *PLoS One* **6**, e19641 (2011).
- 263. Choi, H. J. *et al.* Corrigendum to genome-wide association study in East Asians suggests UHMK1 as a novel bone mineral density susceptibility gene. *Bone* **106**, 211 (2018).
- 264. Taylor, K. C. *et al.* A genome-wide association study meta-analysis of clinical fracture in 10,012 African American women. *Bone Rep.* **5**, 233–242 (2016).
- 265. Koller, D. L. *et al.* Genome-wide association study of bone mineral density in premenopausal European-American women and replication in African-American women. *J. Clin. Endocrinol. Metab.* **95**, 1802–1809 (2010).
- 266. Villalobos-Comparán, M. *et al.* A pilot genome-wide association study in postmenopausal Mexican-Mestizo women implicates the RMND1/CCDC170 locus is associated with bone mineral density. *Int. J. Genomics* **2017**, 5831020 (2017).