1

2 **Brief Communication:**

3 **Proteome-scale tagging and functional screening in mammalian cells by ORFtag**

4

5

6 Filip Nemčko[1,2,7], Moritz Himmelsbach[2,3,4,5,7], Vincent Loubiere[1], Ramesh Yelagandula[5],

7 Michaela Pagani[1], Nina Fasching[5], Julius Brennecke[5,*], Ulrich Elling[5,*], Alexander

8 Stark[1,6,*], Stefan L. Ameres[3,4,5,*]

9

10

11 [1] Research Institute of Molecular Pathology (IMP), Vienna BioCenter (VBC), 1030 Vienna,

12 Austria

13 [2] Vienna BioCenter PhD Program, Doctoral School of the University of Vienna and

14 Medical University of Vienna, 1030 Vienna, Austria.

15 [3] Max Perutz Labs, Vienna BioCenter (VBC), Dr.-Bohr-Gasse 9, 1030 Vienna, Austria

16 [4] University of Vienna, Center for Molecular Biology, Department of Biochemistry and Cell

17 Biology, Dr.-Bohr-Gasse 9, 1030 Vienna, Austria

18 [5] Institute of Molecular Biotechnology (IMBA), Vienna BioCenter (VBC), 1030 Vienna,

19 Austria

20 [6] Medical University of Vienna, Vienna BioCenter (VBC), Vienna, Austria

21 [7] These authors contributed equally: Filip Nemčko, Moritz Himmelsbach

22

23 * Co-corresponding authors: julius.brennecke@imba.oeaw.ac.at,

24 ulrich.elling@imba.oeaw.ac.at, stark@starklab.org, stefan.ameres@univie.ac.at

**Abstract**

Determining protein function in a systematic manner is a key goal of modern biology, but remains challenging with current approaches. Here, we present ORFtag, a versatile, cost-effective and highly efficient method for the massively-parallel tagging and functional interrogation of proteins at proteome scale. Using mouse embryonic stem cells, we showcase ORFtag's utility through screens for transcriptional activators, repressors and post-transcriptional regulators. Each screen finds known and novel regulators, including long ORFs not accessible to other methods, revealing that Zfp574 is a highly selective transcriptional activator and that oncogenic fusions frequently function as transactivators.

**Main text**

Proteins are pivotal in nearly all cellular processes, but their biochemical diversity often hinders systematic protein function studies. Genetic loss- or gain-of-function screens - such as CRISPR-Cas9, Cas9i and Cas9a screens - are powerful methods for identifying genes involved in specific cellular processes, but typically do not provide direct insight into protein function[1]. They are also often hampered by functional redundancies and the essentiality of many genes. Conversely, sufficiency-based assays allow direct determination of protein-inherent function[2,3]. However, current systematic methods hinge on the delivery and expression of open reading frame (ORF) libraries[4,5], which are not only costly and challenging to maintain but also tend to favor shorter ORFs (<5kb) due to limitations in DNA synthesis, cloning, viral packaging, and delivery into cells[2]. While targeted engineering of native gene locations can bypass these limitations, recent CRISPR-Cas9 techniques for systematic gene tagging have only scaled to several hundred genes[6–10].

Here, we present ORFtag, a versatile approach that allows for the massive, parallel, and proteome-scale tagging of endogenous ORFs, overcoming critical limitations of current methods. ORFtag is based on insertional elements such as retroviral vectors containing a constitutively active promoter, a selection gene, and a functional tag of interest followed by a splice donor sequence (**Fig. 1a**). Upon large-scale transduction of cultured cells,

55   ORFtag cassettes randomly integrate into the genome and drive the transcription of

56   nearby endogenous gene loci by splicing of the functional tag to splice-acceptor sites

57   downstream of the integration site. With one cassette for each of three open reading

58   frames, ORFtag can be used to generate N-terminal fusions of endogenous ORFs with a

59   wide range of functional tags. A key feature of ORFtag is its compatibility with diverse

60   functional readouts including reporter-based positive selection by fluorescence-activated

61   cell sorting (FACS). In the selected cell population, tagged genes are identified by

62   mapping integration sites using inverse PCR (iPCR) followed by next-generation

63   sequencing (NGS)[11].

64

65   To benchmark the ORFtag approach, we performed three functional screens for

66   transcriptional activators, transcriptional repressors, and post-transcriptional gene

67   regulatory (PTGR) proteins in mouse embryonic stem cells (mESCs), each in two

68   biological replicates (**Fig. 1b**). For the transcriptional activator and repressor screens, we

69   systematically fused proteins to the DNA binding domain of the bacterial Tet Repressor

70   (TetR), enabling their recruitment to TetO binding sites located upstream of an integrated

71   GFP reporter. This reporter contained an inactive minimal promoter or a constitutively

72   active promoter for transcriptional activator and repressor screens, respectively. For the

73   PTGR screen, we tagged proteins with the lambda phage N protein (λN) in order to recruit

74   them to boxB sites located in the 3´ UTR of a constitutively expressed GFP reporter

75   mRNA. To ensure that each cell expresses only one tagged ORF, we transduced reporter

76   cells with retroviruses carrying ORFtag cassettes at low multiplicity of infection (MOI)

77   followed by selection. Cells with altered GFP reporter expression – increased for activator

78   and decreased for repressor and PTGR screens – were isolated by FACS and insertion

79   sites were determined in pool by inverse PCR (iPCR) on ring-ligated short DNA fragments

80   obtained by restriction digest and subsequent NGS[11] (**Fig. 1a**). Finally, we identified gene

81   loci where insertions were statistically over-represented in the sorted samples compared

82   to the non-selected background dataset by assigning each integration to the nearest

83   downstream splice acceptor-containing exons of protein-coding genes (see **Online**

84   **Methods** and **Extended Data Table 1**). For each of the three screens, we found a

85    prominent, screen-specific enrichment of insertions at positive control genes, exemplified

86    by the transcriptional coactivator Yap1 (for the activator screen), the KRAB domain-

87    containing Zfp57 (repressor), and the mRNA deadenylase complex subunit Cnot9

88    (PTGR) (**Fig. 1c**).

89

90    Overall, we identified 139 putative transcriptional activators, 207 repressors, and 77

91    PTGR proteins using stringent thresholds (FDR < 0.1%, log2 odds ratio ≥1, see **Extended**

92    **Data Table 1**). Activator hits include several known transcriptional activators, such as

93    p65, Ep300, subunits of the Mediator complex, and all Kmt2(a-d) histone-

94    methyltransferases, which could not be screened before due to their long ORFs of up to

95    17kb (**Fig. 1d, Extended Data Fig. 1a**). Repressor hits contain 75 KRAB zinc-finger

96    repressors and their corepressor Trim28, HP1 family proteins, H3K9 methyl-transferases,

97    and Polycomb repressive complex components (**Fig. 1d, Extended Data Fig. 1a**).

98    Finally, the PTGR screen identified core components of the microRNA (Ago2, Tnrc6a/b/c)

99    and nonsense-mediated decay (Smg1, Smg9, Upf2) pathways, members of the Ccr4-Not

100    deadenylation complex (Cnot2, Cnot3, Cnot9), as well as translational inhibitors (Eif4e2,

101    Eif4enif1) (**Fig. 1d, Extended Data Fig. 1a**).

102

103    While ORFtag integrations per screen were highly reproducible (**Extended Data Fig. 1b**),

104    the hits from the three distinct assays showed almost no overlap, indicating that ORFtag

105    does not lead to the recurrent and artefactual detection of unspecific genes (**Fig. 2a &**

106    **Extended Data Fig. 1c**). Consistent with this, the activator and repressor screen hits

107    were strongly enriched for proteins containing activating and repressive domains,

108    respectively, and both protein sets share a significant enrichment for known transcription

109    factors (TFs) (**Fig. 2b**). Similarly, only PTGR hits were enriched for known RNA-binding

110    proteins (**Fig. 2b**). Moreover, the genes identified by the three screens were enriched for

111    distinct gene ontology (GO) terms, protein domains and subcellular locations, all of which

112    are consistent with their associated functions (**Fig. 2c**). Overall, these results indicate that

113    ORFtag is compatible with diverse functional assays and delivers assay-specific hits.

114

115    To experimentally validate the screen results at the level of protein-inherent functionality,

116    we individually cloned and transduced eight hits from each screen, fused to the respective

117    TetR or λN tags, and tested whether they were sufficient to regulate the respective

118    reporters. All candidates tested, including hits that were not previously linked to the

119    respective biological processes, could be validated in recruitment assays together with

120    previously known regulators, confirming that ORFtag screens are highly specific and have

121    low false positive rates (**Fig. 2d**). For example, recruitment of the annotated cytoskeletal

122    protein Pxn or the uncharacterized protein 1700102P08Rik was sufficient to strongly

123    activate transcription, while the E3 ubiquitin ligase Trim8 and the uncharacterized protein

124    Msantd3 were sufficient to repress transcription. In turn, the neuronal activity-associated

125    protein Maco1 and the E3 ubiquitin ligase Trim13 were sufficient to repress reporter gene

126    expression when recruited to the 3'UTR of an mRNA.

127

128    To assess the potential of ORFtag in assigning cellular roles to uncharacterized proteins,

129    we sought to investigate the endogenous function of the zinc-finger protein Zfp574, which

130    ORFtag specifically identified as a transcriptional activator. Utilizing the auxin-inducible

131    degron system (**Fig. 2e**), we show that depletion of Zfp574 results in a significant growth

132    defect (**Fig. 2f**), indicating that Zfp574 is essential for cellular fitness. Rapid depletion of

133    Zfp574 followed by PRO-seq further revealed that Zfp574 functions strictly as a

134    transcriptional activator in accordance with the ORFtag results (39 genes go down, 0

135    genes go up upon depleting Zfp574 at FDR≤0.05 and FC≥2) (**Fig. 2g**). Cut&Run for

136    Zfp574 identified 140 binding sites genome-wide, the majority (87.9%) of which are

137    located in promoter-proximal positions (+- 500 bp around the gene transcription start

138    sites), and transcription of the promoter-bound genes was strongly affected upon Zfp574

139    depletion (**Fig. 2h, i**). Thus, Zfp574 is a novel selective transcriptional activator that

140    specifically binds and activates a small set of genes that support cell fitness. Taken

141    together, our results demonstrate that ORFtag, coupled with functional assays, provides

142    a robust and powerful method for the high-throughput assignment of protein function.

143

144    Some identified hits may regulate gene expression in ORFtag assays without necessarily

145    doing so endogenously. This underscores the distinction between a protein's inherent

146    biochemical function (as evaluated here) and its role within the cell. In fact, the process

147    of tagging and/or chromatin- or RNA-tethering can potentially overwrite a protein's usual

148    cellular function and change its localization within the cell (e.g. signaling peptides can be

149    bypassed, replaced or overwritten by ORFtag). Importantly, these hits remain valuable as

150    their ability to activate/repress gene expression *in principle* is highly relevant, e.g. in

151    cancer, when chromosomal re-arrangements create oncogenic fusion proteins. Indeed,

152    among our hits is the ortholog of the oncogene C3orf62, recently described by a tethering-

153    based approach to be an activator[2]. We compared the ORFtag hits systematically to their

154    human orthologs and found oncogenes to be enriched among the activators and – more

155    weakly – the repressors but not the post-transcriptional regulators (**Fig. 2j**). These include

156    for example Zc3h7b and D630045J12Rik (KIAA1549in human) that can function as

157    activators, and Gm10324 that can function as a repressor, highlighting that oncogenic

158    fusions can recruit unrelated genes to function in gene regulation.

159

160    Having established that retroviral integration sites represent indeed successful ORF

161    tagging events that score in functional assays, we conducted a systematic and critical

162    assessment of ORFtag's ability to comprehensively and reproducibly tag proteins.

163    Comparison of the retroviral integration sites from six independent transductions,

164    performed in three different laboratories, revealed that, regardless of the protein tag used,

165    experimental cassettes integrated in a similar distribution across the genome and the

166    number of insertions per genomic region correlated highly (PCC≥0.84) (**Extended Data**

167    **Fig. 1b**). ORFtag integrations were enriched near transcription start sites (TSS), a well-

168    known feature of retroviral vectors[11], enabling the tagging of near-full-length proteins (**Fig.**

169    **3a**). Assigning each integration to a gene locus indicated that we were able to tag at least

170    83.7% of all mouse protein-coding genes with a median count of 15 integrations per gene,

171    given the scale and sequencing depth of our screens (**Fig. 3b, 3c**). The tagged genes

172    include those with large open reading frames yielding high molecular weight proteins (**Fig.**

173    **3d**). Indeed, in contrast to ORFeome-based approaches that are biased towards short

174  ORFs, ORFtag is not influenced by gene length (**Fig. 3e)**. Moreover, the retroviral ORFtag

175  cassette allowed the tagging of ORFs that exhibited different endogenous expression

176  levels, including >59% of genes that are normally not expressed in mESCs (**Fig. 3f**).

177  Importantly, also the hits identified in the three functional screens include genes of varying

178  lengths and expression levels (**Fig. 3e & 3f**).

179

180  A limitation of ORFtag lies in its inability to functionally probe intronless genes and first

181  exons, due to them lacking splice acceptor sites. However, it is worth noting that 45.6%

182  of first exons are non-coding and that among protein-coding first exons, the median

183  encoded peptide length is 31 amino acids short. As a result, only 12.8% of first exons

184  contain annotated protein domains (**Extended Data Fig. 1d**). Intronless genes, which

185  cannot be tagged, constitute only a small fraction of protein coding genes (5.9%). These

186  are dominated by a few protein families, including histones and various sensory receptors

187  (**Extended Data Fig. 1e**), leaving 94.1% of protein-coding genes as potentially taggable

188  by ORFtag. We also note that certain genes may not be accessible to ORFtag screens if

189  cellular fitness is sensitive to changes in their expression levels.

190

191  In conclusion, ORFtag is an easy-to-implement functional genomics tool that enables

192  cost-effective proteome-scale functional screens. Due to its modularity, ORFtag can be

193  combined with various functional tags and a wide range of applications, including bio-

194  imaging, targeting proteins to various organelles, and protein-protein interaction studies.

195  Notably, while ORFtag utilizes N-terminal tagging, it can be adapted for C-terminal or

196  internal tagging, broadening the scope of proteins and applications that can be explored.

197  These alternative tagging approaches would mirror the tagged genes' endogenous

198  expression levels and are thus limited to genes expressed in the particular cell line.

199  Finally, ORFtag can be readily employed in cellular systems of various model organisms

200  without the need to generate species-specific resources. This adaptability and versatility

201  make ORFtag a promising tool for advancing functional genomics research.

202

**Figure Legends**

**Fig.1: ORFtag is a versatile tool for proteome-wide functional assays**

**a,** Overview of the ORFtag approach. The ORFtag cassette is embedded in a retroviral backbone and contains a constitutively active promoter, selection gene, a tag, and a splice-donor site. Upon transduction, the ORFtag cassette randomly integrates into the genome and prompts splicing to a downstream splice acceptor, ultimately producing a tagged protein. **B,** Schematic view of 3 different screens for transcriptional activators (green), repressors (blue) or PTGRs (yellow). **C,** Genome browser screenshots of ORFtag integration sites (vertical lines) in positive (+, top) or negative (-, bottom) strand direction, before (non-selected, grey) and after FACS selection at the genomic locus of each one activator (Yap1, green), repressor (Zfp57, blue) and PTGR (Cnot9, yellow) hit emerging from ORFtag screens in mESCs. Log2 odds-ratio (log2OR) and false discovery rate (FDR) are indicated. **D,** Volcano plots highlighting known (black circles, names) and validated (marked red, see Fig. 2d) hits for the 3 screens.

**Fig.2: ORFtag interrogates protein function with high specificity**

**a,** Overlap between activator (green), repressor (blue) and PTGR (yellow) hits. **b,** Enrichment of screen hits for human homologous genes with annotated DNA-binding, activation, or repressive domains and RNA-binding proteins. **c,** Top enriched protein domains, biological process and cellular compartment GO terms for activator, repressor and PTGR hits. **d,** Independent validation of select screen hits. GFP intensity measured by flow cytometry in reporter cell lines stably expressing the indicated full-length proteins fused to TetR (Activator, Repressor) or $\lambda$N (PTGR); Wilcoxon test, ***p≤1×10$^{-3}$. Refer to Fig. 1d for the position of the hits in the volcano plot. **e**, Schematic view of Zfp574 rapid depletion using Auxin-Inducible Degron (AID). The rapid depletion of Zfp574 upon 3-indoleacetic acid (IAA) treatment shown by Western Blot. **f**, Cell viability timecourse in the presence (-IAA, in grey) or absence of Zfp574 (+IAA, in red). Shown are two biological replicates. **g**, MA plot showing PRO-seq fold-changes (log2) upon Zfp574 6h depletion. Signficantly up- (0) or down-regulated (39) genes are highlighted in red. **h**, PRO-seq fold-changes (log2) of not-bound versus Zfp574 promoter-bound genes upon Zfp574

233    depletion; Wilcoxon test, ****p≤1×10⁻⁵. **i,** Zfp574 Cut&Run and PRO-seq screenshots at

234    the Rpl10 locus. **j,** Enrichment of screen hits for genes that were identified as part of

235    oncogenic fusions.

236

237    **Fig.3: Scope and limitations of massive parallel protein tagging using ORFtag**

238    **a,** Distribution of ORFtag integrations around TSSs of mouse protein-coding genes. **b,**

239    Saturation curve displaying the relationship between the fraction of tagged proteins and

240    the number of determined integration sites. **c,** Fraction of genes showing at least one

241    integration in the combined background sample. **d,** Western blot against the FLAG tag

242    assessing the tagging pattern in mESC lysate before (-) and after (+) ORFtag

243    transduction. **e,** Ratio of protein coding genes that were successfully tagged using

244    ORFtag (ORFtag, pink) or were hits in any of the three screens (ORFtag hits, purple),

245    normalized by the distribution found across the whole mouse genome (Genome, dashed

246    line). Human ORFeome is shown for comparison (ORFeome, light grey). See material

247    and methods for further details. **f,** Ratio of intron-containing protein coding genes that

248    were successfully tagged using ORFtag (ORFtag, pink) or were hits in any of the three

249    screens (ORFtag hits, purple), normalized by the distribution found across the whole

250    mouse genome (Genome, dashed line).

251

252    **Extended Data Legends**

253    **Extended Data Fig.1: Evaluation of ORFtag integrations at global scale**

254    **a,** STRING protein-protein interaction networks between activator/ repressor/ PTGR hits.

255    Node communities were highlighted using a color code (Louvain method), and their size

256    is proportional to the Odd ratio of the corresponding hit. Only the hits showing at least

257    one interaction with another hit are shown. **b,** Dendrogram of Pearson's Correlation

258    Coefficients between unsorted (triangles) and sorted (round) replicates from each

259    functional screen. All background (input) samples show high PCC (≥ 0.85). **c,** Scatter

260    plots displaying a pairwise comparison of -log10(FDR) values for screened genes across

261    different assays. **d,** Protein family enrichment of intronless genes, whose majority belongs

262    to few protein families. **e,** Fraction of exons that are non-coding (in white), code either for

263 no known protein domain (in shades of pink) or for a high confidence PFAM protein

264 domain (Domain-containing CDS, in blue). Importantly, most first exons are non-coding

265 or do not contain specific protein domains, fostering the use of ORFtag for a wide range

266 of functional studies.

267

268 **Extended Data Table 1: Identification of activator, repressor and PTGR hits**

269 For each gene locus, raw counts, odd ratio (log2) and the associated FDR are shown for

270 the 3 different screens. The last "hit" column specifies whether a locus was considered

271 as a hit (TRUE) or not (FALSE).

272

273 **References**

274 1. Nemčko, F. & Stark, A. Proteome-scale identification of transcriptional activators in
275   human cells. *Mol Cell* **82**, 497–499 (2022).
276 2. Alerasool, N., Leng, H., Lin, Z. Y., Gingras, A. C. & Taipale, M. Identification and
277   functional characterization of transcriptional activators in human cells. *Mol Cell* **82**, 677-
278   695.e7 (2022).
279 3. Luo, E. C. *et al.* Large-scale tethered function assays identify factors that regulate mRNA
280   stability and translation. *Nat Struct Mol Biol* **27**, 989–1000 (2020).
281 4. Wiemann, S. *et al.* The ORFeome Collaboration: A genome-scale human ORF-clone
282   resource. *Nat Methods* **13**, 191–192 (2016).
283 5. Yang, X. *et al.* A public genome-scale lentiviral expression library of human ORFs. *Nat*
284   *Methods* **8**, 659–661 (2011).
285 6. Reicher, A., Koren, A. & Kubicek, S. Pooled protein tagging, cellular imaging, and in situ
286   sequencing for monitoring drug action in real time. *Genome Res* **30**, 1846–1855 (2020).
287 7. Serebrenik, Y. V., Sansbury, S. E., Kumar, S. S., Henao-Mejia, J. & Shalem, O. Efficient
288   and flexible tagging of endogenous genes by homology-independent intron targeting.
289   *Genome Res* **29**, 1322–1328 (2019).
290 8. Schmid-Burgk, J. L., Höning, K., Ebert, T. S. & Hornung, V. CRISPaint allows modular
291   base-specific gene tagging using a ligase-4-dependent mechanism. *Nat Commun* **7**,
292   (2016).
293 9. Yarnall, M. T. N. *et al.* Drag-and-drop genome insertion of large sequences without
294   double-strand DNA cleavage using CRISPR-directed integrases. *Nat Biotechnol* **41**, 500–
295   512 (2023).
296 10. Sansbury, S. E., Serebrenik, Y. V, Lapidot, T., Burslem, G. M. & Shalem, O. Pooled
297    tagging and hydrophobic targeting of endogenous proteins for unbiased mapping of
298    unfolded protein responses. *bioRxiv* 2023.07.13.548611 (2023)
299    doi:10.1101/2023.07.13.548611.
300 11. Elling, U. *et al.* A reversible haploid mouse embryonic stem cell biobank resource for
301    functional genomics. *Nature* **550**, 114–118 (2017).
302

303

## Methods

### Cell culture conditions

All experiments presented here were carried out in diploid mouse embryonic stem cells (mESCs) that were derived from originally haploid HMSc2 termed AN3-12[11]. The mESCs were cultivated without feeders in high-glucose-DMEM (Sigma-Aldrich) supplemented with 13.5% fetal bovine serum (Sigma-Aldrich), 2 mM L-glutamine (Sigma-Aldrich), 1x Penicillin-Streptomycin (Sigma-Aldrich), 1x MEM non-essential amino acid solution (Gibco), 1mM sodium pyruvate (Sigma-Aldrich), 50 mM β-mercaptoethanol (Merck) and in-house produced recombinant LIF. Virus packaging cell lines, Lenti-X 293T (Takara), and PlatinumE (Cell Biolabs), were grown according to the manufacturer's instructions. All cell lines were cultured at 37°C and 5% CO2 and regularly tested for mycoplasma contamination.

### Reporter cell lines

The reporter cell line for the Repressor screen was established previously[12] and contains the reporter construct inserted into the expression-stable locus on Chr15 that is compatible with the Flp recombinase-mediated cassette exchange (RMCE). Reporter cell line for the Activator screen was generated by RMCE as follows – $5x10^6$ cells were electroporated with a mix of 10 $\mu$g of plasmid containing constructs flanked by FRT/F3 sites, and 6 $\mu$g of plasmid expressing Flp, using a Maxcyte STX electroporation device (GOC-1) and the Opt5 program. Seven days after the transfection, cells were sorted and clonal cell lines were generated. Cell lines were genotyped using integration-site specific PCRs and Sanger sequencing. The Activator reporter construct (Addgene, this study) contains the PuroR-IRES-GFP reporter under the control of the minimal promoter derived from the MYLPF gene (chr16:30374730–30374857 +, hg38)  that was shown to have a low basal expression and high inducibility[13]. Upstream of the promoter are 7x TetO sites flanked by the loxP sites.

331

332  The reporter cell line for the PTGR screen was created by nucleofection of haploid AN3-
333  12 mESCs with 500ng of the reporter construct and 10 $\mu$g of a Tol2 transposase encoding
334  plasmid using the Mouse ES Cell Nucleofector Kit (Lonza) according to the
335  manufacturer's protocol using an Amaxa Nucleofector (Lonza). The PTGR reporter
336  construct (Addgene, this study) encodes for PuroR-IRES-GFP followed by 10 boxB sites
337  that are flanked by two loxP sites under the control of a PGK promoter. Cells were
338  subsequently selected using 1 $\mu$g/ml Puromycin (Gibco) followed by single clone
339  selection. Single cell clones were afterwards transduced with a retroviral vector for the
340  expression of pMSCV_hygro_CreERT2 and selected with 250 $\mu$g/ml Hygromycin (Roche)
341  followed by single cell clone selection.

342

343  ORFtag screens
344  The ORFtag viral constructs were derived from the ecotropic Retro-EGT construct[11] that
345  includes the sequence features necessary for the inverse-PCR protocol (see below).
346  Furthermore, the construct contains constitutively active PGK promoter that drives the
347  expression of a NeoR resistance gene separated from a tag by the IRES sequence. The
348  tag contained either TetR with an N-terminally located nuclear localization signal
349  (Activator and Repressor screens; this study) or LambdaN domain (PTGR screen; this
350  study). Additionally, the tag contained 2x GGGS-linker followed by the BC2-tag and
351  3xFLAG-tag. Finally, the ORFtag construct contains a consensus splice donor motif
352  followed by a part of the Hprt intron (chrX:53020400-53020556 +, mm10). In order to tag
353  genes in all three possible coding frames, three constructs were used that contain either
354  0, 1 or 2 additional nucleotides upstream of the consensus splice motif.

355

356  Every ORFtag screen was performed in two independent replicates. Retroviral constructs
357  carrying ORFtag cassette were packed in PlatinumE cell lines using polyethylenimine
358  (PEI) reagent as described previously[11]. Reporter cell lines (100-150 million cells) were
359  transduced with packaged retrovirus in the presence of 6 $\mu$g/ml polybrene (Sigma) and at
360  low transduction efficiency (< 20%) to ensure only one virus per cells. Cells were
361  harvested 24 hours later and plated in medium containing 0.1 mg/ml G418 (Gibco) for

362 selection of transduced cells. Selection was continued until all cells on the control plate

363 died (4-5 days), after which 40 million cells were processed as non-selected background

364 for mapping of genomic integrations (see below). The remaining cells were sorted for

365 GFP-positive (Activator screens) or GFP-negative (Repressor screens) populations using

366 BD FACSAria III or IIu cell sorters (BD Biosciences) and processed for mapping of

367 genomic integrations (see below). For the PTGR screen a five-sort strategy was applied

368 to enrich cells that show a tethering dependent repression of reporter gene expression.

369 Cells with a GFP expression equal to the lowest 10 percent of GFP expression observed

370 after selection were sorted using BD FACSAria III and expanded thereafter. Additionally,

371 non-sorted cells were maintained for gating of the consecutive sorts. Two additional sorts

372 for cells with GFP expression similar to the lowest 10% of GFP signal observed in the

373 non-sorted cells were performed and again expanded in-between the sorts. A fourth sort

374 was performed for cells with a GFP expression equal to the lowest 5% of GFP signal

375 observed in the non-sorted cells. After expansion, the cells were treated with 500nM 4-

376 Hydroxytamoxifen (Sigma) to induce Cre-mediated recombination and to flox the boxB

377 sites of the reporter construct and hence to revert the tethering. Thereafter a final sort

378 was performed to select a cell population with a GFP expression equal to the highest 70%

379 of GFP expressing cells.

380

381 Transduction efficiency was measured by plating 10,000 cells on a 15-cm dish and

382 selecting with G418 (Gibco). A control plate with 1,000 cells was also plated without

383 selection. After 10 days, colonies were counted and transduction efficiency was

384 calculated as the number of colonies on the selected plate divided by the total number of

385 cells plated (10 times the number of colonies on the control plate).

386

387 <u>Mapping of genomic integrations by next-generation sequencing</u>

388 Genomic locations of ORFtag integrations were mapped using modified inverse-PCR

389 followed by next generation sequencing (iPCR-NGS) protocol[11]. Genomic DNA was

390 prepared by lysing cell pellets in lysis buffer (10 mM Tris-HCl pH 8.0, 5 mM EDTA, 100

391 mM NaCl, 1% SDS, 0.5 mg/ml proteinase K) at 55°C overnight. Following a 2-hour RNase

392　A treatment (Qiagen, 100 mg/ml, 1:1,000 dilution) at 37°C, two extractions using

393　phenol:chloroform:isoamyl alcohol and one extraction using chloroform:isoamyl alcohol

394　were carried out. The samples then underwent two separate digestion reactions (with up

395　to 4 $\mu$g of genomic DNA) using NlaIII and MseI enzymes (NEB) at 37°C overnight,

396　followed by purification using a Monarch PCR&DNA Cleanup Kit (NEB). Ring-ligation was

397　carried out using T4 DNA ligase (NEB) at 16 °C overnight, followed by heat-inactivation

398　(65°C, 15 min) and linearization using SbfI-HF (NEB) at 37°C for 2 h. The digests were

399　then purified using a Monarch PCR&DNA Cleanup Kit (NEB) and amplified using firstly a

400　nested PCR reaction with KAPA HiFi HotStart ReadyMix (Roche), and a specific primer

401　pair　　　　　　　　　　　　　　　　　　　　　　　(TGCAGGACCGGACGTGACTGGAGTTC*A,

402　TGCAGGACGATGAGCAGAGCCAGAACC*A) for 16 cycles. After cleanup with AMPure

403　XP Reagent (Beckman Coulter, 1:1 ratio beads:PCR), iPCR amplification was carried out

404　with　KAPA　HiFi　HotStart　ReadyMix　(Roche),　and　a　specific　primer　pair

405　(AATGATACGGCGACCACCGAGATCTACACGAGCCAGAACCAGAAGGAACTTGA*C,

406　CAAGCAGAAGACGGCATACGAGAT　　　　　　　　　　　　　　　　　　[custom-barcode]

407　GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT)　for　18　cycles.　Afterwards,

408　amplified libraries were size selected for a range of 400-800 bp using SPRIselect beads

409　(Beckman Coulter). NGS was performed on an Illumina NextSeq550 or Ilumina HiSeq

410　2500 sequencer according to the manufacturers' protocols with custom first-read primer

411　(1:1　　mix　　of　　GAGTGATTGACTACCCGTCAGCGGGGGTCTTTCA　　and

412　TGAGTGATTGACTACCCACGACGGGGGTCTTTCA).

413

414　<u>Immunoprecipitation</u>

415　To confirm expression of tagged proteins, the PTGR reporter mESCs, transduced with

416　the ORFtag construct as well as non-transduced cells, were lysed in lysis buffer (50mM

417　TRIS HCl pH: 7.5, 150mM NaCl, 0.1% SDS, 1% Triton-X-100, 0.5% NP-40, 0.5mM EDTA

418　supplemented　with　Proteinase　Inhibitor　(Roche)　and　protein　concentration　was

419　determined photometrically using the Protein Assay Dye Reagent Concentrate (BioRad),

420　according　to　the　manufacturer's　protocol　and　photometric　measurement　at　595nm.

421　Tagged proteins were captured using 80 $\mu$l of in-house BC2-nanobody coupled magnetic

422    beads from 1mg total protein. Bound proteins were eluted by resuspension of the beads

423    in 1x SDS-sample buffer and incubated at 95°C for 5 minutes. Further details about

424    Western blotting can be found below.

425

426    <u>Individual recruitment validations</u>

427    To validate Activator hits, the candidates were amplified by PCR from mESC cDNA and

428    inserted into retroviral constructs that comprises the PGK promoter that drives the

429    expression of a PuroR resistance gene and a tag separated by the IRES sequence. The

430    tag contains TetR, along with an N-terminal nuclear localization signal, a 2x GGGS-linker,

431    a BC2-tag, and a 3xFLAG-tag, followed by the tested candidate. Retroviral constructs

432    were packed in PlatinumE cell lines (see above), and reporter cell lines (170,000 cells)

433    were transduced in the presence of 6 $\mu$g/ml polybrene (Sigma). Cells were harvested 24

434    hours later and plated in medium containing 1 $\mu$g/ml Puromycin (InvivoGen) to select for

435    transduced cells. After five days of selection, the reporter expression was analyzed on an

436    LSR Fortessa (BD) flow cytometer. For processing and visualization, FlowJo and R

437    package flowCore (v2.12.2) was used.

438

439    In order to validate Repressor and PTGR hits, PCR was used to amplify the candidates

440    from mESC cDNA, and lentiviral plasmids were created as fusion proteins containing

441    TetR/lamdaN-Candidate-P2A-mCherry coding sequence under the control of an EF1a

442    promoter. For the validation of Trim71, cDNA excluding the fragment encoded in exon 1

443    (Trim71dE1) was cloned into the aforementioned lentiviral plasmid. A fragment encoding

444    for the silencing domain of human Tnrc6b (Tnrc6b-SD) was expressed using the same

445    lentiviral plasmid as above as a positive control for the validation of PTGR hits. Lentivirus

446    was produced in Lenti-X 293T cells as in (Ref.[12]). Reporter cells were then transduced

447    with the virus in the presence of 8$\mu$g/ml polybrene (Santa Cruz Biotechnology, SACSC-

448    134220). After 7 days of transduction, reporter expression was analyzed on an LSR

449    Fortessa (BD) flow cytometer. Reporter cells transduced with recruitment constructs were

450    gated based mCherry expression. For processing and visualization, FlowJo and R

451    package flowCore (v2.12.2) was used.

452

453     <u>AID cell line generation</u>

454     A parental cell line expressing the E3 ligase for AID was generated by inserting a cassette

455     into the expression-stable locus on Chr15 that is compatible with the Flp recombinase-

456     mediated cassette exchange in mESCs (RMCE, see "Reporter cell lines" section). The

457     construct contained EF1alpha- ARF16- HA- P2A- OsTir1- 3xMyc- T2A- mCherry-

458     SV40_polA site flanked by the FRT/F3 sites. The clonal Tir1 parental cell line was

459     genotyped using integration-site specific PCRs and Sanger sequencing.

460

461     To generate the N-terminally AID-tagged Zfp574 cell line, $5x10^6$ Tir1 parental cells were

462     transfected with 10 $\mu$g of plasmid (Ref.[14]) that expresses Cas9 and the gRNA against a

463     target locus (CTTGCTGCTGCCATGACTG) and 5 $\mu$g of plasmid with a knock-in cassette

464     containing Blasticidin-P2A-V5-AID-GGGS flanked by 20 bp microhomology arms (Ref.[14])

465     using a Maxcyte STX electroporation device (GOC-1) and the Opt5 program. Two days

466     after the transfection, cells were selected for knock-ins with 10 $\mu$g/mL Blasticidin

467     (ThermoFisher), individual clones were genotyped using knock-in-site specific PCRs and

468     Sanger sequencing. Potential candidates were investigated by western blotting against

469     the integrated V5-tag (Thermo Fisher, R960-25) with or without 500 uM 3-indoleacetic

470     acid (IAA, Merc) treatment.

471

472     <u>Western blotting</u>

473     Cells ($3x10^6$) were collected, centrifuged at 300g for 5 min, washed with 1xPBS and lysed

474     in 100 $\mu$l RIPA buffer containing protease inhibitor (Roche) and Benzonase (Sigma

475     Aldrich). For complete lysis, cells were incubated on ice for 30 min and sonicated for five

476     minutes (30 sec on/off, Diagenode Bioruptor). Afterwards, samples were centrifugated for

477     5 min at full speed and 4 °C, and supernatants were supplemented with 20 $\mu$l 4x Laemmli

478     buffer with 10% β-mercaptoethanol. Samples were boiled for 5 min at 98°C.

479

480     Proteins were resolved on SDS-PAGE on a 4-15% Mini-PROTEAN TGX gel (BioRad)

481     and transferred to an Immobilon-P PVDF membrane (Merck Millipore) using a wet-

482   chamber system. Tagged proteins were detected using mouse α-Flag M2 (Sigma Aldrich

483   F3165, 1:10,000), mouse α-V5-tag (Thermo Fisher R960-25, 1:1,000), or rabbit α-β-

484   tubulin (Addgene ab6046, 1:10,000) as primary and HRP-α-Mouse (Cell Signaling, 7076,

485   1:10,000) or HRP-α-Rabbit (Cell Signaling, 7074, 1:10,000) as secondary antibody and

486   imaged using ClarityTM Western ECL Substrate (BioRad) with a ChemiDocTM Imaging

487   System (BioRad) using ImageLab v5.1.1 (BioRad).

488

489   Cell viability timecourse

490   For growth curve assays, AID-tagged cell line (mCherry positive, see AID cell line

491   generation) was mixed at 1:1 ratio with WT cells, split into control (-IAA) and treatment

492   (+IAA, Merc, 500 uM) group and cultured in a 24-well cell culture plate. The ratio between

493   mCherry positive and negative cell was quantified every 24hrs by Flow Cytometry (iQue

494   Screener PLUS, Intellicyt).

495

496   PRO-seq

497   For each condition, $1x10^7$ AID-Zfp574 cells were collected and nuclei were isolated after

498   6h of 500uM IAA treatment or no treatment (two biological replicates per condition).

499   Spike-in control (S2 *Drosophila* cells; 1% of mESC cells) were added at the level of nuclei

500   permeabilization step. The next steps of the PRO-seq protocol were performed as in

501   (Ref.[15]) with a single modification: the nuclear run-on was performed at 37 °C for 3 min.

502

503   Cut&Run

504   For each biological replicate, $1x10^6$ cells from the AID-Zfp574 cell line or the Tir1 parental

505   cell line were used. The Tir1 parental cell line is used as Input, each experiment was

506   performed in two biological replicates. The protocol was performed as in (Ref.[16]) with a

507   V5-tag antibody (Thermo Fisher, R960-25) that was added to a final dilution of 1:100.

508

509   Bioinformatic analyses

510   All bioinformatic analyses were performed in R (v4.2.0, https://www.R-project.org/).

511   Computations on genomic coordinate were conducted using the GenomicRanges

512  (v1.50.1)[17] and the data.table (https://CRAN.R-project.org/package=data.table) R

513  packages. All box plots depict the median (line), upper and lower quartiles (box) ±1.5x

514  interquartile range (whiskers); outliers not shown.

515

516  *Processing of ORFtrap screens*

517  First, iPCR reads from sorted and background (non-selected) samples were trimmed

518  using Trim galore (v0.6.0) with default parameters to remove Illumina adapters. Then,

519  trimmed reads were aligned to the mm10 version of the mouse genome using bowtie2[18]

520  with default parameters (for paired-end sequenced samples, only first mate reads were

521  considered), before removal of duplicated and low mapping quality reads (mapq<=30)

522  using samtools (v1.9)[19]. Mapped insertions were assigned to the closest downstream

523  exon junction – with a maximum distance of 200kb – based on GENCODE annotations

524  of the mouse genome (vM25). Finally, insertion counts were aggregated per gene. Of

525  note, only exons from protein-coding transcripts were considered, except for the first exon

526  of each transcript, which might not contain splicing acceptor sites. Consequently,

527  intronless genes – for which none of the isoforms contain a spliced intron – were not

528  considered.

529

530  Background replicates showed reproducible gene counts (PCC ≥0.84) and therefore were

531  merged, and genes with at least one insertion were considered as putatively tagged.

532  Finally, genes showing significantly more insertions in sorted samples compared to

533  merged background samples were identified using one-tailed fisher's exact test

534  (alternative= "greater") on merged biological replicates. Of note, only genes with at least

535  3 unique insertions in sorted samples were considered. Obtained p-values were corrected

536  for multiple testing using the FDR method and genes showing an FDR<0.001 and a log2

537  Odd Ratio≥1 were classified as hits.

538

539  *Protein-protein interaction networks*

540  For each functional assay, STRING protein-protein interaction between hits were

541  retrieved using the STRINGdb R package (v2.10.0, database version 11.0). Finally, only

542 the hits showing at least one protein-protein interaction with another hit with a combined

543 score ≥900 were considered.

544

545 *CDS length bias*

546 To assess whether ORFtag is biased towards short ORFs, we stratified intronic protein

547 coding genes based on their shortest CDS length (< 2.5kb, 2.5-5kb and longer than 5kb).

548 Then, we compared how tagged genes (with at least one insertion in background

549 samples) and hits (union from the three screens) were distributed between these groups,

550 using all intronic protein coding genes as a reference. For example, to compute the

551 normalized ratio of tagged genes for the <2.5kb group, we used the following formula:

552 normalized ratio= ([tagged genes with CDS<2.5kb]/[total tagged genes])/([intronic protein

553 genes with CDS<2.5kb]/[total intronic protein coding genes]). To allow side-by-side

554 comparison, we also considered ORFs from the human ORFeome that Alerasool and

555 colleagues were able to transfect and detect[2].

556

557 *Gene expression bias*

558 To assess whether transcriptionally inactive mouse genes could be assayed using

559 ORFtag, we used publicly available data from the same mESC cell line (GSE99971)[20].

560 For each intronic protein coding gene, mean TPM was computed across three RNA-Seq

561 replicates (only protein-coding genes were considered). Genes with a mean TPM of 0

562 were classified as inactive and active genes were further stratified into quartiles. Then,

563 we compared how tagged genes (with at least one insertion in background samples) and

564 hits (union from the three screens) were distributed between these groups, using all

565 intronic protein coding genes as a reference. For example, to compute the normalized

566 ratio of tagged genes for the inactive group, we used the following formula: normalized

567 ratio= ([tagged genes with TPM=0]/[total tagged genes])/([intronic protein genes with

568 TPM=0]/[total intronic protein coding genes]).

569

*Enrichment analysis of expected protein functions*

To assess whether hits were enriched for genes with expected functions, we collected publicaly available lists of human TF genes (Ref.[21]), human genes containing activation or repressive-domains (Ref.[22], only genes containing sufficient ('S' or 'N and S') and high confidence ('H') domains were considered), human genes containing RNA-binding domains (RBPbase[23], only the genes identified in at least two different cell lines were used) and human fusion oncoproteins (COSMIC database v97, Ref.[24]). Screen hits were first assigned to their human orthologs using MGI[25] homology data. For each functional assay, we assessed whether relevant categories were enriched among the hits using one-tailed fisher's exact test (alternative= "greater"), and the total number of intronic protein coding genes as background.

*GO terms and protein domains enrichment*

Biological Process (BP), Molecular Process (MF) and Cellular Compartment (CC) Gene Ontology (GO) terms were obtained from the org.Mm.eg.db (v3.15.0) R package. Protein domains were retrieved from the EnsDb.Mmusculus.v79 R package (v2.99.0). For each functional assay, GO terms and protein domains that were over-represented among hits were identified using one-tailed fisher's exact test (alternative= "greater"), using all intronic protein coding genes as background. Obtained p-values were corrected for multiple testing using the FDR method and features with an FDR<0.05 were considered as significantly enriched. Of note, small categories containing less than five genes in total and categories with less than three matching hits were not considered. Finally, top 8-10 enriched GO terms and proteins domains were plotted for each functional assay.

*Protein family enrichment*

To identify protein families enriched among intronless genes (for which none of the isoforms contain a spliced intron), annotations were retrieved from the EnsDb.Mmusculus.v79 R package (v2.99.0). Enriched protein families were identified using one-tailed fisher's exact test (alternative= "greater"), and the total number of protein

599   coding genes as background. Obtained p-values were corrected for multiple testing using

600   the FDR method, and the protein families with an FDR<0.05 were plotted.

601

602   *Analysis of first exons*

603   For the analysis of first exons, first exons containing a predicted CDS were classified as

604   either short (≤20aa) or long (>20aa). Then, manually-curated Pfam-A domains from

605   UCSC[26] were used to discriminate first exon CDSs containing a know protein domain

606   (e.g. coding for at least 10% of a full Pfam domain) or not.

607

608   *Gene annotation for PRO-seq analysis*

609   To obtain a non-redundant set of genes for quantification of PRO-seq signals, we

610   collected all coding and long non-coding transcripts from Ensembl v.100 for the mm10

611   version of the mouse genome, excluding transcripts shorter than 300 bp. When several

612   transcript isoforms shared the same annotated transcription start site (TSS), only the

613   longest isoform was retained. Next, TSS positions were corrected using FANTOM5[27]

614   CAGE TSS clusters: for each unique annotated TSS, we identified the strongest CAGE

615   TSS within a 1kb window centered on the annotated TSS, excluding the coding sequence.

616   Finally, for each CAGE TSS, only the full length of the nearest transcript was used to

617   count overlapping reads (see next section).

618

619   *PRO-seq analysis*

620   PRO-seq libraries were sequenced in  paired-end mode with 36-bp read length. To

621   eliminate PCR duplicates, an 8-bp long unique molecular identifier (UMI) was

622   incorporated at the 5′ end of the reads during the sample processing. Before mapping,

623   the UMI was separated, and the Illumina adapters were trimmed using cutadapt v.1.18.

624   Only reads with a length greater than 10 bp were then mapped using Bowtie v.1.2.2 [28],

625   initially to the mm10 version of the mouse genome. The mapping allowed for up to 2

626   mismatches and reported only the best alignment (-m 1 --best --strata) for each read. To

627   ensure the counting of unique nascent RNA molecules, reads that mapped to the same

628   genomic location were collapsed based on their UMIs, allowing for up to 1 mismatch. To

629  create the PRO-seq coverage signal with the exact positions of RNA pol II molecules,

630  only the first nucleotide of each read (i.e. the 3' end of nascent transcripts) was considered

631  and the strand swapped to match the transcription direction. A non-redundant CAGE-

632  corrected gene set was used to count the number of UMI-collapesed 1nt-long mapped

633  PRO-seq reads that overlap them (see the "Gene annotation for PRO-seq analysis"

634  section). Differential analysis was performed using DESeq2[29] (v.1.22.2) and significantly

635  up- or down-regulated genes were selected using FDR<0.05, log2 Fold Change $\geq 1$

636  threshold.

637

638  *Cut&Run analysis*

639  Single-end 50-bp long reads were mapped to the mm10 genome using bowtie v.0.12.9,

640  allowing up to 3 mismatches and only uniquely mapping reads were retained. Afterwards,

641  peaks were called for each individual replicate, as well as for the combined replicates

642  against their respective input, using Macs2 v.2.1.2.1, with following settings: -f BEDPE -

643  g mm -B --nomodel --extsize 300 --SPMR. The Macs2 generated BedGraph files that

644  contain normalized coverage were converted into BigWig using bedGraphToBigWig.

645  Given the high correlation between two replicates (PCC of 0.613 at a common set of

646  peaks), only the merged sample was used for assigning bound genes if the peak was

647  localized within +- 500 bp around the gene transcription start sites.

648

649  **Data availability**

650  The raw sequencing data are available from GEO (https://www.ncbi.nlm. nih.gov/geo/)

651  under accession number GSE225972.

652

653  **Code availability**

654  All custom scripts that were generated for this study were made publicly available at

655  https://github.com/vloubiere/ORFtag_2023.

656

657

658

## References

12. Moussa, H. F. *et al.* Canonical PRC1 controls sequence-independent propagation of Polycomb-mediated gene silencing. *Nat Commun* **10**, 1–12 (2019).

13. Haberle, V. *et al.* Transcriptional cofactors display specificity for distinct types of core promoters. *Nature* **570**, 122–126 (2019).

14. Neumayr, C. *et al.* Differential cofactor dependencies define distinct types of human enhancers. *Nature* **606**, 406–413 (2022).

15. Serebreni, L. *et al.* Functionally distinct promoter classes initiate transcription via different mechanisms reflected in focused versus dispersed initiation patterns. *EMBO J* **42**, 1–24 (2023).

16. Hendy, O. *et al.* Developmental and housekeeping transcriptional programs in Drosophila require distinct chromatin remodelers. *Mol Cell* **82**, 3598-3612.e7 (2022).

17. Lawrence, M. *et al.* Software for Computing and Annotating Genomic Ranges. *PLoS Comput Biol* **9**, 1–10 (2013).

18. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**, 357–359 (2012).

19. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *Gigascience* **10**, 1–4 (2021).

20. Herzog, V. A. *et al.* Thiol-linked alkylation of RNA to assess expression dynamics. *Nat Methods* **14**, 1198–1204 (2017).

21. Lambert, S. A. *et al.* The Human Transcription Factors. *Cell* **172**, 650–665 (2018).

22. Soto, L. F. *et al.* Compendium of human transcription factor effector domains. *Mol Cell* **82**, 514–526 (2022).

23. Hentze, M. W., Castello, A., Schwarzl, T. & Preiss, T. A brave new world of RNA-binding proteins. *Nat Rev Mol Cell Biol* **19**, 327–341 (2018).

24. Tate, J. G. *et al.* COSMIC: The Catalogue Of Somatic Mutations In Cancer. *Nucleic Acids Res* **47**, D941–D947 (2019).

25. Blake, J. A. *et al.* Mouse Genome Database (MGD): Knowledgebase for mouse-human comparative biology. *Nucleic Acids Res* **49**, D981–D987 (2021).

26. Finn, R. D. *et al.* The Pfam protein families database. *Nucleic Acids Res* **38**, 211–222 (2009).

27. Forrest, A. R. R. *et al.* A promoter-level mammalian expression atlas. *Nature* **507**, 462–470 (2014).

28. Langmead, B., Trapnell, C., Pop, M. & Salzberg, S. L. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol* **10**, (2009).

29. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol* **15**, 1–21 (2014).

## Acknowledgements

**Author contributions**

715 F.N. and M.H. implemented the ORFtag method and protocols. F.N. performed the
716 activator screen, M.H. the PTGR screen, and R.Y. and V.L. the repressor screen. F.N.,
717 M.H. and R.Y. performed candidate validation experiments. F.N. performed all Zfp574
718 follow-up experiments. V.L., F.N. and A.S. developed the bioinformatic pipeline. V.L. and
719 F.N. performed NGS data and downstream analyses. N.F. and M.P. helped with the
720 experiments. U.E. and S.L.A conceptualized the ORFtag approach. S.L.A., A.S., U.E. and
721 J.B. coordinated and supervised the work. All authors wrote the manuscript.
722

**Competing interests**

724 S.L.A. is co-founder, advisor, and member of the board of QUANTRO Therapeutics
725 GmbH. U.E. is co-founder of JLP Health and VIVERITA as well as advisor to TANGO
726 Therapeutics. N.F. is employed by QUANTRO Therapeutics GmbH. The other authors
727 declare no competing interests.

# Figure 1

# Figure 2

# Figure 3

# Extended Data Figure 1