

1 **Identification of known and novel long non-coding RNAs potentially responsible**
2 **for the effects of BMD GWAS loci**

3
4 Abdullah Abood^{1,2}, Larry Mesner^{1,3}, Will Rosenow¹, Basel M. Al-Barghouthi^{1,2}, Nina
5 Horwitz⁵, Elise F. Morgan⁴, Louis C. Gerstenfeld⁵, Charles R. Farber¹⁻³

6
7 ¹ Center for Public Health Genomics, School of Medicine, University of Virginia,
8 Charlottesville, VA 22908

9 ² Department of Biochemistry and Molecular Genetics, School of Medicine, University of
10 Virginia, Charlottesville, VA 22908

11 ³ Department of Public Health Sciences, School of Medicine, University of Virginia,
12 Charlottesville, VA 22908

13 ⁴ Department of Mechanical Engineering, Boston University, Boston, MA, 02215

14 ⁵ Department of Orthopaedic Surgery, Boston University, Boston, MA, 02215

15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30

31 Correspondence to:

32
33
34
35
36
37
38
39

Charles R. Farber
E-mail: crf2s@virginia.edu
Center for Public Health Genomics
University of Virginia
P.O. Box 800717
Charlottesville, VA 22908, USA
Tel. 434-243-8584

40 **Abstract:**

41
42 Osteoporosis, characterized by low bone mineral density (BMD), is the most common
43 complex disease affecting bone and constitutes a major societal health problem.
44 Genome-wide association studies (GWASs) have identified over 1100 associations
45 influencing BMD. It has been shown that perturbations to long non-coding RNAs
46 (lncRNAs) influence BMD and the activities of bone cells; however, the extent to which
47 lncRNAs are involved in the genetic regulation of BMD is unknown. Here, we combined
48 the analysis of allelic imbalance (AI) in human acetabular bone fragments with a
49 transcriptome-wide association study (TWAS) and expression quantitative trait loci
50 (eQTL) colocalization analysis using data from the Genotype-Tissue Expression (GTEx)
51 project to identify lncRNAs potentially responsible for GWAS associations. We identified
52 27 lncRNAs in bone that are located in proximity to a BMD GWAS association and
53 harbor SNPs demonstrating AI. Using GTEx data we identified an additional 31
54 lncRNAs whose expression was associated (FDR correction <0.05) with BMD through
55 TWAS and had a colocalizing eQTL (regional colocalization probability (RCP) >0.1). The
56 58 lncRNAs are located in 43 BMD associations. To further support a causal role for the
57 identified lncRNAs, we show that 23 of the 58 lncRNAs are differentially expressed as a
58 function of osteoblast differentiation. Our approach identifies lncRNAs that are
59 potentially responsible for BMD GWAS associations and suggest that lncRNAs play a
60 role in the genetics of osteoporosis.

61 **Introduction:**

62

63 Osteoporosis is characterized by low bone mineral density (BMD) and deteriorated
64 structural integrity which leads to an increased risk of fracture ^{1,2}. In the U.S. alone, 12
65 million individuals have been diagnosed with osteoporosis, contributing to over 2 million
66 fractures per year ³. This number is expected to nearly double by 2025, resulting in
67 approximately \$26 billion in health care expenditures ³.

68

69 BMD is one of the strongest predictors of fracture ⁴ and is a highly heritable quantitative
70 trait ($h^2 = 0.5-0.8$) ⁵⁻⁸. As a result, the majority of genome-wide association studies
71 (GWASs) conducted for osteoporosis have focused on BMD. The largest BMD GWAS
72 performed to date used the UK BioBank (N~420K) and identified 1103 associations
73 influencing heel estimated BMD (eBMD) ⁹. One of the main challenges of BMD GWAS
74 is that the majority (>90%) of associations implicate non-coding variants that lie in
75 intronic or intergenic regions suggesting they have a role in gene regulation. This has
76 made it difficult to pinpoint causal genes and highlights the need for follow-up studies ¹⁰.
77 In addition, few studies have systematically evaluated non-coding transcripts as
78 potential causal genes.

79

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

89

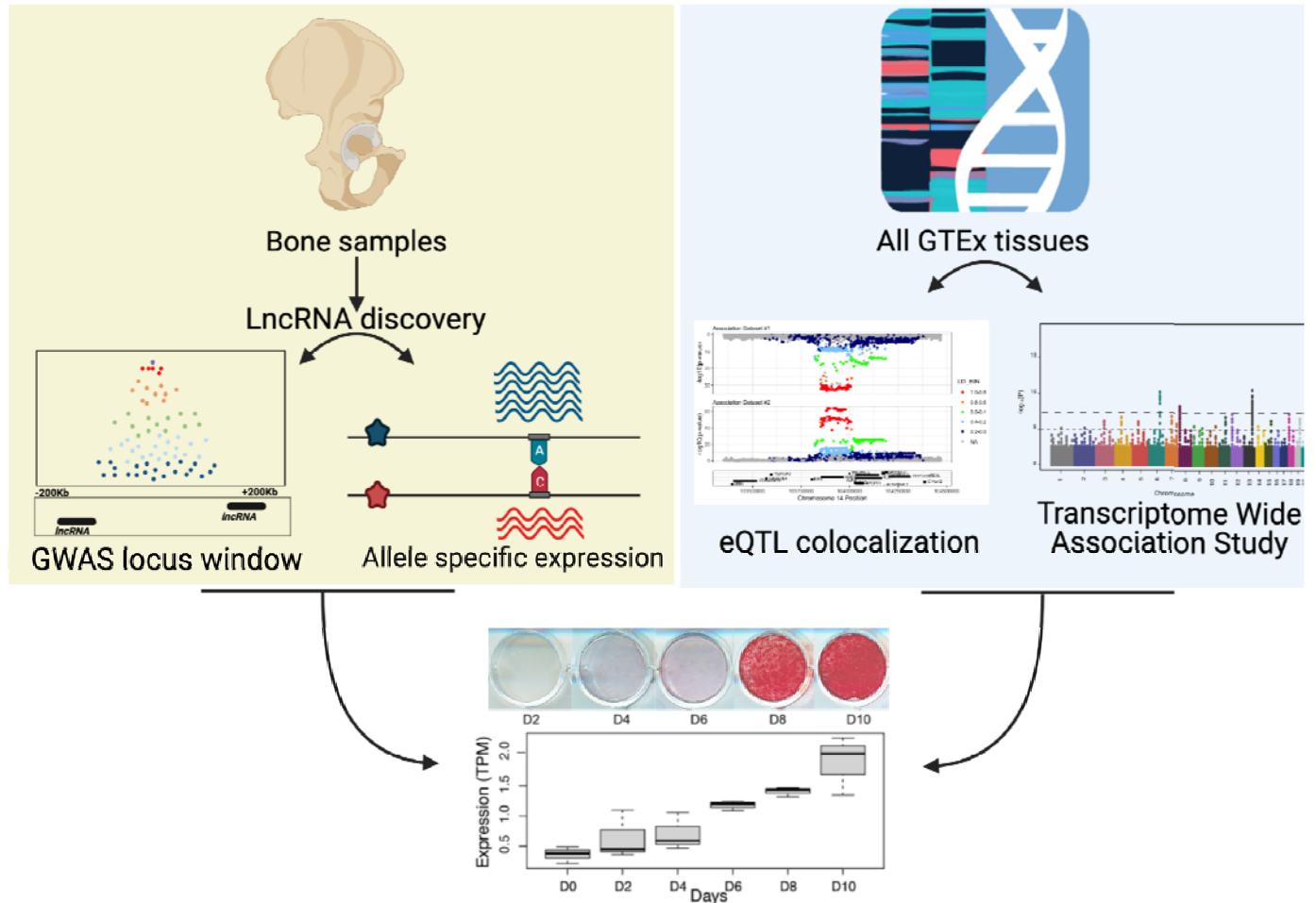
90 Although understudied in the context of GWAS ¹³, there is increasing evidence
91 suggesting that lncRNAs are causal for a subset of associations identified by GWAS. A
92 recent analysis of data from the Genotype-Tissue Expression (GTEx) project identified
93 690 potentially causal lncRNAs underlying associations influencing risk of a wide range
94 of diseases ¹³. Additionally, there is emerging evidence implicating lncRNAs in the
95 genetics of BMD ¹⁶⁻¹⁸. For example, a study reported 575 differentially expressed
96 lncRNAs between high and low BMD groups in Caucasian women, 26 of which regulate
97 protein-coding genes that are potentially causal in BMD GWAS ¹⁹. Additionally, a recent
98 BMD single nucleotide polymorphism (SNP) prioritization analysis implicated lncRNAs
99 as potential causal mediators ²⁰. Together these studies suggest that lncRNAs may play
100 an important role in the genetic regulation of bone mass.

101
102 In recent years, a number of approaches have been developed that utilize
103 transcriptomics data to inform GWAS, including the analysis of allelic imbalance (AI),
104 transcriptome-wide association studies (TWASs), and expression quantitative trait loci
105 (eQTL) colocalization²¹. AI results from the cis-regulatory effects (i.e., local eQTL) that
106 can be tracked using heterozygous coding SNPs. In transcriptome-wide association
107 studies (TWASs) the genetic component of gene expression in a reference population is
108 estimated and then imputed in a much larger population. Once gene expression is
109 imputed, genetically regulated gene expression is associated with a disease or disease
110 phenotype²². Most genes identified by TWAS are located in GWAS associations for that
111 disease and, as a result, TWAS can pinpoint genes likely to be causal at GWAS loci.
112 eQTLs are genetic variants associated with changes in gene expression and can be
113 tissue-specific or shared across multiple tissues. eQTL colocalization tests whether the
114 change in gene expression and the change in a trait of interest are driven by the same
115 shared genetic variant(s). All three approaches, alone or in combination, have been
116 successfully used to pinpoint potential causal disease genes at GWAS associations.

117
118 Here, we identified lncRNAs that are potentially responsible for the effects of BMD
119 GWAS associations by first applying AI to bone samples and, next, applying TWAS and
120 eQTL colocalization to gene expression data from GTEx. Through both approaches we
121 identified 58 lncRNAs with evidence of being causal BMD GWAS genes. We further
122 prioritized these lncRNAs by identifying those that were differentially expressed as a
123 function of osteoblast differentiation. Together, these results highlight the potential
124 importance of lncRNAs as candidate causal BMD GWAS genes.

125 126 **Results**

127
128 In this study, we used two approaches to identify lncRNAs that potentially underlie BMD
129 GWAS associations. In the first approach, we quantified known and novel lncRNAs
130 using RNA-seq data from human bone fragments and identified lncRNAs located in
131 proximity of a BMD GWAS association and harboring SNPs demonstrating AI. In the
132 second approach, we leveraged GTEx to identify lncRNAs across a large number of
133 tissues and cell-types whose expression was significantly associated with BMD by
134 TWAS and regulated by an eQTL which colocalized with a BMD association. **Figure 1**
135 provides an overview of our study.



136
137 *Figure 1: Overview of the study. We conducted de novo lncRNA discovery using RNA-seq data on human acetabular*
138 *bone fragments from 17 patients. We then identified known and novel lncRNAs located in GWAS associations that*
139 *were influenced by Allelic Imbalance (AI) (yellow box). We applied Transcriptome Wide Association Study (TWAS)*
140 *and colocalization on eQTL data from 49 Genotype-Tissue Expression (GTEx) project tissues (blue box). We*
141 *assessed the role of lncRNAs reported by both approaches in osteogenic differentiation using RNAseq data from the*
142 *human fetal osteoblast (hFOB) cell line at six time points across differentiation (bottom panel).*

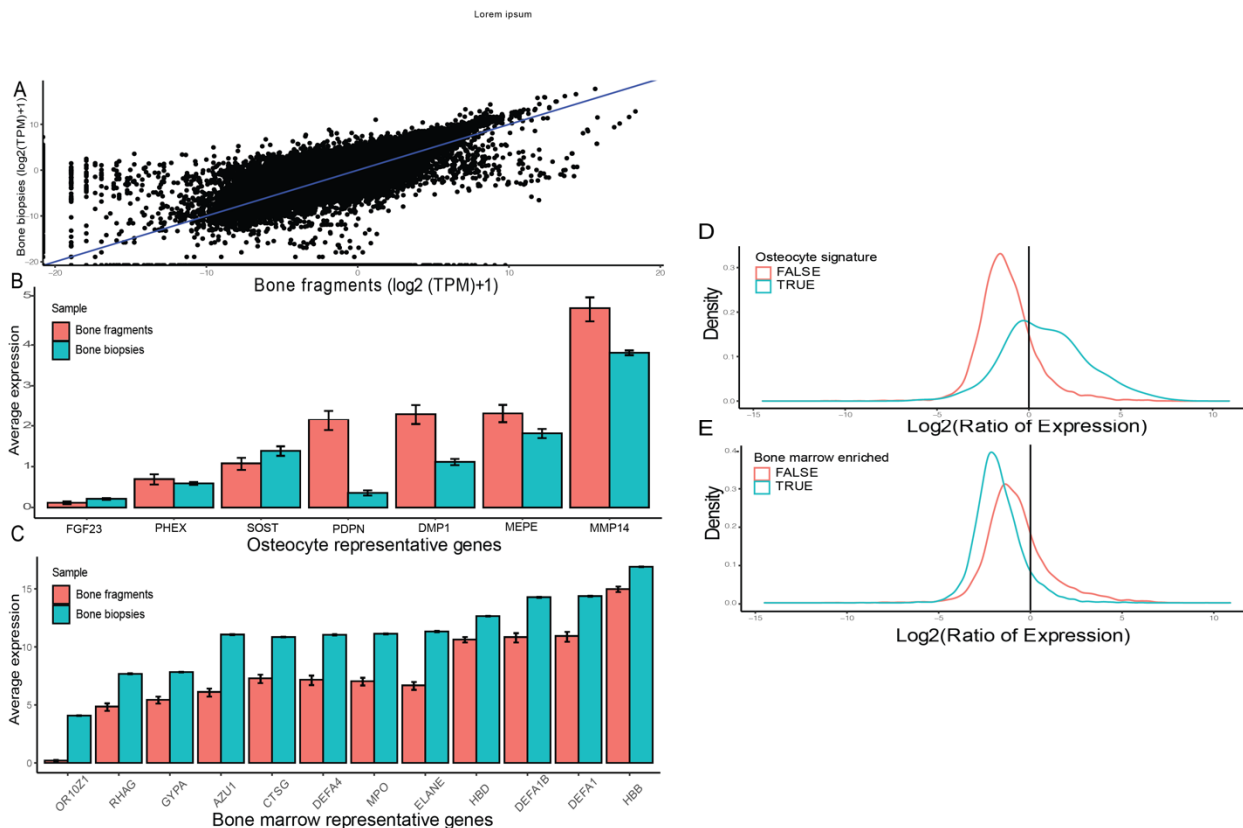
143 **Generation of bone expression data from bone fragments**

144

145 To identify potentially causal lncRNAs in a BMD relevant tissue, we generated total
146 RNA-seq (ribo-depleted) data on bone fragments isolated from acetabular reamings
147 from patients undergoing hip arthroplasty (N=17; 5 males and 12 females; ages 43 to
148 80). In contrast to most gene expression data generated on bone which are typically
149 from biopsies that contain marrow, we were able to remove the marrow leaving purified
150 trabecular and cortical bone. We hypothesized that the acetabular bone fragments
151 consisted primarily of late-stage osteoblasts/osteocytes²³, allowing us to characterize
152 lncRNAs enriched in these cell types. To confirm that the acetabular samples were
153 enriched in osteocytes, we compared these data to published RNA-seq data on bone
154 biopsies²⁴. Farr et al. generated RNA-seq data on 58 iliac crest needle biopsies from

155 healthy women containing both bone and marrow. Average transcripts per million (TPM)
156 across all samples in both experiments were highly correlated (**Figure 2A**, $r=0.845$, $P <$
157 2.2×10^{-16}). Importantly, differential expression analysis between the two datasets
158 showed that the top 1000 genes with the largest fold change increase in the bone
159 fragment samples compared to bone biopsy samples were enriched in Gene Ontology
160 (GO) terms such as “skeletal system development” ($FDR=4.01 \times 10^{-3}$) and “extracellular
161 matrix organization” ($FDR=4.11 \times 10^{-5}$).

162
163 To support the notion that our samples are unique in osteocyte enrichment, we used
164 data from a recent study that identified an osteocyte gene signature consisting of 1239
165 genes in mice and their orthologs in humans²⁵. The ratio of expression (bone fragment
166 samples / bone biopsy samples) was used. A ratio value > 1 indicates that gene
167 expression is higher in the bone fragment samples relative to the bone biopsy samples.
168 In contrast, a ratio value < 1 indicates that the gene is highly expressed in bone biopsy
169 samples relative to bone fragment samples. We expect to see that osteocyte signature
170 genes show ratio values > 1 and marrow enriched genes show ratio values < 1 . The
171 osteocyte signature genes showed a median ratio of 1.72 (62% of osteocyte signature
172 genes ratio > 1). Additionally, the ratio of expression of genes enriched in bone marrow
173 showed a median of 0.27 (91% of marrow enriched genes ratio < 1). The distribution of
174 osteocyte signature genes ratio values showed a significant median shift (Wilcoxon test,
175 $P < 2.2 \times 10^{-16}$) (**Figure 2D**), and the opposite pattern was observed for the bone
176 marrow enriched genes (Wilcoxon test, $P < 2.2 \times 10^{-16}$) (**Figure 2E**). These data
177 suggest that the purified acetabular bone fragments are enriched for late
178 osteoblasts/osteocytes compared to iliac crest biopsies.



179
 180 *Figure 2: Enrichment of osteocyte marker genes in bone fragment samples (used in this study) compared to bone*
 181 *biopsy samples in the literature. A) Overall gene expression is highly correlated between the RNA-seq data*
 182 *generated in both studies ($r=0.845$, $P < 2.2 \times 10^{-16}$) 24 B) Gene expression of osteocyte marker genes reported in 23*
 183 *showing enrichment in the bone fragments samples (this study) relevant to bone biopsies. C) Gene expression of*
 184 *bone marrow enriched genes reported in www.proteinatlas.org/ showing higher expression in bone biopsy samples.*
 185 *D) Osteocyte signature genes reported in Youlten et al. 25 are highly expressed in bone fragment samples relative to*
 186 *bone biopsies E) Bone marrow enriched genes reported in 25 are highly expressed in bone biopsy samples*
 187 *compared to bone fragment samples.*

188 Identifying novel lncRNAs in purified acetabular bone fragments

189
 190 Given the paucity of bone transcriptomics data in the literature, and the tissue-specific
 191 nature of lncRNA expression, we hypothesized that many bone/osteocyte specific
 192 lncRNAs would not be present in current sequence databases. Additionally, ~50% of
 193 lncRNAs do not possess a poly-A tail modification²⁶ and most RNA-seq data is
 194 generated after poly-A selection. Therefore, in order to capture a more comprehensive
 195 profile of lncRNAs in bone, we implemented a lncRNA discovery step to identify putative
 196 “novel” lncRNA transcripts using the computational algorithm CPAT²⁷. Across the 17
 197 bone samples we identified 6612 known lncRNAs and 2440 novel lncRNAs
 198 (Supplementary tables 1 and 2). The mean length of novel lncRNAs was 30.3 Kb and
 199 median length of 11.8 Kb. These values were comparable to the mean length of known
 200 lncRNAs expressed in the bone samples (mean = 35.4 Kb; median = 4.7 Kb).

201

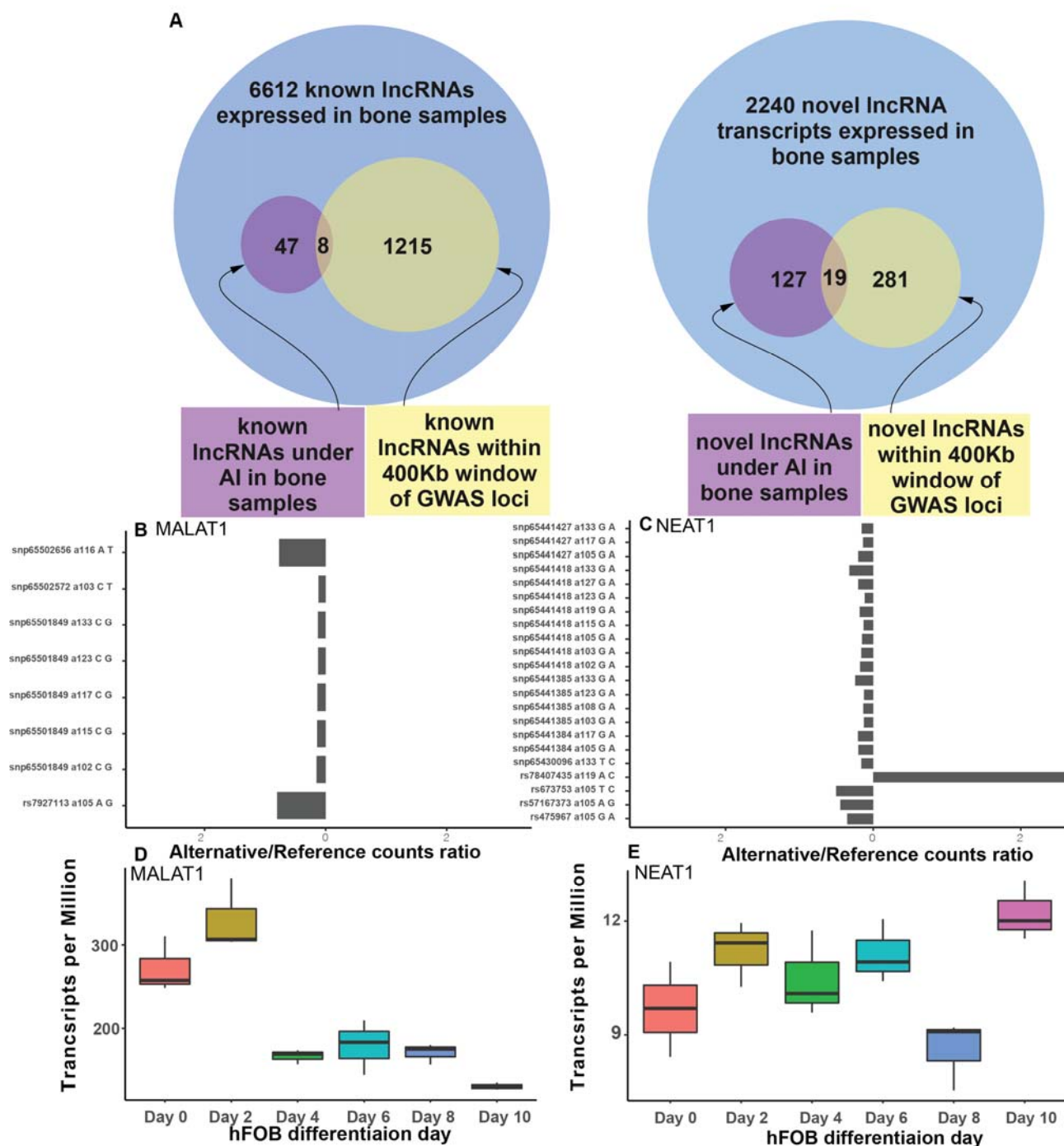
202 Identifying potentially casual lncRNAs in bone

203

204 For lncRNAs to be considered potentially causal in bone, we identified those that are
205 both located in proximity of a BMD GWAS association and regulated by AI. We
206 hypothesized that such genes may be causal for their respective associations because
207 of the potential to be regulated by an eQTL which colocalizes with a BMD association.
208 Of the 9,052 lncRNAs (2440 novel and 6612 known) we quantified in acetabular bone,
209 1,496 lncRNAs (~17% of expressed lncRNAs) were found within a 400Kb window (\pm
210 200Kb from the lncRNA start site) of each of 1103 GWAS associations previously
211 identified by Morris et al.⁹.

212

213 Next, we identified heterozygous coding variants that demonstrated significant evidence
214 of AI within lncRNAs. Of the total number of lncRNAs we identified, 174 (47 known, 127
215 novel; ~2% of expressed lncRNAs) had at least one SNP demonstrating AI in at least
216 one of the 17 bone fragment samples. Out of the 174, 27 (15.5%; 8 known, 19 novel)
217 were located in proximity of a GWAS association (**Figure 3A, Supplementary Table 3**).



218
 219
 220 *Figure 3: Identification of lncRNAs located within eBMD GWAS associations, are under AI in acetabular bone, and*
 221 *are differentially expressed in hFOBs. A) Venn diagram showing the number of known and novel lncRNAs within*
 222 *proximity of GWAS loci, implicated by AI, and implicated by both approaches. B) lncRNA MALAT1 AI plot showing*
 223 *the ratio of reads aligning to the alternative SNP relative to the reference SNP in eight of the bone fragments*
 224 *samples where the gene is under AI. C) lncRNA NEAT1 AI plot showing the ratio of reads aligning to the alternative*
 225 *SNP relative to the reference SNP in ten of the bone fragments samples where the gene is under AI. rs78407435 is not in*
 226 *LD with the rest of the SNPs in the region and this is likely the reason it shows a different direction of effect. D)*
 227 *Expression of MALAT1 across hFOB differentiation points. E) Expression of NEAT1 across hFOB differentiation*

228 **Identifying putatively causal lncRNAs by leveraging GTEx**

229
230 Next, we sought to leverage non-bone data to identify potentially causal lncRNAs. To do
231 this, we integrated 1103 BMD GWAS loci⁹ with GTEx (v8) eQTL data by coupling
232 TWAS²⁸ using S-MultiXScan²⁹ and Bayesian colocalization analysis using fastENLOC
233³⁰. The rationale behind using GTEx data is genes that are shared in multiple tissues
234 and showing a colocalizing eQTL with BMD GWAS data can be potentially causal in
235 bone tissue as well. Our TWAS analysis resulted in 333 significant lncRNA-BMD
236 associations (FDR correction < 0.05). Our colocalization analysis yielded 48 lncRNAs
237 with a colocalizing eQTL (RCP > 0.1) in at least one GTEx tissue. There were 31
238 lncRNAs significant in both the TWAS and eQTL colocalization analysis
239 **(Supplementary Table 4)**.

240 241 **Most identified lncRNAs are the only potential causal mediators implicated by** 242 **TWAS/eQTL colocalization in their respective GWAS associations**

243
244 To determine if the lncRNAs listed in **Supplementary Table 4** are the strongest
245 candidates in their respective GWAS associations, we evaluated a recent report of
246 protein coding genes that used the same approach³¹. Five out of the 31 lncRNAs
247 (*LINC01116*, *LINC01117*, *SNHG15*, *LINC01290*, *LINC00665*) have a protein coding
248 gene with a colocalizing eQTL (*HOXD8*, *HOXD9*, *MYO1G*, *NACAD*, *EMP2*, *ZFP14*,
249 *ZFP82*) within 1 Mb of the lncRNA start site (**Supplementary Table 5**). Upon further
250 investigation of the RCP values, some of the lncRNAs showed higher RCP than their
251 protein coding gene counterpart. For example, *LINC01290* had a higher RCP in lung
252 tissue (0.4992) compared to its counterpart *EMP2* (0.2227). On the other hand, the
253 same lncRNA has a lower RCP value (0.1498) than *EMP2* (0.6089) in breast and
254 mammary gland tissue. However, for the remaining lncRNAs, this analysis provides
255 support that the lncRNA alone is the potential causal mediator in the region as we show
256 no evidence of protein coding colocalization within 1 Mb distance of the start site of the
257 lncRNA.

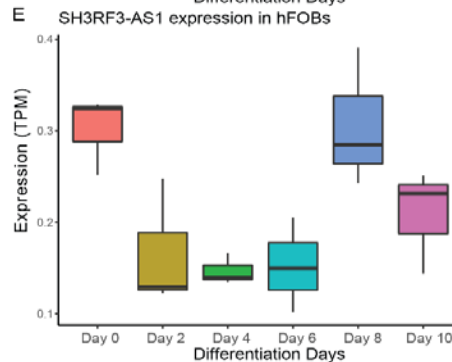
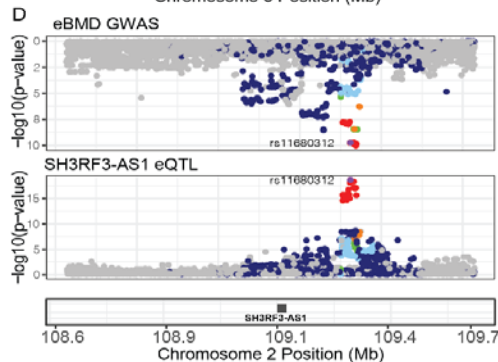
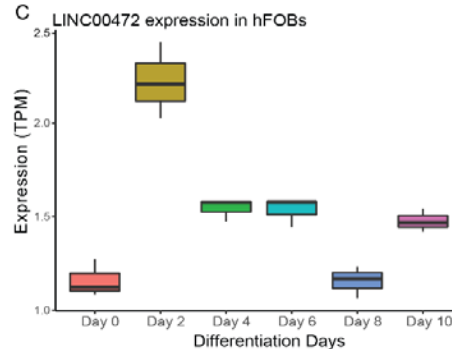
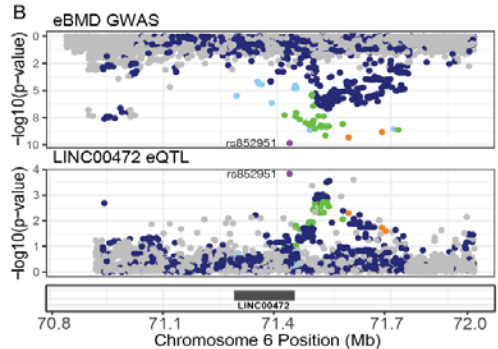
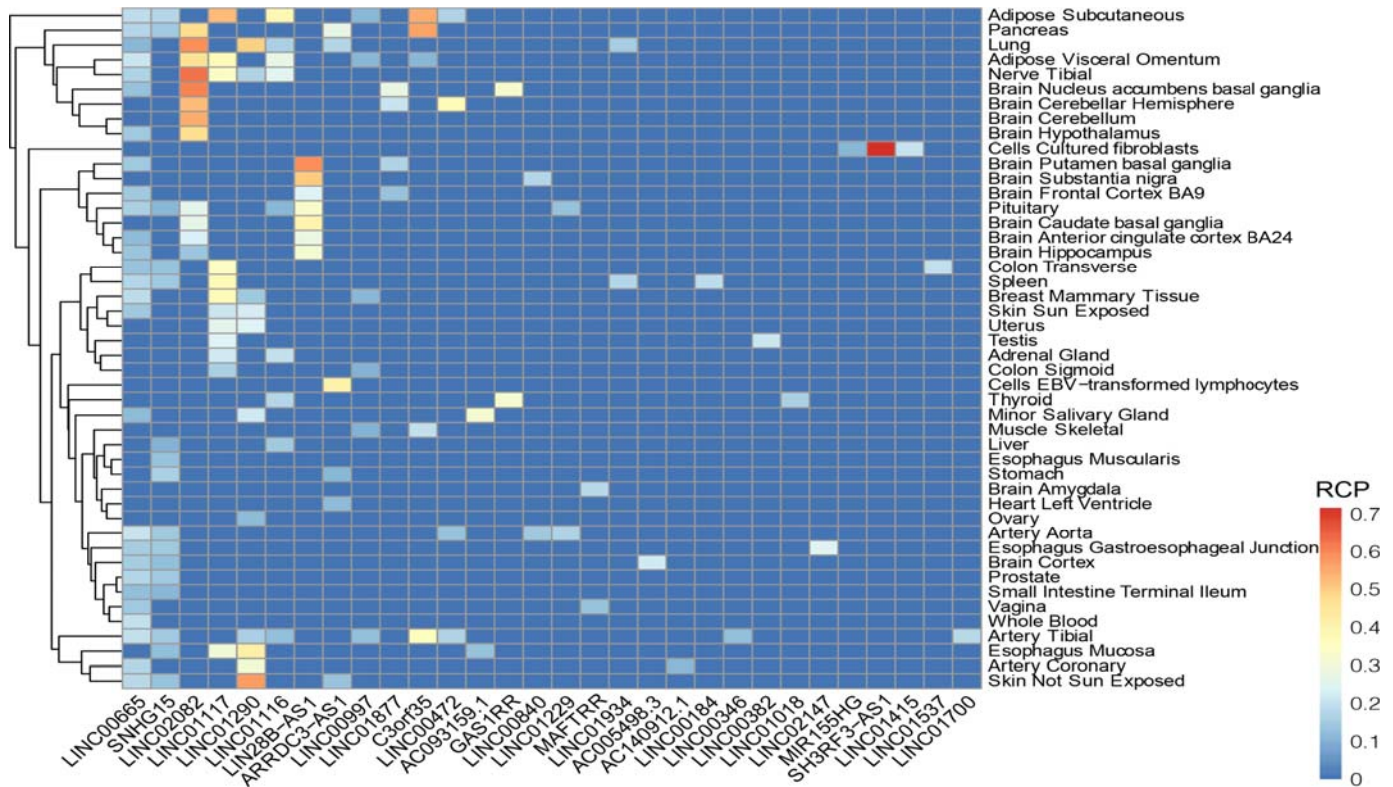
258 259 **Many identified lncRNAs are differentially expressed as a function of osteoblast** 260 **differentiation**

261
262 To provide further support for the hypothesis that these lncRNAs mediate GWAS
263 associations, we measured their expression as a function of osteoblast differentiation in
264 hFOB. We performed total RNA-seq at six hFOB differentiation time-points (Days 0, 2,
265 4, 6, 8, and 10). Of the 27 lncRNAs implicated in the analysis of AI, all eight known
266 lncRNAs were differentially expressed (FDR<0.05). On the other hand, none of the
267 novel lncRNAs were differentially expressed (**Supplementary Table 3**). Examples of

268 the identified genes include *MALAT1* and *NEAT1* (**Figure 3B** and **3C**), which were
269 differentially expressed in hFOBs and showed evidence of AI in 8 and 10 of the 17
270 acetabular bone samples, respectively. There were four unique SNPs in the exonic
271 regions of *MALAT1* (**Figure 3B**) that were heterozygous in at least one of the 17
272 individuals (with a maximum of 8 individuals). All four SNPs showed higher expression
273 in the alternative allele relative to the reference allele. The expression of *MALAT1* gene
274 decreased as the cell differentiated into a mineralizing state. Additionally, there were
275 nine unique SNPs reported in the exonic regions of *NEAT1* that were heterozygous in at
276 least one of the 17 individuals (with a maximum of 10 individuals). Of the nine, eight
277 showed higher expression associated with the alternative allele compared to the
278 reference allele. The remaining SNP was associated with the opposite pattern and this
279 was likely due to it being the only SNP not in high LD with the others ($R^2 = 0.0021$).
280 *NEAT1* showed significant increase in expression around day 10 in hFOBs.

281
282 We assessed the expression of lncRNAs identified by GTEx TWAS/eQTL colocalization
283 in osteoblast differentiation using the same approach in the previous section. Out of the
284 31 lncRNAs identified by TWAS/eQTL colocalization, 15 were found to be differentially
285 expressed (*LINC00184*, *SH3RF3-AS1*, *LINC01116*, *LINC01934*, *C3orf35*, *LINC01018*,
286 *ARRDC3-AS1*, *LINC00472*, *SNHG15*, *GAS1RR*, *LINC00840*, *LINC01537*, *LINC00346*,
287 *LINC01415*, *MIR155HG*). In general, the expression of those genes in hFOBs was low
288 compared to the lncRNAs reported in the AI section. Examples include *SHR3F3-AS1*
289 and *LINC00472*, which were regulated by colocalizing eQTL (**Figure 4 B and D**) and
290 were differentially expressed in hFOBs. (**Figure 4 C and E**). *SH3RF3-AS1* was shown
291 to have the highest RCP value overall (RCP= 0.72) and in only one GTEx tissue
292 (cultured fibroblasts) (**Figures 4A and 4D, Table 2**). While the gene was differentially
293 expressed across hFOB differentiation points, it had a very low overall level of
294 expression (**Figure 4E**). The pattern of expression decreased during mid differentiation
295 points with spikes in early and late points (**Figure 4E**). *LINC00472* was shown to have a
296 colocalizing eQTL in four GTEx tissues with the highest RCP value in brain cerebellar
297 hemisphere (RCP = 0.37) (**Figures 4A and 4B, Table 2**). The gene also showed a
298 moderate level of expression in hFOBs with an average of 1.5 TPM (**Figure 4C**). The
299 expression of *LINC00472* peaked at day 2 and then declined (**Figure 4C**).

300



301

302

303

304

305

306

307

Figure 4: lncRNAs implicated by eQTL colocalization and TWAS are potential causal mediators of BMD GWAS loci. A) Heatmap 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.

308 Discussion

309

310 In this study, we interrogated BMD GWAS loci and identified known and novel lncRNAs
311 as potential causal mediators. We identified potentially important lncRNA using two
312 different approaches. First, we identified novel and known lncRNAs in a unique
313 transcriptomic bone dataset that were localized in GWAS loci and demonstrated AI.
314 Second, we implicated additional lncRNAs by leveraging GTEx and identifying eQTLs in
315 non-bone tissues that colocalized with eBMD GWAS loci whose expression was
316 associated with eBMD via TWAS. We also assessed differential expression across the
317 time course of hFOB differentiation to provide more evidence of a potential causal role
318 for these lncRNAs.

319

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

326

327 A total of eight lncRNAs (*NEAT1*, *MALAT1*, *DLEU2*, *LINC01578*, *CARMN*, *AC011603.3*,
328 *PXN-AS1*, *AC020656.1*) were found to be within a 400 Kb window of an eBMD GWAS
329 locus and were also differentially expressed across hFOB differentiation time points.
330 Many of these lncRNAs have been demonstrated to play a role in bone. For example,
331 *NEAT1* has been reported to stimulate osteoclastogenesis via sponging miR-7³² and
332 *NEAT1*/miR-29b-3p/BMP1 axis promotes osteogenic differentiation in human bone
333 marrow-derived mesenchymal stem cells³³. In addition, *MALAT1* has been shown to
334 influence BMD³⁴. *MALAT1* acts as a sponge of miR-34c to promote the expression of
335 *SATB2*. *SATB2* then acting to reduce the ALP activity of osteoblasts and mineralized
336 nodules formation³⁴. A recent study has shown that *LINC01578* (referred to as
337 *CHASERR* in this study) represses chromodomain Helicase DNA Binding Protein 2
338 (*Chd2*). A model for *Chd2* loss of function by the International Mouse Phenotyping
339 Consortium (IMPC)³⁵ reported that these mice exhibit significant decreased body
340 weight and length, skeletal abnormalities, abnormal bone structure, decreased fat levels
341 and bone mineral density³⁶. Lastly, *DLEU2* expression has been shown to be inversely
342 correlated with BMD in a study involving postmenopausal Caucasian women³⁷. The
343 remaining four lncRNAs have not been reported to date to have a role in bone and
344 should be further pursued.

345

346 In our second analysis, we reported 15 lncRNAs implicated jointly by colocalization,
347 TWAS, and differential expression analysis. We show one example of the 15 lncRNAs

348 reported *SH3RF3-AS1* in **(Figure 4A)**. Most of these lncRNAs have not been shown
349 previously in the literature to have a role in bone biology. However, *LINC00472* **(Figure**
350 **4B)** has been experimentally shown to influence osteogenic differentiation by sponging
351 miR-300 which in turn increases the expression of *Fgfr2* in mice³⁸. These preliminary
352 results provide more evidence to the potential causal role of these lncRNAs in
353 osteoporosis.

354
355 This study is not meant to be comprehensive as we are limited by the number of
356 samples and are not suitably powered to identify eQTLs and apply TWAS/colocalization
357 analysis. However, due to the scarcity of population-level bone transcriptomic dataset,
358 and the lack of bone cell or tissue data in GTEx, our study is an attempt to
359 systematically leverage the available datasets to capture a subset of lncRNAs that we
360 think are potentially causal. As mentioned, some of these lncRNAs have been
361 implicated experimentally outside of this study. Moreover, lncRNAs under AI and within
362 proximity of GWAS loci may not be causal as they could be false positives because they
363 are not prioritized via a systems analysis like colocalization. Another limitation of our
364 study is that we evaluated their expression as a function of osteoblast differentiation;
365 however, it is likely that some of the lncRNAs, if truly causal, impact BMD via a function
366 in other cell-types (e.g., osteoclasts). Future studies should focus on enhancing these
367 results by generating transcriptomic and eQTL datasets from bone and other bone cell
368 types, using network approaches to aid in the prioritization of lncRNAs, and
369 experimentally validating the role of specific lncRNAs.

370
371 In this study, we were able to use multiple systems genetics approaches on two
372 transcriptomic datasets (acetabular bone and GTEx) to identify lncRNAs that are
373 potentially responsible for the effects of some BMD GWAS loci. This is the first study to
374 our knowledge that evaluated the role of lncRNAs in mediating the effect of BMD GWAS
375 loci from a genome-wide perspective. We combined osteoblast differentiation samples
376 and the literature to provide experimental evidence in previous studies to support the
377 causal mediator list we generated from our analysis. These results highlight the
378 importance of studying other aspects of the transcriptome to identify potential drug
379 targets for osteoporosis and bone fragility.

380
381 **Data availability statement:**
382 Analysis code is available on GitHub [https://github.com/aa9gj/lncRNA_publication].
383 Raw samples are submitted to Gene Expression Omnibus
384 [<https://www.ncbi.nlm.nih.gov/geo/>] reference number [GSE186922].

385
386 **Acknowledgements:**

387 Research reported in this publication was supported in part by the National Institute of
388 Arthritis and Musculoskeletal and Skin Diseases of the National Institutes of Health
389 under Award Number AR071657 to Charles R. Farber, Louis C. Gerstenfeld, and Elise
390 F. Morgan, and Abdullah Abood was supported in part by a National Institutes of Health,
391 Biomedical Data Sciences Training Grant (5T32LM012416). The authors acknowledge
392 Dr. Emily Farber for generating RNAseq data on bone fragments. We thank the IMPC
393 for accessibility to BMD data in knockout mice (www.mousephenotype.org). The
394 Genotype-Tissue Expression (GTEx) Project was supported by the Common Fund of
395 the Office of the Director of the National Institutes of Health, and by NCI, NHGRI,
396 NHLBI, NIDA, NIMH, and NINDS. The data used for the analyses described in this
397 manuscript were obtained from the GTEx Portal on 6/30/20.

398

399 **Methods**

400

401 **Patient demographics**

402

403 All human specimen collection was performed in accordance with IRB approval from our
404 institution (IRB number H-32517). Acetabular reaming from 17 Boston Medical Center
405 (BMC) patients (ages 43-80 year) undergoing elective hip arthroplasty were collected:
406 12 Females and 5 Males; 8 Black, 8 White, and 1 Hispanic. This demographic mix
407 reflects the population serviced by BUMC, which is an urban safety-net hospital.

408

409 **RNA extraction**

410

411 Bone fragments were isolated from the 17 patients. Total RNA was isolated from bone
412 fragments as previously described in ³⁹. Total RNA-Seq libraries were constructed from
413 bone as well as hFOB RNA samples using Illumina TruSeq Stranded Total RNA with
414 Ribo-Zero Gold sample prep kits. Constructed libraries contained all RNAs greater than
415 100 nt (both unpolyadenylated and polyadenylated) minus cytoplasmic and
416 mitochondrial rRNAs. Samples were sequenced to achieve a minimum of 50 million
417 reads 2 x 75 bp paired-end reads on an Illumina NextSeq500.

418

419 **Human fetal osteoblast (hFOB) cell line culture**

420

421 hFOB 1.19 cells (ATCC #CRL-11372) were cultured at 34C and differentiated at 39.5C
422 as recommended with the following modifications. Growth media: Minimal Essential
423 Media (MEM, Gibco 10370-021) supplemented with 10% Fetal Bovine Serum (FBS,
424 Atlantic Biological S12450), 1% Glutamax (Gibco 35050-061), 1% Pen Strep (Gibco
425 15140-122). Differentiation Media: MEM alpha (Gibco 12571-063) supplemented with
426 10% FBS, 1% Glutamax, 1% Pen Strep, 50ug/ul Ascorbic Acid (Sigma A4544-25G),

427 10mM beta-Glycerophosphate (Sigma G9422-100G), 10nM Dexamethasone (sigma
428 D4902-25MG). RNA was isolated from $\sim 0.5 \times 10^6$ cells at days 0, 2, 4, 6, 8 and 10 of
429 differentiation as recommended (RNAeasy Minikit, Qiagen 74106). Mineralized nodule
430 formation was measured by staining cultures with Alizarin Red (40 mM, pH 5.6; Sigma
431 A5533-25G). Reported results were obtained from three biological replicate
432 experiments.

433

434 **RNA sequencing and Differential Gene Expression analysis**

435

436 Computational analysis of RNA sequencing data for the 17 bone samples, Farr et al.²⁴
437 and the hFOB samples were performed using a custom bioinformatics pipeline. Briefly,
438 FastqQC (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) and RSeQC⁴⁰
439 were used to assess the quality of raw reads. Adapter trimming was completed using
440 Trimmomatic⁴¹. Sequences were aligned to the GENCODE v34⁴² reference genome
441 using the SNP and splice aware aligner HISAT2⁴³. Genome assembly and abundances
442 in transcripts per million (TPM) were quantified using StringTie⁴⁴. Differential
443 expression analysis for the hFOB differentiation experiment was performed using
444 DEseq2⁴⁵ across all six differentiation time points using analysis of deviance
445 (ANODEV) which is conceptually similar to analysis of variance (ANOVA). Differential
446 expression analysis for the comparison between this study's samples and the samples
447 in the literature was performed using DEseq2⁴⁵ standard approach.

448

449 **lncRNA discovery**

450

451 The Coding Potential Assessment Tool (CPAT)²⁷ was used to assess the protein-
452 coding potential of the novel transcripts assembled. In short, CPAT is a machine
453 learning algorithm trained on a set of known human lncRNAs to identify novel putative
454 lncRNAs based on shared sequence features. We used all known lncRNAs in the latest
455 human genome assembly (GRCh38) as the training set. Novel transcripts with coding
456 probability < 0.367 are regarded as lncRNAs in accordance with software authors.
457 Novel lncRNAs with TPM < 1 were regarded as noise and discarded.

458

459 **Allelic Imbalance analysis**

460

461 Reads were aligned to the GENCODE v34⁴² reference genome using the SNP and
462 splice aware aligner HISAT2⁴³. The resultant BAM files were then used as input for
463 variant calling using the GATK pipeline⁴⁶. Briefly, duplicate reads were identified using
464 MarkDuplicates. Next, reads spanning introns were reformatted using SplitNCigarReads
465 to match the DNA aligner conventions. Then base quality recalibration was performed to
466 detect and correct for patterns of systematic errors in the base quality scores. Finally,

467 the variant calling and filtration step was performed using HaplotypeCaller. The
468 resultant vcf file included only known and novel snps and reference bias was corrected
469 using WASP⁴⁷. Briefly, mapped reads that overlap SNPs are identified. For each read
470 that overlaps a SNP, its genotype is swapped with that of the other allele and it is re-
471 mapped. If a re-mapped read fails to map to exactly the same location, it is discarded.
472 The resultant corrected BAM and filtered VCF files were used as input for GATK
473 ASEReadCounter to provide a table of filtered base counts at heterozygous sites for
474 allele specific expression. Bases with a read depth less than 20 were discarded. In
475 order to determine significance, a binomial test was performed and only heterozygous
476 sites with FDR corrected p-value of <0.05 were considered significant.

477

478 **Transcriptome Wide Association Studies**

479

480 We conducted a transcriptome-wide association study by integrating genome-wide
481 SNP-level association summary statistics from a bone mineral density GWAS⁹ with
482 GTEx version 8 gene expression QTL data from 49 tissue types. We used the S-
483 MultiXcan²⁹ approach for this analysis, to correlate gene expression across tissues to
484 increase power and identify candidate susceptibility genes. Gene-level associations
485 were identified at FDR correction < 0.05 and were further filtered using fastENLOC
486 (described in for a regional colocalization probability > 0.1 in at least one tissue type.

487

488 **Bayesian colocalization analysis**

489

490 We used fastENLOC, a faster implementation of ENLOC³⁰ to perform Bayesian
491 colocalization analysis. We integrated summary statistics from the most recent (and
492 largest) eBMD GWAS⁹ and eQTL data from 49 GTEx tissues⁴⁸. We used the
493 recommended regional colocalization probability (RCP) threshold of >0.1 as indication
494 of significant overlap between SNP and eQTL.

495 References

- 496
- 497 1. NIH Consensus Development Panel on Osteoporosis Prevention, Diagnosis, and
498 Therapy. Osteoporosis prevention, diagnosis, and therapy. *JAMA* **285**, 785–795
499 (2001).
 - 500 2. Office of the Surgeon General (US). *Bone Health and Osteoporosis: A Report of*
501 *the Surgeon General*. (Office of the Surgeon General (US), 2010).
 - 502 3. Burge, R. *et al.* Incidence and economic burden of osteoporosis-related fractures in
503 the United States, 2005-2025. *J. Bone Miner. Res.* **22**, 465–475 (2007).
 - 504 4. Johnell, O. *et al.* Predictive value of BMD for hip and other fractures. *J. Bone Miner.*
505 *Res.* **20**, 1185–1194 (2005).
 - 506 5. Smith, D. M., Nance, W. E., Kang, K. W., Christian, J. C. & Johnston, C. C., Jr.
507 Genetic factors in determining bone mass. *J. Clin. Invest.* **52**, 2800–2808 (1973).
 - 508 6. Arden, N. K., Baker, J., Hogg, C., Baan, K. & Spector, T. D. The heritability of bone
509 mineral density, ultrasound of the calcaneus and hip axis length: a study of
510 postmenopausal twins. *J. Bone Miner. Res.* **11**, 530–534 (1996).
 - 511 7. Slemenda, C. W. *et al.* The genetics of proximal femur geometry, distribution of
512 bone mass and bone mineral density. *Osteoporos. Int.* **6**, 178–182 (1996).
 - 513 8. Richards, J. B., Zheng, H.-F. & Spector, T. D. Genetics of osteoporosis from
514 genome-wide association studies: advances and challenges. *Nat. Rev. Genet.* **13**,
515 576–588 (2012).
 - 516 9. Morris, J. A. *et al.* An atlas of genetic influences on osteoporosis in humans and
517 mice. *Nat. Genet.* **51**, 258–266 (2019).
 - 518 10. Claussnitzer, M. *et al.* A brief history of human disease genetics. *Nature* **577**, 179–
519 189 (2020).
 - 520 11. Zhang, J., Hao, X., Yin, M., Xu, T. & Guo, F. Long non-coding RNA in
521 osteogenesis: A new world to be explored. *Bone Joint Res.* **8**, 73–80 (2019).
 - 522 12. Marchese, F. P., Raimondi, I. & Huarte, M. The multidimensional mechanisms of
523 long noncoding RNA function. *Genome Biol.* **18**, 206 (2017).
 - 524 13. de Goede, O. M. *et al.* Long non-coding RNA gene regulation and trait associations
525 across human tissues. *bioRxiv* 793091 (2019) doi:10.1101/793091.
 - 526 14. Silva, A. M. *et al.* Long noncoding RNAs: a missing link in osteoporosis. *Bone Res*
527 **7**, 10 (2019).
 - 528 15. Nardocci, G. *et al.* Identification of a novel long noncoding RNA that promotes
529 osteoblast differentiation. *J. Cell. Biochem.* **119**, 7657–7666 (2018).
 - 530 16. Chen, X.-F. *et al.* An Osteoporosis Risk SNP at 1p36.12 Acts as an Allele-Specific
531 Enhancer to Modulate LINC00339 Expression via Long-Range Loop Formation.
532 *Am. J. Hum. Genet.* **102**, 776–793 (2018).
 - 533 17. Roca-Ayats, N. *et al.* Functional characterization of the C7ORF76 genomic region,
534 a prominent GWAS signal for osteoporosis in 7q21.3. *Bone* **123**, 39–47 (2019).
 - 535 18. Mei, B. *et al.* LncRNA ZBTB40-IT1 modulated by osteoporosis GWAS risk SNPs
536 suppresses osteogenesis. *Hum. Genet.* **138**, 151–166 (2019).
 - 537 19. Zhou, Y. *et al.* Long Noncoding RNA Analyses for Osteoporosis Risk in Caucasian
538 Women. *Calcif. Tissue Int.* **105**, 183–192 (2019).
 - 539 20. Zhang, X., Deng, H.-W., Shen, H. & Ehrlich, M. Prioritization of osteoporosis-
540 associated genome-wide association study (GWAS) single-nucleotide

- 541 polymorphisms (SNPs) using epigenomics and transcriptomics. *JBMR Plus* **5**,
542 e10481 (2021).
- 543 21. Hukku, A. *et al.* Probabilistic Colocalization of Genetic Variants from Complex and
544 Molecular Traits: Promise and Limitations. *Cold Spring Harbor Laboratory*
545 2020.07.01.182097 (2020) doi:10.1101/2020.07.01.182097.
- 546 22. Abood, A. & Farber, C. R. Using “-omics” Data to Inform Genome-wide Association
547 Studies (GWASs) in the Osteoporosis Field. *Curr. Osteoporos. Rep.* (2021)
548 doi:10.1007/s11914-021-00684-w.
- 549 23. Bonewald, L. F. The amazing osteocyte. *J. Bone Miner. Res.* **26**, 229–238 (2011).
- 550 24. Farr, J. N. *et al.* Effects of Age and Estrogen on Skeletal Gene Expression in
551 Humans as Assessed by RNA Sequencing. *PLoS One* **10**, e0138347 (2015).
- 552 25. Youlten, S. E. *et al.* Osteocyte transcriptome mapping identifies a molecular
553 landscape controlling skeletal homeostasis and susceptibility to skeletal disease.
554 *Nat. Commun.* **12**, 2444 (2021).
- 555 26. Cheng, J. *et al.* Transcriptional maps of 10 human chromosomes at 5-nucleotide
556 resolution. *Science* **308**, 1149–1154 (2005).
- 557 27. Wang, L. *et al.* CPAT: Coding-Potential Assessment Tool using an alignment-free
558 logistic regression model. *Nucleic Acids Res.* **41**, e74 (2013).
- 559 28. Gusev, A. *et al.* Integrative approaches for large-scale transcriptome-wide
560 association studies. *Nat. Genet.* **48**, 245–252 (2016).
- 561 29. Barbeira, A. N. *et al.* Integrating Predicted Transcriptome From Multiple Tissues
562 Improves Association Detection. doi:10.1101/292649.
- 563 30. Wen, X., Pique-Regi, R. & Luca, F. Integrating molecular QTL data into genome-
564 wide genetic association analysis: Probabilistic assessment of enrichment and
565 colocalization. *PLoS Genet.* **13**, e1006646 (2017).
- 566 31. Al-Barghouthi, B. M. *et al.* Transcriptome-wide Association Study and eQTL
567 colocalization identify potentially causal genes responsible for bone mineral density
568 GWAS associations. (2021) doi:10.1101/2021.10.12.464046.
- 569 32. Zhang, Y. *et al.* lncRNA Neat1 Stimulates Osteoclastogenesis Via Sponging miR-7.
570 *J. Bone Miner. Res.* **35**, 1772–1781 (2020).
- 571 33. Zhang, Y., Chen, B., Li, D., Zhou, X. & Chen, Z. LncRNA NEAT1/miR-29b-
572 3p/BMP1 axis promotes osteogenic differentiation in human bone marrow-derived
573 mesenchymal stem cells. *Pathol. Res. Pract.* **215**, 525–531 (2019).
- 574 34. Yang, X., Yang, J., Lei, P. & Wen, T. LncRNA MALAT1 shuttled by bone marrow-
575 derived mesenchymal stem cells-secreted exosomes alleviates osteoporosis
576 through mediating microRNA-34c/SATB2 axis. *Aging* vol. 11 8777–8791 (2019).
- 577 35. Muñoz-Fuentes, V. *et al.* The International Mouse Phenotyping Consortium (IMPC):
578 a functional catalogue of the mammalian genome that informs conservation.
579 *Conserv. Genet.* **19**, 995–1005 (2018).
- 580 36. Rom, A. *et al.* Regulation of CHD2 expression by the Chaserr long noncoding RNA
581 gene is essential for viability. *Nature Communications* vol. 10 (2019).
- 582 37. Reppe, S. *et al.* Eight genes are highly associated with BMD variation in
583 postmenopausal Caucasian women. *Bone* **46**, 604–612 (2010).
- 584 38. Guo, H.-L. *et al.* LINC00472 promotes osteogenic differentiation and alleviates
585 osteoporosis by sponging miR-300 to upregulate the expression of FGFR2. *Eur.*
586 *Rev. Med. Pharmacol. Sci.* **24**, 4652–4664 (2020).

- 587 39. Sagi, H. C., Young, M. L., Gerstenfeld, L., Einhorn, T. A. & Tornetta, P. Qualitative
588 and quantitative differences between bone graft obtained from the medullary canal
589 (with a Reamer/Irrigator/Aspirator) and the iliac crest of the same patient. *J. Bone*
590 *Joint Surg. Am.* **94**, 2128–2135 (2012).
- 591 40. Wang, L., Wang, S. & Li, W. RSeQC: quality control of RNA-seq experiments.
592 *Bioinformatics* **28**, 2184–2185 (2012).
- 593 41. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
594 sequence data. *Bioinformatics* **30**, 2114–2120 (2014).
- 595 42. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse
596 genomes. *Nucleic Acids Res.* **47**, D766–D773 (2019).
- 597 43. Kim, D., Langmead, B. & Salzberg, S. L. HISAT: a fast spliced aligner with low
598 memory requirements. *Nat. Methods* **12**, 357–360 (2015).
- 599 44. Pertea, M. *et al.* StringTie enables improved reconstruction of a transcriptome from
600 RNA-seq reads. *Nat. Biotechnol.* **33**, 290–295 (2015).
- 601 45. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and
602 dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
- 603 46. Castel, S. E., Levy-Moonshine, A., Mohammadi, P., Banks, E. & Lappalainen, T.
604 Tools and best practices for data processing in allelic expression analysis. *Genome*
605 *Biol.* **16**, 195 (2015).
- 606 47. van de Geijn, B., McVicker, G., Gilad, Y. & Pritchard, J. K. WASP: allele-specific
607 software for robust molecular quantitative trait locus discovery. *Nat. Methods* **12**,
608 1061–1063 (2015).
- 609 48. GTEx Consortium. Human genomics. The Genotype-Tissue Expression (GTEx)
610 pilot analysis: multitissue gene regulation in humans. *Science* **348**, 648–660
611 (2015).