

# Identification of Known and Novel Long Noncoding RNAs Potentially Responsible for the Effects of Bone Mineral Density (BMD) Genomewide Association Study (GWAS) Loci

Abdullah Abood,<sup>1,2</sup>  Larry Mesner,<sup>1,3</sup> Will Rosenow,<sup>1</sup> Basel M. Al-Barghouthi,<sup>1,2</sup> Nina Horowitz,<sup>4</sup> Elise F. Morgan,<sup>5</sup> Louis C. Gerstenfeld,<sup>4</sup>  and Charles R. Farber<sup>1,2,3</sup> 

<sup>1</sup>Center for Public Health Genomics, School of Medicine, University of Virginia, Charlottesville, VA, USA

<sup>2</sup>Department of Biochemistry and Molecular Genetics, School of Medicine, University of Virginia, Charlottesville, VA, USA

<sup>3</sup>Department of Public Health Sciences, School of Medicine, University of Virginia, Charlottesville, VA, USA

<sup>4</sup>Department of Orthopaedic Surgery, Boston University, Boston, MA, USA

<sup>5</sup>Department of Mechanical Engineering, Boston University, Boston, MA, USA

## ABSTRACT

Osteoporosis, characterized by low bone mineral density (BMD), is the most common complex disease affecting bone and constitutes a major societal health problem. Genome-wide association studies (GWASs) have identified over 1100 associations influencing BMD. It has been shown that perturbations to long noncoding RNAs (lncRNAs) influence BMD and the activities of bone cells; however, the extent to which lncRNAs are involved in the genetic regulation of BMD is unknown. Here, we combined the analysis of allelic imbalance (AI) in human acetabular bone fragments with a transcriptome-wide association study (TWAS) and expression quantitative trait loci (eQTL) colocalization analysis using data from the Genotype-Tissue Expression (GTEx) project to identify lncRNAs potentially responsible for GWAS associations. We identified 27 lncRNAs in bone that are located in proximity to a BMD GWAS association and harbor single-nucleotide polymorphisms (SNPs) demonstrating AI. Using GTEx data we identified an additional 31 lncRNAs whose expression was associated (false discovery rate [FDR] correction < 0.05) with BMD through TWAS and had a colocalizing eQTL (regional colocalization probability [RCP] > 0.1). The 58 lncRNAs are located in 43 BMD associations. To further support a causal role for the identified lncRNAs, we show that 23 of the 58 lncRNAs are differentially expressed as a function of osteoblast differentiation. Our approach identifies lncRNAs that are potentially responsible for BMD GWAS associations and suggest that lncRNAs play a role in the genetics of osteoporosis. © 2022 The Authors. *Journal of Bone and Mineral Research* published by Wiley Periodicals LLC on behalf of American Society for Bone and Mineral Research (ASBMR).

**KEY WORDS:** OSTEOPOROSIS; HUMAN ASSOCIATION STUDIES; OSTEOCYTES; OSTEOBLASTS

## Introduction

Osteoporosis is characterized by low bone mineral density (BMD) and deteriorated structural integrity that leads to an increased risk of fracture.<sup>(1,2)</sup> In the United States alone, 12 million individuals have been diagnosed with osteoporosis, contributing to over 2 million fractures per year.<sup>(3)</sup> This number is expected to nearly double by 2025, resulting in approximately \$26 billion in health care expenditures.<sup>(3)</sup>

BMD is one of the strongest predictors of fracture<sup>(4)</sup> and is a highly heritable quantitative trait ( $h^2 = 0.5-0.8$ ).<sup>(5-8)</sup> As a result, the majority of genomewide association studies (GWASs) conducted for osteoporosis have focused on BMD. The largest BMD GWAS performed to date used the UK BioBank ( $N \approx 420,000$ ) and identified 1103 associations influencing heel estimated BMD (eBMD).<sup>(9)</sup> One of the main challenges of BMD GWAS is that the majority (>90%) of associations implicate non-coding variants that lie in intronic or intergenic regions,

This is an open access article under the terms of the [Creative Commons Attribution](#) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

Received in original form November 4, 2021; revised form April 26, 2022; accepted June 4, 2022.

Address correspondence to: Charles R. Farber, Center for Public Health Genomics, University of Virginia, P.O. Box 800717, Charlottesville, VA 22908, USA.

E-mail: [crf2s@virginia.edu](mailto:crf2s@virginia.edu)

Additional Supporting Information may be found in the online version of this article.

*Journal of Bone and Mineral Research*, Vol. 37, No. 8, August 2022, pp 1500–1510.

DOI: 10.1002/jbmr.4622

© 2022 The Authors. *Journal of Bone and Mineral Research* published by Wiley Periodicals LLC on behalf of American Society for Bone and Mineral Research (ASBMR).

suggesting they have a role in gene regulation. This has made it difficult to pinpoint causal genes and highlights the need for follow-up studies.<sup>(10)</sup> In addition, few studies have systematically evaluated noncoding transcripts as potential causal genes.

The largest and most functionally diverse family of noncoding transcripts are long noncoding RNAs (lncRNAs). lncRNAs are transcripts longer than 200 nucleotides and have no coding potential.<sup>(11)</sup> The majority of lncRNAs share sequence features with protein-coding genes including a 3' poly-A tail, a 5' methyl cap, and an open reading frame.<sup>(12)</sup> However, their expression is low and heterogenous, and they show intermediate to high tissue specificity.<sup>(13)</sup> Aberrant expression of lncRNAs has been linked to diseases such as osteoporosis.<sup>(14)</sup> Additionally, there is accumulating evidence suggesting their involvement in key regulatory pathways, including osteogenic differentiation.<sup>(11,15)</sup>

Although understudied in the context of GWAS,<sup>(13)</sup> there is increasing evidence suggesting that lncRNAs are causal for a subset of associations identified by GWAS. A recent analysis of data from the Genotype-Tissue Expression (GTEx) project identified 690 potentially causal lncRNAs underlying associations influencing risk of a wide range of diseases.<sup>(13)</sup> Additionally, there is emerging evidence implicating lncRNAs in the genetics of BMD.<sup>(16-18)</sup> For example, a study reported 575 differentially expressed lncRNAs between high and low BMD groups in white women, 26 of which regulate protein-coding genes that are potentially causal in BMD GWAS.<sup>(19)</sup> Additionally, a recent BMD single nucleotide polymorphism (SNP) prioritization analysis implicated lncRNAs as potential effector transcripts.<sup>(20)</sup> Together these studies suggest that lncRNAs may play an important role in the genetic regulation of bone mass.

In recent years, a number of approaches have been developed that utilize transcriptomics data to inform GWAS, including the analysis of allelic imbalance (AI), transcriptome-wide association studies (TWASs), and expression quantitative trait loci (eQTL) colocalization.<sup>(21)</sup> AI results from the cis-regulatory effects (ie, local eQTL) that can be tracked using heterozygous coding SNPs. In TWASs the genetic component of gene expression in a reference population is estimated and then imputed in a much larger population. Once gene expression is imputed, genetically regulated gene expression is associated with a disease or disease phenotype.<sup>(22)</sup> Most genes identified by TWAS are located in GWAS associations for that disease and, as a result, TWAS can pinpoint genes likely to be causal at GWAS loci.<sup>(23,24)</sup> eQTLs are genetic variants associated with changes in gene expression and can be tissue-specific or shared across multiple tissues. eQTL colocalization tests whether the change in gene expression and the change in a trait of interest are driven by the same shared genetic variant(s). All three approaches, alone or in combination, have been successfully used to pinpoint potential causal disease genes at GWAS associations.

Here, we identified lncRNAs that are potentially responsible for the effects of BMD GWAS associations by first applying AI to bone samples and, next, applying TWAS and eQTL colocalization to gene expression data from GTEx. Through both approaches we identified 58 lncRNAs with evidence of being causal BMD GWAS genes. We further prioritized these lncRNAs by identifying those that were differentially expressed as a function of osteoblast differentiation. Together, these results highlight the potential importance of lncRNAs as candidate causal BMD GWAS genes.

## Methods

### Patient demographics

All human specimen collection was performed in accordance with institutional review board (IRB) approval from our institution (IRB number H-32517). Acetabular reaming from 17 Boston Medical Center (BMC) patients (ages 43–80 years) undergoing elective hip arthroplasty were collected: 12 females and 5 males; 8 black, 8 white, and 1 Hispanic. This demographic mix reflects the population serviced by Boston University Medical Campus (BUMC), which is an urban safety-net hospital.

### RNA extraction

Bone fragments were isolated from the 17 patients. Total RNA was isolated from bone fragments as described in Sagi and colleagues.<sup>(25)</sup> Total RNA sequencing (RNAseq) libraries were constructed from bone as well as human fetal osteoblast (hFOB) RNA samples using Illumina TruSeq Stranded Total RNA with Ribo-Zero Gold sample prep kits (Illumina, San Diego, CA, USA). Constructed libraries contained all RNAs greater than 100 nucleotides (nt) (both unpolyadenylated and polyadenylated) minus cytoplasmic and mitochondrial ribosomal RNAs (rRNAs). Samples were sequenced to achieve a minimum of 50 million reads  $2 \times 75$ -basepair (bp) paired-end reads on an Illumina NextSeq500 (Illumina).

### hFOB cell line culture

hFOB 1.19 cells (American Type Culture Collection [ATCC], Manassas, VA, USA; #CRL-11372) were cultured at 34C and differentiated at 39.5C as recommended with the following modifications. Growth media: minimal essential media (MEM; Gibco, Grand Island, NY, USA; 10370-021) supplemented with 10% fetal bovine serum (FBS; Atlantic Biologicals, Morrisville, NC, USA; S12450), 1% Glutamax (Gibco; 35050-061), 1% Pen Strep (Gibco; 15140-122). Differentiation media: MEM alpha (Gibco; 12571-063) supplemented with 10% FBS, 1% Glutamax, 1% Pen Strep, 50  $\mu\text{g}/\mu\text{L}$  Ascorbic Acid (Sigma-Aldrich, St. Louis, MO, USA; A4544-25G), 10mM beta-Glycerophosphate (Sigma-Aldrich; G9422-100G), 10nM Dexamethasone (Sigma-Aldrich; D4902-25MG). RNA was isolated from  $\sim 0.5 \times 10^6$  cells at days 0, 2, 4, 6, 8, and 10 of differentiation as recommended (RNAeasy Minikit; QIAGEN, Valencia, CA, USA; 74106). Mineralized nodule formation was measured by staining cultures with Alizarin Red (40mM, pH 5.6; Sigma-Aldrich; A5533-25G). Reported results were obtained from three biological replicate experiments.

### RNA sequencing and differential gene expression analysis

Computational analysis of RNA sequencing data for the 17 bone samples, Farr and colleagues<sup>(26)</sup> and the hFOB samples were performed using a custom bioinformatics pipeline. Briefly, FastQC (Babraham Bioinformatics, Cambridge, UK; <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) and RseqQC<sup>(27)</sup> were used to assess the quality of raw reads. Adapter trimming was completed using Trimmomatic.<sup>(28)</sup> Sequences were aligned to the GENCODE v34<sup>(29)</sup> reference genome using the SNP and splice aware aligner HISAT2.<sup>(30)</sup> Genome assembly and abundances in transcripts per million (TPM) were quantified using StringTie.<sup>(31)</sup> Differential expression analysis for the hFOB differentiation experiment was performed using Deseq2<sup>(32)</sup> across all

six differentiation time points using analysis of deviance (ANODEV) which is conceptually similar to analysis of variance (ANOVA). Differential expression analysis for the comparison between this study's samples and the Farr and colleagues<sup>(26)</sup> samples was performed using Deseq2<sup>(32)</sup> standard approach.

### lncRNA discovery

The Coding Potential Assessment Tool (CPAT)<sup>(33)</sup> was used to assess the protein-coding potential of the novel transcripts assembled. In short, CPAT is a machine learning algorithm trained on a set of known human lncRNAs to identify novel putative lncRNAs based on shared sequence features. We used all known lncRNAs in the latest human genome assembly (GRCh38) as the training set. Novel transcripts with coding probability <0.367 are regarded as lncRNAs in accordance with software authors. Novel lncRNAs with TPM <1 were regarded as noise and discarded.

### AI analysis

Reads were aligned to the GENCODE v34<sup>(29)</sup> reference genome using the SNP and splice aware aligner HISAT2.<sup>(30)</sup> The resultant BAM files were then used as input for variant calling using the GATK pipeline.<sup>(34)</sup> Briefly, duplicate reads were identified using MarkDuplicates. Next, reads spanning introns were reformatted using SplitNCigarReads to match the DNA aligner conventions. Then base quality recalibration was performed to detect and correct for patterns of systematic errors in the base quality scores. Finally, the variant calling and filtration step was performed using HaplotypeCaller. The resultant VCF file included only known and novel SNPs and reference bias was corrected using WASP.<sup>(35)</sup> Briefly, mapped reads that overlap SNPs are identified. For each read that overlaps a SNP, its genotype is swapped with that of the other allele and it is re-mapped. If a re-mapped read fails to map to exactly the same location, it is discarded. The resultant corrected BAM and filtered VCF files were used as input for GATK ASEReadCounter to provide a table of filtered base counts at heterozygous sites for allele specific expression. Bases with a read depth less than 20 were discarded. In order to determine significance, a binomial test was performed and only heterozygous sites with false discovery rate (FDR)-corrected *p* value of <0.05 were considered significant.

### TWAS

We conducted a TWAS by integrating genomewide SNP-level association summary statistics from a BMD GWAS<sup>(9)</sup> with GTEx version 8 gene expression QTL data from 49 tissue types. We used the S-MultiXcan<sup>(36)</sup> approach for this analysis, to correlate gene expression across tissues to increase power and identify candidate susceptibility genes. Gene-level associations were identified at FDR correction <0.05 and were further filtered using fastENLOC (a faster implementation of ENLOC<sup>(37)</sup>) regional colocalization probability >0.1 in at least one tissue type.

### Bayesian colocalization analysis

We used fastENLOC to perform Bayesian colocalization analysis. We integrated summary statistics from the most recent (and largest) eBMD GWAS<sup>(9)</sup> and eQTL data from 49 GTEx tissues.<sup>(38)</sup> We used the recommended regional colocalization probability (RCP) threshold of >0.1 as indication of significant overlap between SNP and eQTL.

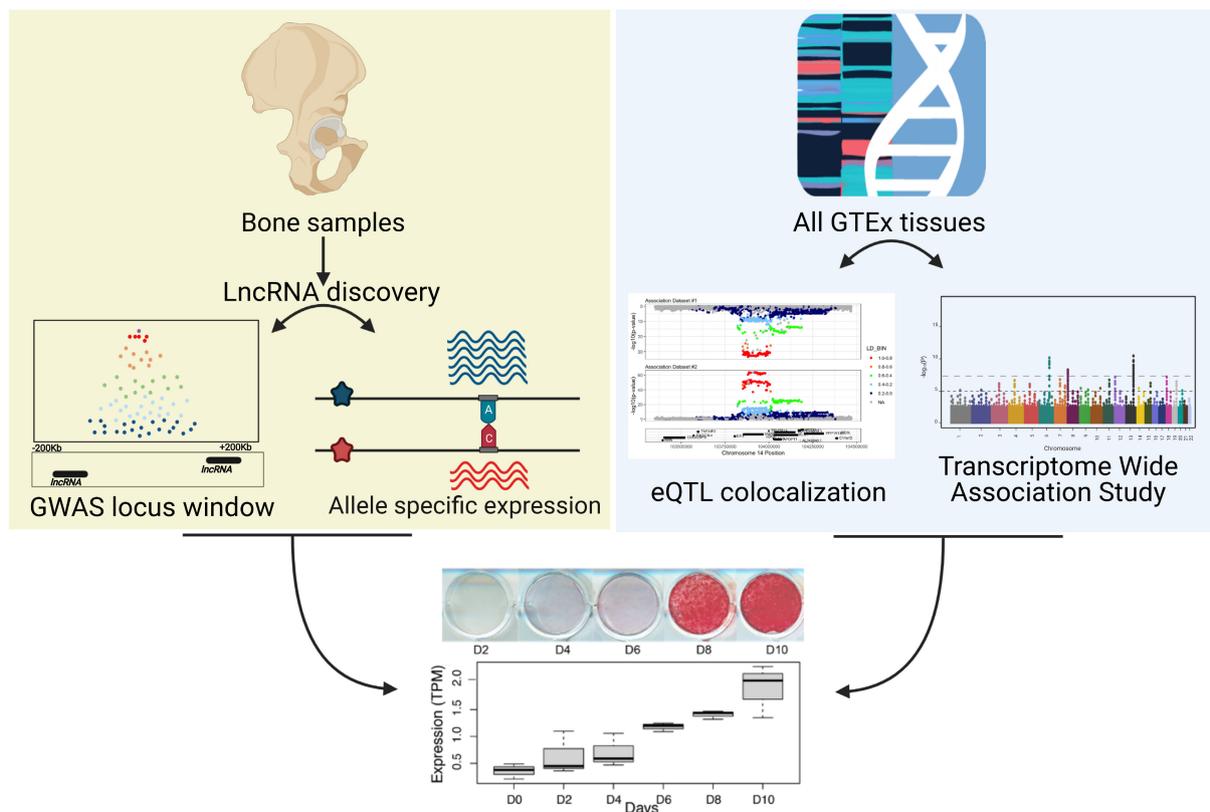
## Results

In this study, we used two approaches to identify lncRNAs that potentially underlie BMD GWAS associations. In the first approach, we quantified known (lncRNAs that have been reported in the GENCODE database) and novel lncRNAs using RNA sequencing (RNAseq) data from human bone fragments and identified lncRNAs located in proximity of a BMD GWAS association and harboring SNPs demonstrating AI. In the second approach, we leveraged GTEx to identify lncRNAs across a large number of tissues and cell-types whose expression was significantly associated with BMD by TWAS and regulated by an eQTL that colocalized with a BMD association. Figure 1 provides an overview of our study.

### Generation of bone expression data from bone fragments

To identify potentially casual lncRNAs in a BMD relevant tissue, we generated total RNAseq (ribo-depleted) data on bone fragments isolated from acetabular reamings from patients undergoing hip arthroplasty (*n* = 17; 5 males and 12 females; ages 43 to 80 years). The acetabular reamings were comprised primarily of bone and marrow with a small number of contaminating cartilage fragments. In contrast to most gene expression data generated on bone which are typically from biopsies that contain marrow, we were able to remove the marrow leaving purified trabecular and cortical bone. We hypothesized that the acetabular bone fragments consisted primarily of late-stage osteoblasts/osteocytes,<sup>(39)</sup> allowing us to characterize lncRNAs enriched in these cell types. To confirm that the acetabular samples were enriched in osteocytes, we compared these data to published RNAseq data on bone biopsies.<sup>(26)</sup> Farr and colleagues<sup>(26)</sup> generated RNAseq data on 58 iliac crest needle biopsies from healthy women containing both bone and marrow. Average transcripts per million (TPM) across all samples in both experiments were highly correlated (Fig. 2A,  $r^2 = 0.845$ ,  $p < 2.2 \times 10^{-16}$ ). Importantly, differential expression analysis between the two datasets showed that the top 1000 genes with the largest fold change increase in the bone fragment samples compared to bone biopsy samples were enriched in Gene Ontology (GO) terms such as "skeletal system development" (FDR =  $4.01 \times 10^{-3}$ ) and "extracellular matrix organization" (FDR =  $4.11 \times 10^{-5}$ ).

To support the notion that our samples are unique in osteocyte enrichment, we used data from a recent study that identified an osteocyte gene signature consisting of 1239 genes in mice and their orthologs in humans.<sup>(40)</sup> The ratio of expression (bone fragment samples/bone biopsy samples) was used. A ratio value >1 indicates that gene expression is higher in the bone fragment samples relative to the bone biopsy samples. In contrast, a ratio value <1 indicates that the gene is highly expressed in bone biopsy samples relative to bone fragment samples. We expect to see that osteocyte signature genes show ratio values >1 and marrow enriched genes show ratio values <1. The osteocyte signature genes showed a median ratio of 1.72 (62% of osteocyte signature genes ratio >1). Additionally, the ratio of expression of genes enriched in bone marrow showed a median of 0.27 (91% of marrow enriched genes ratio <1). The distribution of osteocyte signature genes ratio values showed a significant median shift (Wilcoxon test,  $p < 2.2 \times 10^{-16}$ ) (Fig. 2D), and the opposite pattern was observed for the bone marrow enriched genes (Wilcoxon test,  $p < 2.2 \times 10^{-16}$ ) (Fig. 2E). In addition, we compared the expression of osteocyte-specific genes reported in Bonewald<sup>(39)</sup> (Fig. 2B) and bone marrow enriched



**Fig. 1.** Overview of the study. We conducted de novo lncRNA discovery using RNAseq data on human acetabular bone fragments from 17 patients. We then identified known and novel lncRNAs located in GWAS associations that were influenced by AI (yellow box). We applied TWAS and colocalization on eQTL data from 49 GTEx project tissues (blue box). We assessed the role of lncRNAs reported by both approaches in osteogenic differentiation using RNA-seq data from the hFOB cell line at six time points across differentiation (bottom panel). AI = allelic imbalance; GTEx = genotype-tissue expression; hFOB = human fetal osteoblast; TWAS = transcriptome-wide association study.

genes reported in ([www.proteinatlas.org](http://www.proteinatlas.org)) (Fig. 2C). In addition, during the isolation, care was also taken to remove the cartilage fragments. We repeated the analysis for cartilage marker genes and found a modest reduction ( $p = 0.035$ ) of expression in our samples.<sup>(41)</sup> The difference was more modest likely due to a significant overlap in the expression of these genes in both cartilage and bone/osteoblasts. Altogether, these data suggest that the purified acetabular bone fragments are enriched for late osteoblasts/osteocytes and are more marrow depleted compared to iliac crest biopsies.

### Identifying novel lncRNAs in purified acetabular bone fragments

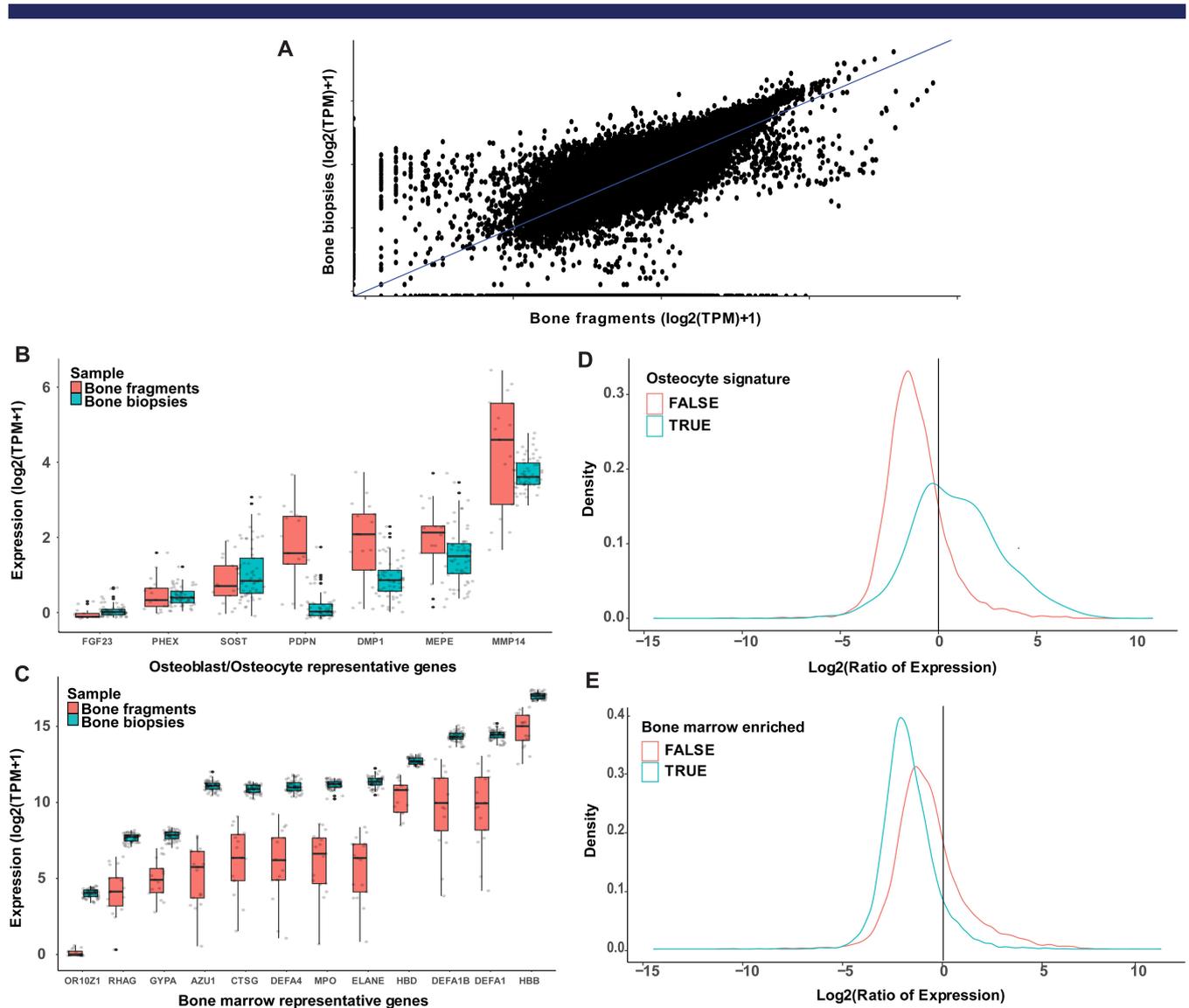
Given the paucity of bone transcriptomics data in the literature, and the tissue-specific nature of lncRNA expression, we hypothesized that many bone/osteocyte-specific lncRNAs would not be present in current sequence databases. Additionally, ~50% of lncRNAs do not possess a poly-A tail modification<sup>(42)</sup> and most RNAseq data is generated after poly-A selection. Therefore, in order to capture a more comprehensive profile of lncRNAs in bone, we implemented a lncRNA discovery step to identify putative “novel” lncRNA transcripts using the computational algorithm CPAT.<sup>(33)</sup> Across the 17 bone samples we identified 6612 known lncRNAs and 2440 novel lncRNAs (Tables S1 and S2).

The mean length of novel lncRNAs was 30.3 kilobases (kb) and median length of 11.8 kb. These values were comparable to the mean length of known lncRNAs expressed in the bone samples (mean = 35.4 kb; median = 4.7 kb).

### Identifying potentially causal lncRNAs in bone

For lncRNAs to be considered potentially causal in bone, we identified those that are both located in proximity of a BMD GWAS association and regulated by AI. We hypothesized that such genes may be causal for their respective associations because of the potential to be regulated by an eQTL which colocalizes with a BMD association. Of the 9,052 lncRNAs (2440 novel and 6612 known) we quantified in acetabular bone, 1496 lncRNAs (~17% of expressed lncRNAs) were found within a 400-kb window ( $\pm 200$  kb from the lncRNA start site) of each of 1103 GWAS associations previously identified by Morris and colleagues.<sup>(9)</sup> The rationale behind choosing this genomic distance was based on findings in Vösa and colleagues,<sup>(43)</sup> where they showed that 92% of lead cis-eQTLs are within 100 kb of the transcription start site. Therefore, this window was extended to ensure we captured the majority of all cis-eQTL effects.

Next, we identified heterozygous coding variants that demonstrated significant evidence of AI within lncRNAs. None of the heterozygous coding SNPs used to assess AI were in linkage



**Fig. 2.** Enrichment of osteocyte marker genes in bone fragment samples (used in this study) compared to bone biopsy samples in the literature. (A) Overall gene expression is highly correlated between the RNAseq data generated in both studies ( $r^2 = 0.845$ ,  $p < 2.2 \times 10^{-16}$ ); Farr and colleagues.<sup>(26)</sup> (B) Gene expression of osteocyte marker genes reported in Bonewald<sup>(39)</sup> showing enrichment in the bone fragments samples (this study) relevant to bone biopsies. (C) Gene expression of bone marrow enriched genes reported in The Human Protein Atlas ([www.proteinatlas.org/](http://www.proteinatlas.org/)) showing higher expression in bone biopsy samples. (D) Osteocyte signature genes reported in Youlten and colleagues<sup>(40)</sup> are highly expressed in bone fragment samples relative to bone biopsies (Wilcoxon test,  $p < 2.2 \times 10^{-16}$ ) (E) Bone marrow enriched genes reported in Youlten and colleagues<sup>(40)</sup> are highly expressed in bone biopsy samples compared to bone fragment samples (Wilcoxon test,  $p < 2.2 \times 10^{-16}$ ).

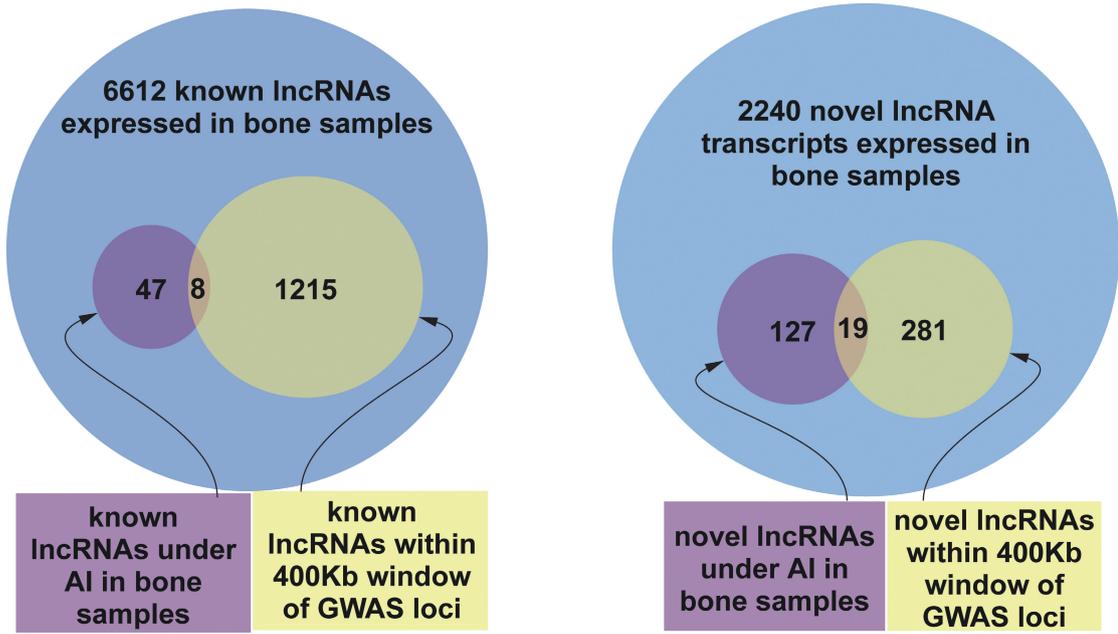
disequilibrium (LD) ( $r^2 < 0.05$ ) with a lead BMD GWAS SNP, which is expected because these SNPs were only used to measure AI and not necessarily functionally associated with lead GWAS SNPs. Of the total number of lncRNAs we identified, 174 (47 known, 127 novel; ~2% of expressed lncRNAs) had at least one SNP demonstrating AI in at least one of the 17 bone fragment samples. Out of the 174, 27 (15.5%; 8 known, 19 novel) were located in proximity of a GWAS association (Fig. 3A, Table S3). It is expected that we find a low number of lncRNAs (known or novel) under AI relative to the number of expressed lncRNAs within 400 kb of GWAS loci. Reasons for our expectation include the absence of an exonic heterozygous SNP because some

lncRNAs that do not have an exonic heterozygous SNPs in LD with a regulatory SNP within the 17 acetabular bone samples will be missing from the intersection. Additionally, lncRNAs in general are lowly expressed; therefore, the power to identify lncRNAs under AI is lower than that of protein-coding genes.

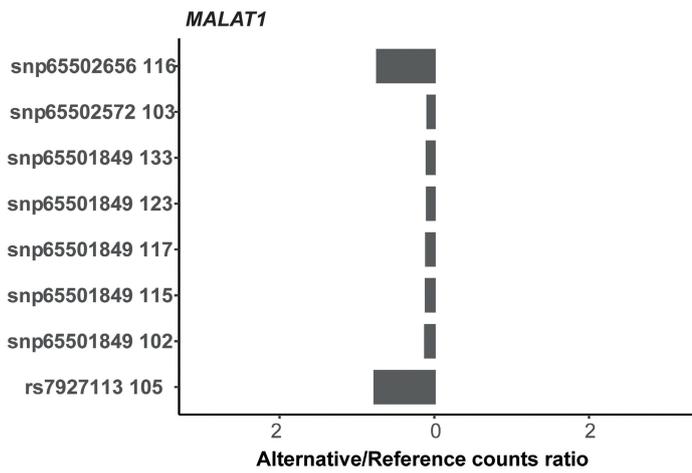
#### Identifying putatively causal lncRNAs by leveraging GTEx

Next, we sought to leverage non-bone data to identify potentially causal lncRNAs. To do this, we integrated 1103 BMD GWAS loci<sup>(9)</sup> with GTEx (v8) eQTL data by coupling TWAS<sup>(44)</sup> using S-MultiXScan<sup>(36)</sup> and Bayesian colocalization analysis using

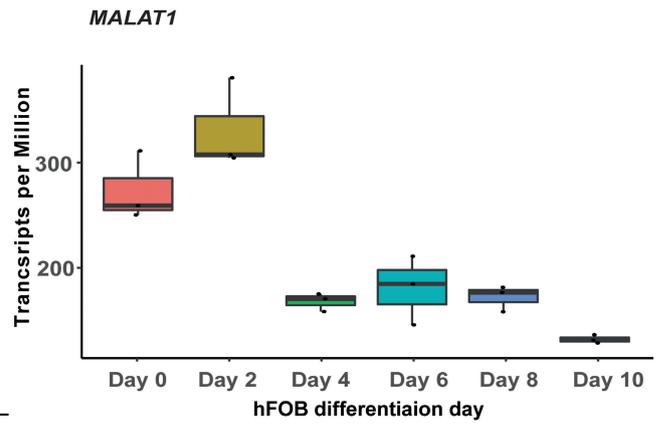
A



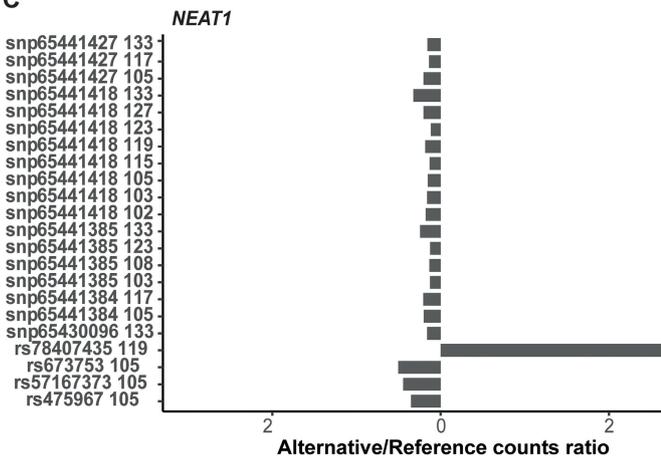
B



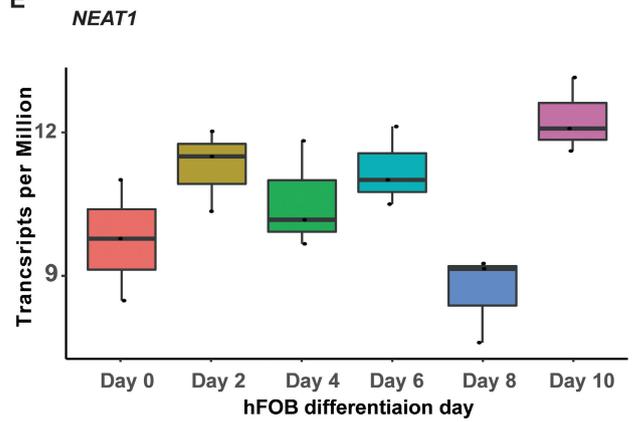
D



C



E



(Figure legend continues on next page.)

fastENLOC.<sup>(37)</sup> The rationale behind using GTEx data is genes that are shared in multiple tissues and showing a colocalizing eQTL with BMD GWAS data can be potentially causal in bone tissue as well. Our TWAS analysis resulted in 333 significant lncRNA-BMD associations (FDR correction <0.05), which constitute 5% of all known lncRNAs that are expressed in the acetabular samples. Our colocalization analysis yielded 48 lncRNAs with a colocalizing eQTL (regional colocalization probability [RCP] >0.1) in at least one GTEx tissue. These lncRNAs with a colocalizing eQTL make up <1% of the known expressed lncRNAs in the acetabular bone samples. There were 31 lncRNAs (<1%) significant in both the TWAS and eQTL colocalization analysis (Table S4).

Most identified lncRNAs are the only potential effector transcripts implicated by TWAS/eQTL colocalization in their respective GWAS associations

To determine if the lncRNAs listed in Table S4 are the strongest candidates in their respective GWAS associations, we evaluated a recent report of protein coding genes that used the same approach.<sup>(45)</sup> Five out of the 31 lncRNAs (*LINC01116*, *LINC01117*, *SNHG15*, *LINC01290*, *LINC00665*) have a protein coding gene with a colocalizing eQTL (*HOXD8*, *HOXD9*, *MYO1G*, *NACAD*, *EMP2*, *ZFP14*, *ZFP82*) within 1 megabase (Mb) of the lncRNA start site (Table S5). Upon further investigation of the RCP values, some of the lncRNAs showed higher RCP than their protein coding gene counterpart. For example, *LINC01290* had a higher RCP in lung tissue (0.4992) compared to its counterpart *EMP2* (0.2227). On the other hand, the same lncRNA has a lower RCP value (0.1498) than *EMP2* (0.6089) in breast and mammary gland tissue. However, for the remaining lncRNAs, this analysis provides support that the lncRNA alone is the potential effector transcript in the region because we show no evidence of protein coding colocalization within 1 Mb distance of the start site of the lncRNA.

Many identified lncRNAs are differentially expressed as a function of osteoblast differentiation

To provide further support for the hypothesis that these lncRNAs mediate GWAS associations, we measured their expression as a function of osteoblast differentiation in human fetal osteoblasts (hFOBs). We performed total RNAseq at six hFOB differentiation time-points (days 0, 2, 4, 6, 8, and 10). Of the 27 lncRNAs implicated in the analysis of AI, all eight known lncRNAs were differentially expressed (FDR <0.05). On the other hand, none of the novel lncRNAs were differentially expressed (Table S3). Examples of the identified genes include *MALAT1* and *NEAT1* (Fig. 3B,C), which were differentially expressed in hFOBs and showed evidence of AI in 8 and 10 of the 17 acetabular bone samples, respectively. There were four unique SNPs in the exonic regions of *MALAT1* (Fig. 3B) that were heterozygous in at least one of the

17 individuals (with a maximum of eight individuals). All four SNPs showed higher expression in the alternative allele relative to the reference allele. The expression of *MALAT1* gene decreased as the cell differentiated into a mineralizing state (Fig. 3D). Additionally, there were nine unique SNPs reported in the exonic regions of *NEAT1* that were heterozygous in at least one of the 17 individuals (with a maximum of 10 individuals). Of the nine, eight showed higher expression associated with the alternative allele compared to the reference allele. The remaining SNP was associated with the opposite pattern and this was likely due to it being the only SNP not in high LD with the others ( $r^2 = 0.0021$ ). *NEAT1* showed significant increase in expression around day 10 in hFOBs (Fig. 3E).

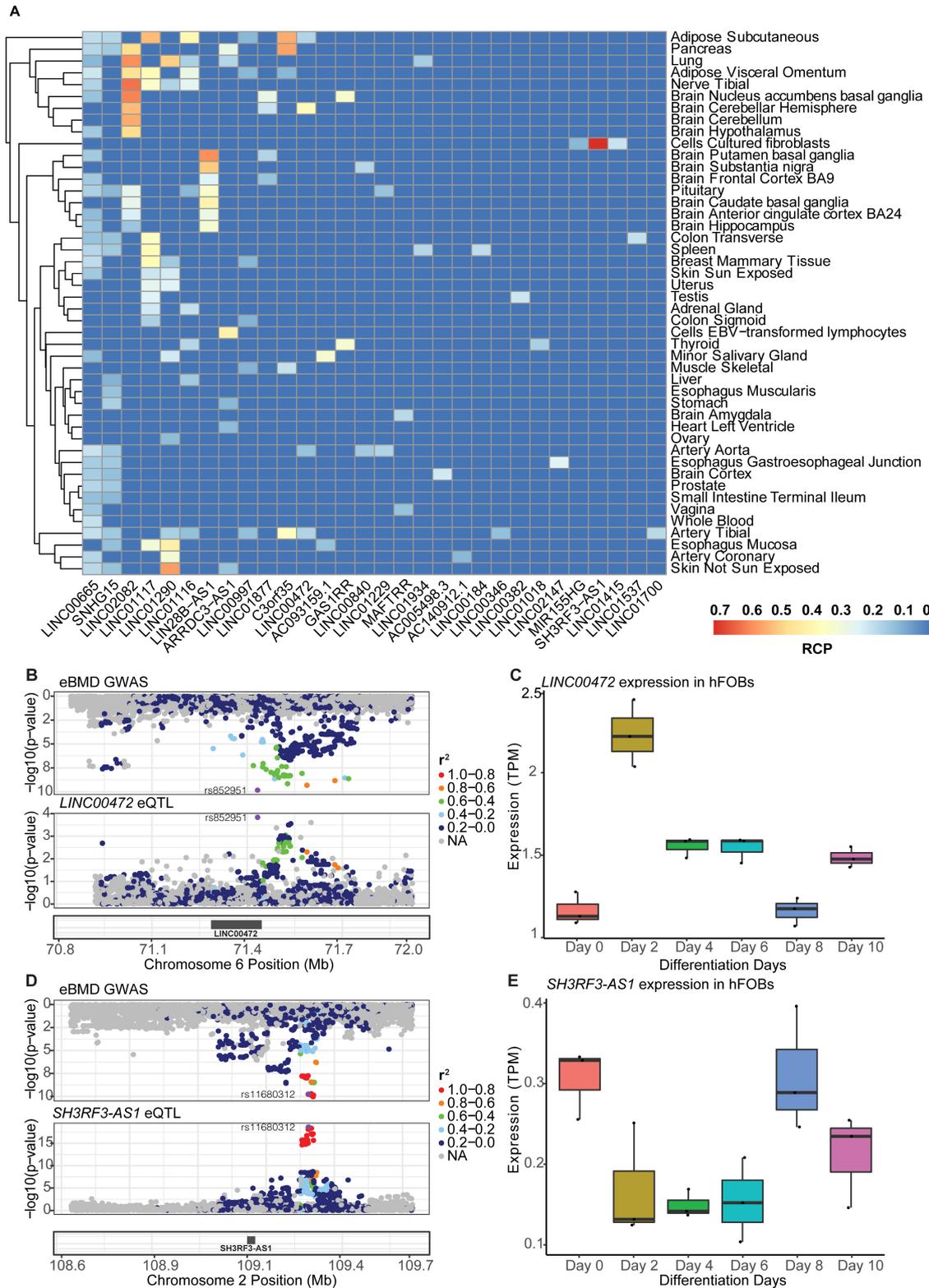
We assessed the expression of lncRNAs identified by GTEx TWAS/eQTL colocalization in osteoblast differentiation using the same approach in the previous section. Out of the 31 lncRNAs identified by TWAS/eQTL colocalization, 15 were found to be differentially expressed (*LINC00184*, *SH3RF3-AS1*, *LINC01116*, *LINC01934*, *C3orf35*, *LINC01018*, *ARRDC3-AS1*, *LINC00472*, *SNHG15*, *GAS1RR*, *LINC00840*, *LINC01537*, *LINC00346*, *LINC01415*, *MIR155HG*). In general, the expression of those genes in hFOBs was low compared to the lncRNAs reported in the AI section. Examples include *SHR3F3-AS1* and *LINC00472*, which were regulated by colocalizing eQTL (Fig. 4B,D) and were differentially expressed in hFOBs (Fig. 4C,E). *SHR3F3-AS1* was shown to have the highest RCP value overall (RCP = 0.72) and in only one GTEx tissue (cultured fibroblasts) (Fig. 4A,D, Table S4). Although the gene was differentially expressed across hFOB differentiation points, it had a very low overall level of expression (Fig. 4E). The pattern of expression decreased during mid differentiation points with spikes in early and late points (Fig. 4E). *LINC00472* was shown to have a colocalizing eQTL in four GTEx tissues with the highest RCP value in brain cerebellar hemisphere (RCP = 0.37) (Fig. 4A,B, Table S4). The gene also showed a moderate level of expression in hFOBs with an average of 1.5 TPM (Fig. 4C). The expression of *LINC00472* peaked at day 2 and then declined (Fig. 4C).

## Discussion

In this study, we interrogated BMD GWAS loci and identified known and novel lncRNAs as potential effector transcripts. We identified potentially important lncRNA using two different approaches. First, we identified novel and known lncRNAs in a unique transcriptomic bone dataset that were localized in GWAS loci and demonstrated AI. Second, we implicated additional lncRNAs by leveraging GTEx and identifying eQTLs in non-bone tissues that colocalized with eBMD GWAS loci whose expression was associated with eBMD via TWAS. We also assessed differential expression across the time course of hFOB differentiation to

(Figure legend continued from previous page.)

**Fig. 3.** Identification of lncRNAs located within eBMD GWAS associations, are under AI in acetabular bone, and are differentially expressed in hFOBs. (A) Venn diagram showing the number of known and novel lncRNAs within proximity of GWAS loci, implicated by AI, and implicated by both approaches. (B) lncRNA *MALAT1* AI plot showing the ratio of reads aligning to the alternative SNP relative to the reference SNP in eight of the bone fragments samples where the gene is under AI. (C) lncRNA *NEAT1* AI plot showing the ratio of reads aligning to the alternative SNP relative to the reference SNP in 10 of the bone fragments samples where the gene is under AI. rs78407435 is not in LD with the rest of the SNPs in the region and this is likely the reason it shows a different direction of effect. (D) Expression of *MALAT1* across hFOB differentiation points. (E) Expression of *NEAT1* across hFOB differentiation points. AI = allelic imbalance; hFOB = human fetal osteoblast.



**Fig. 4.** lncRNAs implicated by eQTL colocalization and TWAS are potential effector transcripts of BMD GWAS loci. (A) Heat map showing colocalization events in GTEx tissues. (B) lncRNA *LINC00472* colocalization plot showing colocalization of eBMD GWAS locus with eQTL from brain cerebellar hemisphere with RCP of 0.37 (C) Differential expression of *LINC00472* across hFOB differentiation points (D) lncRNA *SH3RF3-AS1* colocalization plot showing colocalization of eBMD GWAS locus with GTEx fibroblasts eQTL data with RCP of 0.72 (E) Differential expression of *SH3RF3-AS1* across hFOB differentiation points. hFOB = human fetal osteoblast.

provide more evidence of a potential causal role for these lncRNAs.

In the first approach, we set out to perform transcriptomics on a unique sets of bone samples in order to identify novel lncRNAs in bone, provide deeper coverage for known lncRNA identification, and apply AI analysis. The bone samples that exist in the literature are from bone biopsies, and as we show in the Results section, they are less enriched in bone-relevant genes compared to the dataset produced by the bone fragments used in this study.

A total of eight lncRNAs (*NEAT1*, *MALAT1*, *DLEU2*, *LINC01578*, *CARMN*, *AC011603.3*, *PXN-AS1*, *AC020656.1*) were found to be within a 400-kb window of an eBMD GWAS locus and were also differentially expressed across hFOB differentiation time points. Many of these lncRNAs have been demonstrated to play a role in bone. For example, *NEAT1* has been reported to stimulate osteoclastogenesis via sponging *miR-7*<sup>(46)</sup> and the *NEAT1/miR-29b-3p/BMP1* axis promotes osteogenic differentiation in human bone marrow-derived mesenchymal stem cells.<sup>(47)</sup> In addition, *MALAT1* has been shown to influence BMD.<sup>(48)</sup> *MALAT1* acts as a sponge of miR-34c to promote the expression of *SATB2*. *SATB2* then acts to reduce the alkaline phosphatase (ALP) activity of osteoblasts and mineralized nodules formation.<sup>(48)</sup> A recent study<sup>(49)</sup> has shown that *LINC01578* (referred to as *CHASERR* in this study) represses chromodomain Helicase DNA Binding Protein 2 (*Chd2*). A model for *Chd2* loss of function by the International Mouse Phenotyping Consortium (IMPC)<sup>(50)</sup> reported that these mice exhibit significant decreased body weight and length, skeletal abnormalities, abnormal bone structure, decreased fat levels, and BMD.<sup>(49)</sup> Last, *DLEU2* expression has been shown to be inversely correlated with BMD in a study involving postmenopausal white women.<sup>(51)</sup> The remaining four lncRNAs have not been reported to date to have a role in bone and should be further pursued.

In our second analysis, we reported 15 lncRNAs implicated jointly by colocalization, TWAS, and differential expression analysis. We show one example of the 15 lncRNAs reported in *SH3RF3-AS1* in Fig. 4A. Most of these lncRNAs have not been shown previously in the literature to have a role in bone biology. However, *LINC00472* (Fig. 4B) has been experimentally shown to influence osteogenic differentiation by sponging miR-300, which in turn increases the expression of *Fgfr2* in mice.<sup>(52)</sup> These preliminary results provide more evidence of the potential causal role of these lncRNAs in osteoporosis.

In this study, we were able to use multiple systems genetics approaches on two transcriptomic datasets (acetabular bone and GTEx) to identify lncRNAs that are potentially responsible for the effects of some BMD GWAS loci. This is the first study to our knowledge that evaluated the role of lncRNAs in mediating the effect of BMD GWAS loci from a genomewide perspective. We combined osteoblast differentiation samples and the literature to provide experimental evidence in previous studies to support the effector transcript list we generated from our analysis. These results highlight the importance of studying other aspects of the transcriptome to identify potential drug targets for osteoporosis and bone fragility.

### Limitations of this study

This study is not meant to be comprehensive because we are limited by the number of samples and are not suitably powered to identify eQTLs and apply TWAS/colocalization analysis. However, due to the scarcity of population-level bone transcriptomic

datasets, and the lack of bone cell or tissue data in GTEx, our study is an attempt to systematically leverage the available datasets to capture a subset of lncRNAs that we think are potentially causal. As mentioned, some of these lncRNAs have been implicated experimentally outside of this study. Moreover, lncRNAs under AI and within proximity of GWAS loci may not be causal as they could be false positives because they are not prioritized via a systems analysis such as colocalization. Another limitation of our study is that we evaluated their expression as a function of osteoblast differentiation; however, it is likely that some of the lncRNAs, if truly causal, impact BMD via a function in other cell-types (eg, osteoclasts). Future studies should focus on enhancing these results by generating transcriptomic and eQTL datasets from bone and other bone cell types, using network approaches to aid in the prioritization of lncRNAs, and experimentally validating the role of specific lncRNAs.

## Acknowledgments

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 CRF, LCG, and EFM, and AA was supported in part by a National Institutes of Health, Biomedical Data Sciences Training Grant (5T32LM012416). The authors acknowledge Emily Farber for generating RNAseq data on bone fragments. We thank the IMPC for accessibility to BMD data in knockout mice ([www.mousephenotype.org](http://www.mousephenotype.org)). 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.

Authors' roles: AA: Conceptualization, Methodology, Formal analysis, Investigation, Validation, Visualization, Writing—original draft, Writing—review & editing. LM: Data curation, Formal analysis, Investigation, Methodology, Visualization, Writing—review & editing. WR: Data curation, Formal analysis, Investigation, Methodology. BMA-B: Investigation, formal analysis. NH: Investigation. EFM: Conceptualization, Funding acquisition, Project administration, Resources, Supervision, Writing—review & editing. LCG: Conceptualization, Funding acquisition, Project administration, Resources, Supervision, Writing—review & editing. CRF: Conceptualization, Funding acquisition, Project administration, Resources, Supervision, Writing—review & editing.

## Author Contributions

**Abdullah Abood:** Conceptualization; formal analysis; investigation; methodology; validation; writing – original draft; writing – review and editing. **Larry Mesner:** Data curation; formal analysis; investigation; methodology; visualization; writing – review and editing. **Will Rosenow:** Data curation; formal analysis; investigation; methodology. **Basel M. Al-Barghouthi:** Formal analysis; investigation. **Nina Horowitz:** Investigation. **Elise F. Morgan:** Funding acquisition; writing – review and editing. **Louis C. Gerstenfeld:** Funding acquisition; resources; supervision; writing – review and editing. **Charles R. Farber:** Conceptualization; funding acquisition; project administration; resources; supervision; writing – review and editing.

## Conflict of Interest

The authors declare that they have no conflicts of interest with the contents of this article.

## Data Availability Statement

The data that support the findings of this study are openly available in Gene Expression Omnibus (GEO) at <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE186922>, reference number GSE186922.

## References

1. NIH Consensus Development Panel on Osteoporosis Prevention, Diagnosis, and Therapy. Osteoporosis prevention, diagnosis, and therapy. *JAMA*. 2001;285:785-795.
2. Office of the Surgeon General (US). Bone Health and Osteoporosis: A Report of the Surgeon General. Rockville, MD: Office of the Surgeon General (US); 2010.
3. Burge R, Dawson-Hughes B, Solomon DH, Wong JB, King A, Tosteson A. Incidence and economic burden of osteoporosis-related fractures in the United States, 2005-2025. *J Bone Miner Res*. 2007;22:465-475.
4. Johnell O, Kanis JA, Oden A, et al. Predictive value of BMD for hip and other fractures. *J Bone Miner Res*. 2005;20:1185-1194.
5. Smith DM, Nance WE, Kang KW, Christian JC, Johnston CC Jr. Genetic factors in determining bone mass. *J Clin Invest*. 1973;52:2800-2808.
6. Arden NK, Baker J, Hogg C, Baan K, Spector TD. The heritability of bone mineral density, ultrasound of the calcaneus and hip axis length: a study of postmenopausal twins. *J Bone Miner Res*. 1996;11:530-534.
7. Slemenda CW, Turner CH, Peacock M, et al. The genetics of proximal femur geometry, distribution of bone mass and bone mineral density. *Osteoporos Int*. 1996;6:178-182.
8. Richards JB, Zheng HF, Spector TD. Genetics of osteoporosis from genome-wide association studies: advances and challenges. *Nat Rev Genet*. 2012;13:576-588.
9. Morris JA, Kemp JP, Youtlen SE, et al. An atlas of genetic influences on osteoporosis in humans and mice. *Nat Genet*. 2019;51:258-266.
10. Claussnitzer M, Cho JH, Collins R, et al. A brief history of human disease genetics. *Nature*. 2020;577:179-189.
11. Zhang J, Hao X, Yin M, Xu T, Guo F. Long non-coding RNA in osteogenesis: a new world to be explored. *Bone Joint Res*. 2019;8:73-80.
12. Marchese FP, Raimondi I, Huarte M. The multidimensional mechanisms of long noncoding RNA function. *Genome Biol*. 2017;18:206.
13. de Goede OM, Nachun DC, Ferraro NM, et al. Population-scale tissue transcriptomics maps long non-coding RNAs to complex disease. *Cell*. 2021;84(10):2633-2648.e19. <https://doi.org/10.1016/j.cell.2021.03.050>.
14. Silva AM, Moura SR, Teixeira JH, Barbosa MA, Santos SG, Almeida MI. Long noncoding RNAs: a missing link in osteoporosis. *Bone Res*. 2019;7:10.
15. Nardocci G, Carrasco ME, Acevedo E, Hodar C, Meneses C, Montecino M. Identification of a novel long noncoding RNA that promotes osteoblast differentiation. *J Cell Biochem*. 2018;119:7657-7666.
16. Chen XF, Zhu DL, Yang M, 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*. 2018;102:776-793.
17. Roca-Ayats N, Martínez-Gil N, Cozar M, et al. Functional characterization of the C7ORF76 genomic region, a prominent GWAS signal for osteoporosis in 7q21.3. *Bone*. 2019;123:39-47.
18. Mei B, Wang Y, Ye W, et al. LncRNA ZBTB40-IT1 modulated by osteoporosis GWAS risk SNPs suppresses osteogenesis. *Hum Genet*. 2019;138:151-166.
19. Zhou Y, Xu C, Zhu W, et al. Long noncoding RNA analyses for osteoporosis risk in caucasian women. *Calcif Tissue Int*. 2019;105:183-192.
20. Zhang X, Deng HW, Shen H, Ehrlich M. Prioritization of osteoporosis-associated genome-wide association study (GWAS) single-nucleotide polymorphisms (SNPs) using epigenomics and transcriptomics. *JBMR Plus*. 2021;5:e10481.
21. Hukku A, Pividori M, Luca F, Pique-Regi R, Im HK, Wen X. Probabilistic colocalization of genetic variants from complex and molecular traits: promise and limitations. *Am J Hum Genet*. 2020;108:25-35. <https://doi.org/10.1016/j.ajhg.2020.11.012>.
22. Abood A, Farber CR. Using “-omics” data to inform genome-wide association studies (GWASs) in the osteoporosis field. *Curr Osteoporos Rep*. 2021;19:369-380. <https://doi.org/10.1007/s11914-021-00684-w>.
23. Li D, Liu Q, Schnable PS. TWAS results are complementary to and less affected by linkage disequilibrium than GWAS. *Plant Physiol*. 2021;186:1800-1811.
24. Wainberg M, Sinnott-Armstrong N, Mancuso N, et al. Opportunities and challenges for transcriptome-wide association studies. *Nat Genet*. 2019;51:592-599.
25. Sagi HC, Young ML, Gerstenfeld L, Einhorn TA, Tornetta P. Qualitative and quantitative differences between bone graft obtained from the medullary canal (with a reamer/irrigator/aspirator) and the iliac crest of the same patient. *J Bone Joint Surg Am*. 2012;94:2128-2135.
26. Farr JN, Roforth MM, Fujita K, et al. Effects of age and estrogen on skeletal gene expression in humans as assessed by RNA sequencing. *PLoS One*. 2015;10:e0138347.
27. Wang L, Wang S, Li W. RSeQC: quality control of RNA-seq experiments. *Bioinformatics*. 2012;28:2184-2185.
28. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics*. 2014;30:2114-2120.
29. Frankish A, Diekhans M, Ferreira AM, et al. GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res*. 2019;47:D766-D773.
30. Kim D, Langmead B, Salzberg SL. HISAT: a fast spliced aligner with low memory requirements. *Nat Methods*. 2015;12:357-360.
31. Pertea M, Pertea GM, Antonescu CM, Chang TC, Mendell JT, Salzberg SL. StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol*. 2015;33:290-295.
32. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol*. 2014;15:550.
33. Wang L, Park HJ, Dasari S, Wang S, Kocher JP, Li W. CPAT: coding-potential assessment tool using an alignment-free logistic regression model. *Nucleic Acids Res*. 2013;41:e74.
34. Castel SE, Levy-Moonshine A, Mohammadi P, Banks E, Lappalainen T. Tools and best practices for data processing in allelic expression analysis. *Genome Biol*. 2015;16:195.
35. van de Geijn B, McVicker G, Gilad Y, Pritchard JK. WASP: allele-specific software for robust molecular quantitative trait locus discovery. *Nat Methods*. 2015;12:1061-1063.
36. Barbeira AN, Pividori M, Zheng J, Wheeler HE, Nicolae DL, Im HK. Integrating predicted transcriptome from multiple tissues improves association detection. *PLoS Genet*. 2019;15:e1007889. <https://doi.org/10.1371/journal.pgen.1007889>.
37. 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*. 2017;13:e1006646.
38. GTEx Consortium. Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science*. 2015;348:648-660.
39. Bonewald LF. The amazing osteocyte. *J Bone Miner Res*. 2011;26:229-238.
40. Youtlen SE, Kemp JP, Logan JG, et al. Osteocyte transcriptome mapping identifies a molecular landscape controlling skeletal homeostasis and susceptibility to skeletal disease. *Nat Commun*. 2021;12:2444.
41. Funari VA et al. Cartilage-selective genes identified in genome-scale analysis of non-cartilage and cartilage gene expression. *BMC Genomics*. 2007;8:165.

42. Cheng J, Kapranov P, Drenkow J. Transcriptional maps of 10 human chromosomes at 5-nucleotide resolution. *Science*. 2005;308:1149-1154.
43. Vösa U, Claringbould A, Westra HJ, et al. Large-scale cis- and trans-eQTL analyses identify thousands of genetic loci and polygenic scores that regulate blood gene expression. *Nat Genet*. 2021;53:1300-1310.
44. Gusev A, Ko A, Shi H, et al. Integrative approaches for large-scale transcriptome-wide association studies. *Nat Genet*. 2016;48:245-252.
45. Al-Barghouthi BM, Rosenow WT, Du KP, et al. Transcriptome-wide association study and eQTL colocalization identify potentially causal genes responsible for bone mineral density GWAS associations. *bioRxiv*. 2021.10.12.464046. <https://doi.org/10.1101/2021.10.12.464046>.
46. Zhang Y, Chen XF, Li J, He F, Li X, Guo Y. lncRNA Neat1 stimulates Osteoclastogenesis via sponging miR-7. *J Bone Miner Res*. 2020;35:1772-1781.
47. Zhang Y, Chen B, Li D, Zhou X, Chen Z. lncRNA NEAT1/miR-29b-3p/BMP1 axis promotes osteogenic differentiation in human bone marrow-derived mesenchymal stem cells. *Pathol Res Pract*. 2019;215:525-531.
48. Yang X, Yang J, Lei P, Wen T. lncRNA MALAT1 shuttled by bone marrow-derived mesenchymal stem cells-secreted exosomes alleviates osteoporosis through mediating microRNA-34c/SATB2 axis. *Aging*. 2019;11:8777-8791.
49. Rom A, Melamed L, Gil N, et al. Regulation of CHD2 expression by the Chaserr long noncoding RNA gene is essential for viability. *Nat Commun*. 2019;10:5092.
50. Muñoz-Fuentes V, Cacheiro P, Meehan TF, et al. The International Mouse Phenotyping Consortium (IMPC): a functional catalogue of the mammalian genome that informs conservation. *Conserv Genet*. 2018;19:995-1005.
51. Reppe S, Refvem H, Gautvik VT, et al. Eight genes are highly associated with BMD variation in postmenopausal Caucasian women. *Bone*. 2010;46:604-612.
52. Guo HL, Wang X, Yang GY, et al. LINC00472 promotes osteogenic differentiation and alleviates osteoporosis by sponging miR-300 to upregulate the expression of FGFR2. *Eur Rev Med Pharmacol Sci*. 2020;24:4652-4664.