1  **GenEditID: an open-access platform for the high-throughput identification of CRISPR**

2  **edited cell clones.**

3  Ying Xue*[1], YC Loraine Tung*[2], Rasmus Siersbaek[3], Anne Pajon[3], Chandra SR

4  Chilamakuri[3], Ruben Alvarez-Fernandez[3], Richard Bowers[3], Jason Carroll[3], Matthew

5  Eldridge[3], Alasdair Russell[3], Florian T. Merkle[§,2,4]

6  1. Department of Endocrinology and Metabolism, Tongji Hospital of Tongji University,

7     Tongji University School of Medicine, Shanghai 200065, China

8  2. Metabolic Research Laboratories and Medical Research Council Metabolic Diseases

9     Unit, Wellcome Trust / Medical Research Council Institute of Metabolic Science,

10    University of Cambridge, Cambridge CB2 0QQ, UK

11 3. Cancer Research UK, Cambridge Institute, University of Cambridge, Cambridge CB2 0RE,

12    United Kingdom.

13 4. Wellcome Trust / Medical Research Council Cambridge Stem Cell Institute, University of

14    Cambridge, Cambridge CB2 0QQ, UK

15 * these authors contributed equally

16 [§] to whom correspondence should be addressed: fm436@medschl.cam.ac.uk

17

18

19

20

21

22

23

24

25   **ABSTRACT**

26   CRISPR-Cas9-based gene editing is a powerful tool to reveal genotype-phenotype

27   relationships, but identifying cell clones carrying desired edits remains challenging. To

28   address this issue we developed GenEditID, a flexible, open-access platform for sample

29   tracking, analysis and integration of multiplexed deep sequencing and proteomic data, and

30   intuitive plate-based data visualisation to facilitate gene edited clone identification. To

31   demonstrate the scalability and sensitivity of this method, we identified KO clones in parallel

32   from multiplexed targeting experiments, and optimised conditions for single base editing

33   using homology directed repair. GenEditID enables non-specialist groups to expand their

34   gene targeting efforts, facilitating the study of genetically complex human disease.

35

36   **KEYWORDS**

37   CRISPR-Cas9; gene editing; GWAS; Illumina sequencing; multiplexed; pluripotent stem cell,

38   LIMS, In-Cell Western

39

40   **BACKGROUND**

41   In the last decade, there has been an explosion of data from the sequencing of human

42   populations, including genome-wide association studies (GWAS) based on DNA microarrays

43   and increasingly also whole exomes and whole genomes (1). These studies have revealed

44   thousands of replicable genetic associations for complex diseases such as diabetes, obesity,

45   Alzheimer's Disease and breast cancer (2-4). However, mechanistically determining how

46   these genetic associations contribute to disease remains challenging. Causal evidence

47   requires careful functional follow-up experiments in model cellular systems, organisms and

48   eventually in humans, but the traditional approach of characterising one gene at a time

49   cannot keep pace with the rate of genetic discovery. Furthermore, many associated variants

50   are non-coding, so the genetic elements responsible for conferring disease risk are often

51   unclear (5, 6). This issue is exemplified by the fat mass and obesity associated (*FTO)* locus,

52    in which intronic SNPs are strongly associated with obesity, largely due to increased food

53    intake (7-10) irrespective of gender, age or ethnicity (11, 12). Despite intense study, the

54    identify of the genetic elements that mediate SNP-associated phenotypes remains

55    controversial. Some studies suggest that effect on appetite might not be driven by the *FTO*

56    gene itself as initially thought, but instead by the nearby genes retinitis pigmentosa GTPase

57    regulator-interacting protein-1 like (*RPGRIP1L*) (13-15), or by iroquois homeobox 3 (*IRX3*)

58    and iroquois homeobox 5  (*IRX5*) (16, 17). Two powerful tools have recently emerged to help

59    meet the challenge of uncovering disease mechanisms from the translating the growing

60    wealth of genetic data: human pluripotent stem cells (hPSCs) and the CRISPR-Cas9 system

61    (18, 19).

62

63    hPSCs facilitate human disease modelling since they can be indefinitely maintained in a

64    pluripotent state and can theoretically be differentiated into any cell type in the body,

65    including disease-relevant cell populations(20, 21). For example, hPSCs cell may be useful

66    in dissecting which genes near *FTO* contribute to increased food intake since they can be

67    differentiated into hypothalamic neurons that are pivotally important regulators of food intake

68    and that express these candidate genes (22, 23). The CRISPR-Cas9 system enables most

69    regions of the human genome to be efficiently edited. It consists of a ribonucleoprotein

70    (RNP) complex, including a Cas9 nuclease that is targeted to specific regions of DNA by an

71    approximately 20-base sequence within a guide RNA (gRNA) by forming a DNA-RNA hybrid

72    with complementary DNA sequences (24-26). For Cas9 isolated from the bacterium

73    *Streptococcus Pyogenes*, if the targeted DNA sequence contains a 3' protospacer adjacent

74    motif (PAM) of NGG, Cas9 will cleave the targeted DNA 3 bases 5' to the start of the PAM

75    site (27, 28) to create a double-strand break (DSB). The abundance of these PAM motifs in

76    the genome allows most genes to be targeted by CRISPR-Cas9 (29). DSBs can be repaired

77    by either the error-prone non-homologous end-joining (NHEJ) pathway which introduces

78    either frame-preserving or frameshift mutations (30), or by the homology-directed repair

79    (HDR) mechanisms which can be harnessed to introduce specific DNA alterations (31).

80    A major challenge in the field is how to effectively identify cell clones that have acquired

81    desired edits. Next-generation sequencing (NGS) of multiplexed pools of amplicons provides

82    an attractive solution to this problem, and CRISPR sequence analysis programmes built

83    around this idea provide visualisations of mutation types (32-34) and frequency (35).

84    However, to the best of our knowledge, there are no resources that provide a complete

85    platform for amplicon generation and barcoding, sample tracking, sequence analysis,

86    integration of distinct data forms (e.g. proteomic and sequencing), and intuitive visualisation

87    to empower investigators in non-specialist labs to pursue high-throughput targeted gene

88    editing.

89    To meet this challenge we developed a semi-automated, open-access, and user-friendly

90    pipeline that captures the nature and frequency of CRISPR-Cas9-induced gene editing to

91    identify cell clones of interest, which we call GenEditID. Briefly, targeted regions of interest

92    are amplified by PCR, barcoded, pooled and sequenced on an Illumina MiSeq. If genes of

93    interest are expressed in the targeted cell type, protein expression data can be integrated

94    with sequencing data to support the identification of knockout (KO) clones. Results are

95    graphically represented to reflect the physical location of the clone on the plate, facilitating

96    rapid and accurate clone recovery and further analysis. Using the *FTO* locus as an example,

97    we provide a roadmap by which the community can use GenEditID to rapidly, affordably, and

98    systematically explore genotype-phenotype relationships for genetically complex human

99    diseases. Furthermore, we demonstrate how the sequencing depth and multiplexed nature

100   of this approach enables targeted gene editing approaches to be optimised in cell

101   populations before embarking on the laborious process of gene targeting and clone picking,

102   allowing users to predict how many clones need to be picked in order to recover their clone

103   of interest.

104

105   **RESULTS**

106   We aimed to develop GenEditID to combine the strengths of laboratory information

107   management system (LIMS)-based sample management with open-access customisable

108   bioinformatic pipelines and a user-friendly graphical data display to facilitate the proliferation

109   of parallel cellular gene editing experiments by non-specialist groups (Fig.1). GenEditID was

110   implemented in Python, allowing basic experimental design and relevant sample details to

111   readily be incorporated into a genome editing report. To establish this platform, we first

112   turned to the estrogen receptor positive (ER+) breast cancer cell line MCF7, which is widely

113   used to study breast cancer biology and is amenable to CRISPR-Cas9 gene editing (36). We

114   targeted the oncogene signal transducer and activator of transcription 3 (*STAT3*), which is

115   expressed in MCF7 cells but remains largely inactive in the absence of extracellular stimuli

116   that trigger phosphorylation of tyrosine 705 (37). This approach allowed us to assess gene

117   editing efficiency and identify successfully edited clones without confounding factors such as

118   changes in the cell proliferation rate in response to deletion of *STAT3*, and to integrate

119   mutually supportive data from protein expression and DNA sequencing.

120   We designed four different gRNAs targeting exons 3 and 4 of *STAT3* (Fig. 2A,

121   Supplementary Fig. S1A) that we cloned into a Cas9 expressing vector (pSpCas9(BB)-2A-

122   GFP, PX458, Addgene#48138) (38), and separately transfected into a clonally-derived

123   MCF7 cell line stably transfected with a vector expressing mStrawberry and luciferase

124   (pCLIP-EF1a-LS). The use of a clonal cell line limited confounding factors associated with

125   comparing clonal edited cell lines with a polyclonal parental cell line. Successfully

126   transfected GFP+ single cells were distributed into 96 well plates using FACS, and clonal

127   colonies were allowed to form. Viable colonies were consolidated into a new set of 96 well

128   plates in duplicate, enabling one plate to be used to expand the clone for future use and the

129   second plate to be used for clone characterisation. We first characterised 107 clones of cells

130   by immunostaining wells with a *STAT3*-specific antibody, as well as a total cell stain used to

131   normalise for differences in cell confluence in a high-throughput Li-Cor In-Cell Western (Fig.

132   2B) (39). Based on ratios of fluorescence intensity in channels corresponding to STAT3

133     abundance and total cell staining, we identified gene-edited clones with reduced STAT3

134     immunostaining relative to non-edited controls (Fig. 2B, white arrow). We validated loss of

135     STAT3 protein for 8 such clones of interest using SDS-PAGE Western blotting

136     (Supplementary Fig. S1B). Next, we tested whether the low protein expression observed in

137     some targeted clones was due to CRISPR-Cas9-induced frameshifts in the *STAT3* genomic

138     sequence and not spurious effects due to stress or clonal selection.

139     To address this question, we extracted genomic DNA from 20 clones with low STAT3 protein

140     expression and then PCR amplified and Sanger sequenced amplicons across the *STAT3*

141     guide RNA target. We found that these clones indeed contained frameshift mutations

142     disrupting both wild-type alleles of *STAT3* (Supplementary Fig. S1C). However, gene edited

143     clones are often mosaic for a large number of alleles due to the persistence of CRISPR-

144     Cas9 upon plasmid transfection (40, 41), and less abundant alleles are difficult to detect by

145     Sanger sequencing. Therefore, we sought to test the extent to which sequencing information

146     predicts STAT3 protein expression status across all clones by developing a bioinformatic

147     pipeline for analysing sequencing reads for many clones in parallel. We reasoned that NGS

148     would provide ample read depth and accuracy to permit multiplexing and the detection of

149     low-abundance mutant alleles. We PCR amplified across the guide RNA target sites of

150     *STAT3* for 96 clones, appended unique barcodes to each clone, pooled the barcoded

151     amplicons, sequenced the pools, and bioinformatically identified amplicons arising from

152     distinct cell clones (Supplementary Fig. S2).

153     Across these clones, we observed a median sequencing depth of 14,1322 reads

154     corresponding to >90% of clones with at least 1000x coverage (Fig. 2C, Supplementary

155     Table 3, Supplementary Fig. S3A), providing ample power to call mutation allele frequencies

156     (Supplementary Fig. S3B). We next developed a sequence analysis pipeline to prioritize cell

157     clones with a high burden of mutations predicted to result in gene loss of function (Fig. 1). To

158     identify clones likely to have complete or near-complete gene KO, we aligned observed

159     sequencing traces to the reference genome and quantified the number of reads

160    corresponding to wild-type or variant sequence. We then omitted variants present at less

161    than 5% abundance, classified remaining variant types (e.g. synonymous, missense, in-

162    frame indels, frameshift indels) and assigned a score based on the likely consequence of

163    each variant type on gene function (Fig. 2D, Supplementary Fig. S4, Supplementary Table

164    3). To determine the total burden of predicted gene-disrupting variants each clone, we

165    calculated a "gene KO score" based on the aggregated product of mutant allele frequency

166    and predicted mutation consequence (see Materials and Methods). To visualise both protein

167    KO scores and gene KO scores, we implemented "heat maps" displaying these data based

168    on the physical location on the plate for each clone (Fig. 2E and 2F). Note that due to

169    differences in plate layout, only a subset of wells from plate 1 of this experiment (Fig. 2B)

170    were submitted for Illumina sequencing (blue circles in Fig. 2E and F). Since some users

171    may prefer to customise the calculation of KO scores or to integrate recently-developed

172    methods for calling and classifying mutation types such as AmpliCan (35), we have made

173    code developed for GenEditID freely available to the community at https://geneditid.github.io/

174    .

175    Next, we reasoned that the integration of gene KO and protein KO scores might provide

176    stronger evidence to support KO clone selection. To test if mutations called by our sequence

177    analysis pipeline predicts gene loss of function at the protein level, we compared protein KO

178    and gene KO scores (Fig. 2G). We found that while control samples and clones identified as

179    wild-type by sequence analysis tended to have similarly high STAT3 protein abundance,

180    clones with a high burden of missense and frameshift mutations had significantly lower ($R^2$ =

181    0.79, P < 0.0001) protein abundance, providing further confirmation of functional gene

182    ablation. We therefore took the product of gene and protein KO scores to calculate an

183    "integrated KO score" (Fig. 2H). These results indicate that count-based bioinformatic

184    analysis of multiplexed NGS data predicts functional gene disruption. While multiple lines of

185    evidence collected by high-throughput methods would be preferable to prioritise KO clone

186    selection, this is often not possible since genes of interest may not be strongly expressed in

187   the cell type used as the basis of gene editing, or appropriate antibodies may be lacking. For

188   example, genes in the *FTO* locus that are implicated in obesity by GWAS are expressed in

189   hypothalamic cells (42) but not are not highly expressed in hPSCs (23).

190   To test whether KO clones could be readily generated and identified across multiple genes

191   in parallel in hPSCs, which are more challenging to edit than cancer cell lines (42, 43), we

192   focused on the *FTO* locus. We first designed gRNAs to introduce double-strand breaks in

193   early constitutive coding exons of the genes *FTO*, *RPGRIP1L*, *IRX3*, and *IRX5* (Fig. 3A)

194   which are physically closest to the obesity-associated SNPs (Supplementary Fig. S5A).

195   Next, we *in vitro*-transcribed four sgRNAs per gene, combined them with purified Cas9

196   protein, and tested their ability to cut PCR-amplified target DNA *in vitro* to select maximally

197   active sgRNAs (Supplementary Fig. S5B and S5C). Since persistent Cas9 and sgRNA

198   expression from plasmids can promote clone mosaicism and off-target activity (40, 41), and

199   since double CRISPR-Cas9-induced strand breaks are cytotoxic to hPSCs (44, 45), we

200   nucleofected hPSCs with Cas9 protein complexed with *in vitro* transcribed sgRNA, which

201   has a half-life of approximately 24 hours in cells (46). hPSCs were re-plated at clonal density

202   of $2 \times 10^4$ cells/cm$^2$ and when colonies emerged we manually picked them into one 96-well

203   plate per targeted gene. After cells had reached approximately 70% confluence, we

204   duplicated plates to allow one plate of clones to be frozen down for later use and the other to

205   provide genomic DNA for our sequencing pipeline. This gene editing workflow is described in

206   detail elsewhere (47).

207   Next, we PCR-amplified CRISPR target sites for each targeted gene and added distinct

208   barcodes to each amplicon (Fig. 3B). The unique combination of amplicon and barcode

209   allowed these 378 samples across 4 distinct amplicons to be combined into a "superpool" as

210   previously described (48) along with an additional 987 samples from distinct experiments.

211   We found that across these amplicons, there was no significant bias in sequencing coverage

212   imposed by the barcodes (Supplementary Fig. S5D) and that the median sequencing

213   coverage for these amplicons was 13,689 reads, with 96% of amplicons having 1000 or

8

214   more reads (Supplementary Table 4). Using the pipeline we had built and tested using data

215   from *STAT3* targeting, we categorised variant types and calculated gene KO scores to

216   permit intuitive graphical data display (Fig. 3D), revealing efficient generation of KO clones

217   for all targeted genes except *RPGRIP1L*. This web-based graphical output allows users to

218   readily access data to retrieve clones of interest from highlighted wells for further analysis

219   from remote locations, such as a tissue culture room.

220   In addition to NHEJ-mediated frameshifts, some groups may want to introduce specific

221   mutations to test the consequence of disease-associated genetic variants or other functional

222   elements. To complement NHEJ-mediated KO of *FTO*, we also introduced a single base

223   mutation to introduce a premature stop codon into an early coding exon of *FTO* by HDR. To

224   this end, we designed a single-stranded oligodeoxynucleotide (ssODN) with sequence

225   complementarity to the target region but carrying the mutation of interest as well as a single-

226   base mutation to ablate the PAM to eliminate further activity of CRISPR-Cas9 at the edited

227   allele (Fig. 3E). Since HDR-mediated methods are inefficient, particularly in hPSCs (49), we

228   harnessed the multiplexing capability and sequencing depth of NGS to optimise conditions

229   for gene editing.

230   First, we identified optimal guides using *in vitro* cutting assays (Supplementary Fig. S5B).

231   Next, we nucleofected hPSCs with CRISPR-Cas9 RNP and systematically increased

232   concentrations of ssODN in biological triplicate. Rather than picking colonies, we extracted

233   gDNA from the nucleofected hPSC populations, barcoded each treatment condition, and

234   performed NGS (Fig. 3F). The sequencing depth provided by NGS allowed us to readily

235   detect the desired edit, revealing that the fraction of correctly edited alleles varied with

236   ssODN concentration, exceeding 10% allele frequency (approximately 20% of cells) for

237   some conditions. However, the relationship between ssODN concentration and editing

238   efficiency was non-linear and appeared to peak at a concentration around 100 pmol ssODN.

239   This result was confirmed in a second independent experiment (Fig. 3G). We also tested the

240   hypothesis that biotinylating Cas9 and adding a streptavidin tag to the biotinylating ssODN to

241    physically link the CRISPR-Cas9 complex to the repair oligo would increase rates of HDR,

242    as has been previously suggested (50). Despite clear Cas9 biotinylation, we found that this

243    strategy unexpectedly abolished Cas9 activity (Supplementary Fig. S6A) and HDR

244    (Supplementary Table 5), suggesting that future advances in gene editing could be rapidly

245    tested and optimised using multiplexed NGS of targeted cell populations.

246

247    **DISCUSSION**

248    We present GenEditID as a flexible, open-access pipeline that combines the benefits of a

249    web-based project management system, a bioinformatic pipeline for data analysis, and a

250    user-friendly graphical data output that allows efficient selections of clones of potential

251    interest. Our aim was to reduce the expense, labour, and requisite expertise for groups (or

252    core facilities) to carry out large-scale gene editing experiments. Here, we discuss the

253    benefits of GenEditID, areas where it could be further developed, and observations about

254    the nature of gene editing we observed in our dataset.

255    GenEditID can be readily customised and updated to incorporate new analysis tools (1) to

256    meet diverse needs. We have therefore published the underlying code in the public domain

257    https://geneditid.github.io/ . In particular, the AmpliCount tool we developed was designed to

258    rapidly analyse thousands of samples in parallel in order to prioritise clones for detailed

259    follow-up analysis. Since sample information is clearly associated with raw sequencing

260    traces by our sample tracking system, users can readily extract and further analyse data

261    from clones of interest at the sequence level (51) or with more specialised bioinformatic tools

262    (32-35), or incorporate these into GenEditID in lieu of AmpliCount.

263    Our analysis with AmpliCount revealed that most hPSC clones did not have simple

264    distribution of heterozygous, compound heterozygous, or homozygous edited alleles but

265    instead had a more complex mixture of alleles across a wide range of frequencies (Fig. 3C,

266    Supplementary Table 4). These "mixed clones" were most likely generated by CRISPR-Cas9

10

267  activity that persisted across several cell divisions soon after transfection, as previously

268  described by several groups (23, 40), despite the fact that we used a ribonucleoprotein

269  complex of Cas9 and *in vitro*-transcribed sgRNA, which has a much shorter half-life in cells

270  than plasmid- or virally-encoded Cas9 (46). These findings highlight the need to carefully

271  analyse NGS data of candidate clones for the presence of mosaicism, and if necessary to

272  perform a round of subcloning to isolate the desired clone.

273  Another advantage of the multiplexed barcoded amplicon sequencing pipeline of GenEditID

274  is that the high sequencing depth enables scalable, multiplexed analysis. In this study, we

275  ran 150 bp paired-end reads (300 cycles) on an Ilumina MiSeq. Assuming a minimum

276  desired read depth of 1000x and one full 96 well plate of clones per gene, >10 or >100

277  genes could be screened in parallel using the MiSeq nano or standard v2 kits, respectively

278  at a sequencing cost of less than 10 cents per clone. The establishment of a GenEditID

279  pipeline at an institution would allow amplicons from different research groups to be pooled,

280  facilitating cheaper and/or more frequent clone analysis. PCR barcoding by liquid handling

281  robots could further increase the throughput and decrease the cost of this approach. The

282  web-based data visualisation tool enables quick analysis of large numbers of clones and

283  modularly integrates data such as protein expression, cell growth, or allelic frequency (Fig.

284  1), increasing confidence in clone selection.

285  The scalability of CRISPR-Cas9-based gene editing combined with streamlined clone

286  selection provides a path to uncover the functional roles of disease-associated genes. This

287  catalogue of genes is rapidly growing due to the proliferation and increased sample size of

288  human population sequencing studies (52) In addition, the emergence of new methods for

289  differentiating hPSCs into diverse cell types promises to enable the interrogation of gene

290  function in disease-relevant cell populations. To illustrate this point, we targeted genes in the

291  *FTO* locus, which was the first locus associated with obesity risk loci by GWAS (53, 54). The

292  locus was subsequently linked to increased energy intake (7-10) but the identity of the

293  genetic elements at this locus (or elsewhere) that mediate obesity risk remains controversial.

294 Expression QTL studies of human cerebellum associate the obesity-linked SNPs to *IRX3*

295 expression (17), but since the cerebellum is not an area of brain normally recognised to be

296 involved in the control of food intake, there is a clear need to analyse brain regions pivotally

297 important for body weight regulation such as the hypothalamus (55), which can now be

298 generated from gene-edited hPCSs (22, 23). In this study, we report >80% mutation

299 efficiency in target genes at the *FTO* locus in hPSCs, of which approximately 50% are

300 frameshift mutations (Fig. 3C and 3D, Supplementary Table 4), demonstrating both the

301 efficiency of the gene editing method and the sensitivity of GenEditID to detect these edited

302 clones. For example, here we detected a clear dose-dependent effect of ssODN

303 concentration on the efficiency of the targeted introduction of a point mutation in *FTO*. As

304 CRISPR-Cas9 technology inexorably develops, we propose that new techniques can be

305 more readily optimised using the tools described here to quantify low-frequency gene editing

306 events within cell populations.

307

308 **CONCLUSION**

309 We developed GenEditID to be a flexible, open-access platform to enable groups to track

310 their samples, analyse and integrate amplicon sequencing and proteomic data. The

311 combined data analysis is intuitively visualised to facilitate edited clone identification. The

312 highly multiplexed approach enables the cost-effective and semi-automated identification of

313 targeted clones and provides a powerful platform to systematically investigate complex

314 human diseases in relevant cell types. We further show that GenEditID can be used to

315 rapidly optimise conditions for editing single bases using homology directed repair. Using

316 *FTO* as a proof of concept, we show how the platform can help the community to begin to

317 elucidate the mechanisms by which genetic variants contribute to human disease.

318

319 **METHODS**

320    **Cell lines and routine cell culture.** The clonal MCF7 breast cancer cell line expressing

321    mStrawberry and luciferase (a gift from Scott Lyons, Cold Spring Harbor Laboratory NY) was

322    grown in DMEM (41966-029, Gibco) supplemented with 10% fetal bovine serum (FBS,

323    10500-064, Gibco), 50 U/ml penicillin and 50 ug/ml streptomycin (15070-063, Gibco) and

324    2mM L-glutamine (25030, Gibco) in a humidified 37℃ incubator with 5% $CO_2$. The HUES9

325    human embryonic stem cell line was grown on tissue culture plates coated wth Geltrex

326    (Thermo Fisher Scientific) in mTeSR1 media (StemCell Technologies) and maintained in a

327    humidified 37℃ incubator with 5% $CO_2$. Medium was changed every 24 hours. Cells were

328    passaged with 1 mM ETDA for routine maintenance in mTeSR media supplemented with 10

329    µM ROCK inhibitor Y-27632 dihydrochloride (DNSK International). Please see below for

330    culture details during gene editing.  The absence of mycoplasma was confirmed using a EZ-

331    PCR Mycoplasma Test Kit (Supplementary Fig. S6B ; Biological Industries, 20-700-20)

332    following the manufacturer's instructions.

333

334    **CRISPR-Cas9-mediated targeting of *STAT3*.** Four different gRNAs with high predicted on-

335    target and low predicted off-target activity targeting exons 3 and 4 of STAT3 (NM_139276)

336    were designed using deskgen (www.deskgen.com). These guides were ordered from Sigma

337    Aldrich and cloned into pSpCas9(BB)-2A-GFP (PX458, Addgene #48138). A stably

338    transfected clonal MCF7 cell line expressing mStrawberry and luciferase (pCLIIP-EF1-LS)

339    was transfected with these vectors. Successfully transfected GFP+ cells were purified by

340    FACS into 6 well plates and allowed to recover for ~1 week. Single mStrawberry+ cells

341    (since the transfection with CRISPR-Cas9 vector was a transient transfection, the cells had

342    lost GFP expression at this point) were then distributed into multiple 96 well plates by FACS.

343    After ~3 weeks, viable clonal colonies were consolidated on new 96 well plates in duplicate,

344    so that one plate could be used for characterising clones and the other plate could be used

345    to expand clones for later use. All sgRNAs and primer sequences are provided in

346    Supplementary Table 1.

347

348 **In-Cell Western for STAT3.** Five days after seeding, the test plate used for clone

349 characterisation was fixed in 3.7% formaldehyde for 20 min at room temperature and

350 subjected to In-cell western using a STAT3-specific antibody (9139, Cell Signalling). Briefly,

351 cells were permeabilised by washing 5x in TBS + 0.1% Triton X-100 (Fisher Scientific,

352 BP151-100) for 5 min at room temperature and then blocked for 1 h with TBS Odyssey

353 blocking buffer (Li-Cor biosciences, 927-50000). Cells were then incubated with a STAT3

354 antibody (9139, Cell Signalling) in blocking buffer + 0.1% tween-20 (P1379, Sigma Aldrich)

355 for 1-2 hours at room temperature or overnight at 4℃ on a shaker. After washing 5x with

356 TBS+0.1% Tween-20 for 5 min, cells were incubated with secondary antibody (Goat anti-

357 mouse, 926-32210, Li-Cor) and a 1:500 dilution of CellTag 700 total cell stain (Li-Cor, 926-

358 41090) for 45min at room temperature on a shaker (50 ul/well). Cells were then washed 4x

359 with TBS+0.1% Tween-20 for 5 min followed by a final wash in TBS. Plates were then

360 analysed using the Odyssey CLx Imaging System (Li-Cor) to obtain measurements for both

361 total cell confluency and STAT3 expression. Signal intensity for STAT3 staining (green

362 channel) was then divided by the signal intensity for the CellTag 700 stain (red channel) for

363 each well to determine a STAT3 abundance ratio, and the resulting ratios were normalised

364 to produce a "protein KO score". The mean of negative control values, representing

365 background staining, was subtracted from all ratios, and the resulting scores were divided by

366 the mean of positive (non-edited) intensity scores, and the resulting values were subtracted

367 from 1 so that WT lines would have scores near 0 and KO lines would have scores near 1.

368

369 **Production and testing of *in vitro*-transcribed gRNA.** CRISPR guide RNAs were

370 designed to target early coding gene regions in constitutive exons of genes in the *FTO* locus

371 using Wellcome Trust Sanger Institute Genome Editing tool

372 (http://www.sanger.ac.uk/htgt/wge/) and the Feng Zhang laboratory's CRISPR design tool

373 (http://crispr.mit.edu/) to maximise on-target and minimize off-target activity. For the

374 production of gRNAs, a 120 nucleotide oligo (Integrated DNA Technologies Inc.) including

375 the SP6 promoter, gRNA sequences, and scaffold region were used as a template for

14

376 synthesis by *in vitro* transcription using the MEGAscript SP6 kit (Thermo Fisher, AM1330) as

377 previously described [19]. The resulting sgRNAs were purified using the E.Z.N.A miRNA

378 purification kit (Omega Bio-tek, R7034-01), eluted in RNase-free water, and stored at -80℃.

379 Since gRNAs vay in their efficacy, we designed at least four gRNAs per gene of interest and

380 tested their relative cutting efficiencies in *in vitro* cleavage assays as previously described

381 (47). We selected the gRNAs that showed activity at the lowest Cas9 concentration at each

382 target gene for transfection in the hPSC cells. All sgRNAs and primer sequences are

383 provided in Supplementary Table 1.

384

385 **Cas9 protein production.** Cas9 proteins were purified by the laboratory of Marko Hyvönen

386 (University of Cambridge) from E. coli expressing *Streptococcus Pyogenes* Cas9 carrying a

387 C-terminal fusion to a hexa-histidine tag from the pET-28b-Cas9-His plasmid (Addgene

388 http://www.addgene.org/47327) (56). The soluble Cas9 protein was purified by a

389 combination of nickel affinity and cation exchange chromatographies. The purified protein

390 was concentrated to approximately 30 μM (4.8 mg/ml) in 20 mM HEPES pH 7.5, 500 mM

391 KCl and 1% sucrose buffer and flash frozen for storage at -80°C.

392

393 **CRISPR-Cas9 ribonucleoprotein (RNP) complex-mediated editing in hESCs.** For gene

394 knockout by NHEJ, 3 μg purified sgRNA was mixed with 4 μg Cas9 protein (final volume <5

395 μl) for 10 min at room temperature to form stable RNP complexes. The complex was then

396 transferred to a 20 μl single-cell suspension of $2 \times 10^5$ hESCs in P3 nucleofection solution

397 and electroporated using Amaxa 4D-Nucleofector™ (Lonza) with program CA137.

398 Transfected cells were seeded onto Geltrex-coated 10 cm dishes containing a pre-warmed

399 1:1 mix of mTeSR1 and hESC medium containing 20% knockout serum replacement

400 (KOSR) and 100 ng/ml bFGF, supplemented with 10 μM ROCK inhibitor. Rock Inhibitor was

401 withdrawn after 24 hours. Single colonies were isolated manually 7-10 days after

402 transfection and seeded into Geltrex-coated 96-well plates in 1:1 medium plus ROCK

15

403    inhibitor, which was withdrawn after 24 hours. A total of 96 individual colonies were picked

404    for each targeted gene, and maintained in 1:1 medium for 10-14 days. Once clones were

405    close to confluent, each of the 96 well plates were duplicated by EDTA passaging to allow

406    parallel cell cryopreservation and genomic DNA extraction as previously described (47).

407    Briefly, cell cryopreservation was performed in hESCs culture media containing a final

408    concentration of 40% FBS and 10% DMSO. Cells were slowly frozen in ice-cold freezing

409    media using Mr. Frosty™ Freezing Container (Thermo Fisher Scientific, 5100-0001). During

410    the thawing of hESCs, media was supplemented with CloneR (StemCell Technologies,

411    05889) for 72 hours.

412

413    **Generation and sequencing of pooled amplicons.** Genomic DNA was extracted using

414    HotShot buffer as previously described (47). The target regions were amplified from gDNA

415    using locus-specific primers to generate amplicons approximately 150-200 bp in length as

416    previously described (47). These "first-round" primers contained universal Fluidigm linker

417    sequences    at    their    5'-end    with    the    following    sequences:    Forward    primer:    5'-

418    acactgacgacatggttctaca -3', Reverse primer: 5' tacggtagcagagacttggtct-3'. Specifically, 20 µl

419    PCR reactions were set up in 96 well plates using 1U FastStart high fidelity polymerase

420    (Roche, 3553361001), 2 µl of extracted gDNA as template, 2 µl 10x HF buffer without $MgCl_2$,

421    0.2 mM dNTPs, 0.2 µM primers, and 4.5 mM $MgCl_2$, and run on the following programme:

422    95°C 2 min, followed by 36 cycles of (95°C 20 sec, 64.4°C 20 sec, 72°C 15 sec), 72°C 3

423    min. In the second round of PCR (indexing PCR), Fluidigm barcoding primers were attached

424    to the amplicons to uniquely identify each clone. 2 µl linker PCR product diluted 1:10 was

425    transferred to another 96-well PCR plate to perform this indexing PCR in 20 µl reactions

426    containing 0.04 µM of Fluidigm barcoding primers (Supplementary Table 1), 2 µl 10x HF

427    buffer    without    $MgCl_2$,    0.2    mM    dNTPs,    4.5    mM    $MgCl_2$,    and    1U    FastStart    high    fidelity

428    polymerase (Roche, 3553361001). The PCR programme was 95°C 2 min, 16 cycles of

429    (95°C 20 sec, 60°C 20 sec, 72°C 25 sec), 72°C 3 min. For sequencing library preparation,

430    barcoded PCR products were combined in equal proportion based on estimation of band

431    intensity on a 2% agarose gel, and the combined pool of PCR products was purified in a

432    single tube using Ampure XP beads (Beckman-Coulter, A63880) at 1:1 (V/V) to the pooled

433    sample, and eluted in 25 µl of water according to the manufacturer's instructions. Library

434    purity was confirmed by nanodrop, and final library concentration was measured using the

435    Qbit fluorometer and diluted to 20 nM. Pooled libraries could be combined with other library

436    pools adjusted to 20 nM, and the resulting "superpool" volume was adjusted to a final

437    volume of 20 µl before sequencing.

438

439    **Introduction of STOP codon in FTO through HDR-mediated repair.** To target a STOP

440    codons to an early coding exon of FTO, we designed a ssODN template of 90 bp in length to

441    be homologous to the target site but to contain single base mismatches to introduce a STOP

442    codon and to ablate the PAM site to prevent re-cutting by CRIPSR/Cas9 as previously

443    described (47)(Supplementary Table 1). ssODNs were synthesized (Integrated DNA

444    Technologies Inc.), dried, and re-suspended in nuclease-free sterile water to a final

445    concentration of 100 µM. Various amounts of ssODNs ranging from 20 pmol to 312.5 pmol

446    were added to RNP complexes for nucleofection as described above. Editing efficiency was

447    determined by sequencing the targeted locus using primers outside of the ssODN at the cell

448    population level rather than in single picked colonies, and then counting the number of reads

449    corresponding to the WT amplicon sequence, or sequences with one or both desired edits.

450

451    **Project tracking within the GenEditID web framework.** We designed a Python-based web

452    framework (web app) of GenEditID to facilitate the tracking of different projects, the tracking

453    of samples within a project, and to facilitate the plate-based data integration and

454    visualisation to help users identify clones of interest. When initiating a project, the user first

455    creates a project via the web app (for a screenshot example of the web app home page see

456    Supplementary Fig. S7) along with comments about the project purpose and design, and

17

457     then submits an Excel configuration file containing the plate and well each sample originated

458     from, the CRISPR sequences used, the primers and barcodes used for NGS analysis, and

459     any other pertinent information (Fig. 1, Supplementary Table 2). This project information is

460     later accessible via the web app and also programmatically, and enables samples to be

461     uniquely identified so that data from different analysis modalities (e.g. sequencing, growth

462     rate, protein abundance) can be loaded into the web app for tracking, integration, and

463     visualisation. We designed the web app to enable integration with a laboratory information

464     management system (LIMS) to automatically trigger a sequencing request when sample

465     information is uploaded. Links to specific LIMS are intentionally omitted from the code

466     published here since we anticipate that different institutions will wish to implement GenEditID

467     within existing systems. Analysis can then be run outside the web app using setup scripts

468     generated from the database, as described in further detail below.

469

470     **Amplicon analysis from NGS data with AmpliCount.** The analysis of barcoded PCR

471     amplicons    is    performed    outside    of    the    web    app    using    modular    scripts

472     (https://geneditid.github.io/ ) that are adaptable to each user's specific requirements. First,

473     FASTQ files associated to the project are retrieved and configuration files are created to link

474     sequencing information with the sample information stored in the GenEditID project

475     configuration file. These de-multiplexed FASTQ files are then either merged, or joined using

476     fastq-join if the target size is larger than the read length.

477     To analyse reads, we developed "ampli_count" (Supplementary Fig. S4), a tool that first

478     finds amplicons using primer pairs and group variants with same sequence. Then we

479     identified and filtered out reads of low quality across a 5bp sliding window (average read

480     quality score < 10). Putative primer dimers were defined as any sequence smaller than the

481     combined size of the forward and reverse primers, plus 10bp, and discarded. To focus

482     downstream analysis on variants that reflect CRISPR-Cas9-induced edits rather than

483     sequencing artefacts, sequences supported by 60 or fewer reads were also discarded. After

484    obtaining "filtered reads" that passed these criteria, amplicon-specific variants were identified

485    using the tool "variant_id" to determine variant type and consequence per site by pairwise

486    alignment to human reference genome ensembl_grch38 using pairwise2 from Biopython

487    (https://github.com/biopython/biopython), and using varcode

488    (https://github.com/openvax/varcode) and pyensembl

489    (https://github.com/openvax/pyensembl) to determine consequence. Since the aim of this

490    analysis was to identify clones that carried a high burden of variants likely to lead to loss of

491    gene function rather than generate a comprehensive description of variants observed upon

492    gene editing, only variants with an overall frequency of 5% or higher were retained for

493    downstream analysis. All filter steps and thresholds are tuneable by the user.

494    Remaining variants were classified according to their predicted consequence

495    (Supplementary Fig. S4). Variants with multiple predicted consequences following sequence

496    alignment were labelled "Complex", or "ComplexFrameShift" if they contained a frameshift.

497    Consequences were then given an impact weighted score based on their predicted effect on

498    gene function, ranging from 0 for wild-type sequences to 1 for the gain of a premature stop

499    codon. Variants were then grouped by similar consequence categories and these combined

500    frequencies were multiplied by an impact weighting score and summed across all

501    consequence categories to yield a "gene KO score" for each allele, where a score of 0 would

502    correspond to all WT sequences, and a score of 1 would correspond to all predicted

503    deleterious variants (Supplementary Fig. S4). Variant classification and weighting for KO

504    score calculation can be readily altered via csv file.

505

506    **Clone score integration and data visualisation.** After computing protein KO scores and

507    gene KO scores, data were loaded back into the GenEditID database to facilitate their

508    integration with stored sample information. Where both scores were available, the

509    "integrated KO score" was calculated by taking the product of the gene and protein KO

510    scores. To facilitate the selection and expansion of candidate KO clones, information about

511    each clones' plate and well position was used to graphically display computed scores as a

512    "heat map" in 96 well plate format.

513

514    **Statistical analyses.** Statistical analyses were performed using Graph Pad Prism version

515    8.1.0. A two-tailed Student's *t*-test was performed to compare knock-in efficiencies among

516    different conditions for variable amount of ssODNs. Unless otherwise stated, data shown

517    represent the results of at least three independent experiments. P-values < 0.05 were

518    considered significant.

519

520    **List of abbreviations**

521    GWAS: genome-wide association studies;

522    FTO: fat mass and obesity associated;

523    RPGRIP1L: retinitis pigmentosa GTPase regulator-interacting protein-1 like;

524    IRX3: iroquois homeobox 3;

525    IRX5: iroquois homeobox 5;

526    hPSCs: human pluripotent stem cells;

527    RNP: ribonucleoprotein;

528    gRNA: guide RNA;

529    PAM: protospacer adjacent motif;

530    DSB: double-strand break;

531    NHEJ: non-homologous end-joining;

532    HDR: homology-directed repair;

533    NGS: Next-generation sequencing;

534    KO: knockout;

535     LIMS: laboratory information management system;

536     ER+: estrogen receptor positive;

537     STAT3: signal transducer and activator of transcription 3;

538     ssODN: single-stranded oligodeoxynucleotide;

539     KOSR: knockout serum replacement;

540     Bio-Cas9: biotinylated form of Cas9

541

542     **DECLARATIONS**

543     **Ethics approval and consent to participate**

544     Not applicable

545

546     **Consent for publication**

547     Not applicable

548

549     **Availability of data and material**

550     The datasets generated and analysed during the current study have been deposited in

551     NCBI's Sequence Read Archive at

552     https://dataview.ncbi.nlm.nih.gov/object/PRJNA543767?reviewer=ktlduo7ptjcsmrajrhnfnj0su

553     2 (knockout of human *STAT3*);

554     https://dataview.ncbi.nlm.nih.gov/object/PRJNA543845?reviewer=h1v7go700g7n1ocftv3hckr

555     k3m (multiplexed knockout of human *RPGRIP1L, FTO, IRX3,* and *IRX5*); and

556     htps://dataview.ncbi.nlm.nih.gov/object/PRJNA545266?reviewer=9ptviqm1j6bpfubtjdherr4r0

557     8 (targeted editing of human *FTO*).

558

559     **Competing interests**

560     The authors declare that they have no competing interests.

561

570

**Authors' contributions**

571

572   YX, YCLT, and FTM conceived the project and wrote the manuscript with contributions from

573   all other authors. RS generated data from *STAT3* targeting with the guidance of JC. YX and

574   YCLT generated data from gene knockout and targeted gene editing of *FTO* and

575   neighboring genes with the guidance of FTM. AP and CSRC generated the bioinformatics

576   plots. AP, CSRC, AB and RAF developed bioinformatic tools under the guidance of ME and

577   AR.

578

584

585

586   **FIGURES AND FIGURE LEGENDS**

587     **Figure 1) GenEditID facilitates the management and interpretation of multiplexed gene**

588     **editing projects.** After the design of CRISPR sequences to target one or multiple genes,

589     sgRNAs are introduced along with Cas9 to a cell line of choice, which is then subjected to

590     single-cell cloning by FACS or manual picking. The details of the CRISPR sequences, the

591     location and plate that identifies each picked clone, the primers used to amplify the targeted

592     region, any barcodes used to distinguish between PCR amplicons, and any other information

593     pertaining to experimental conditions is entered into a centralised project and sample

594     management system to track samples across projects. Next, data from assays across

595     modalities (e.g. protein analysis by In-Cell Western and mutation analysis by barcoded PCR)

596     feeds into GenEditID's where it integrates with its sample tracking features and is analysed,

597     for example to assess the burden of deleterious mutations carried by each clone (gene

598     KOscore). Scores across these analysis modalities are graphically illustrated to facilitate the

599     selection and expansion of cell clones for further analysis.

600

601     **Figure 2) Integration of protein- and sequence-based information to inform knockout**

602     **clone selection. A)** Experimental schematic showing the CRISPR targeting of early

603     constitutive exons of *STAT3* in the breast cancer cell line MCF7. After gene targeting, cells

604     were clonally isolated by FACS into 96-well plates that were duplicated to facilitate clone

605     propagation, protein analysis by In-Cell Western, and the barcoded PCR amplification of the

606     targeted locus to facilitate deep sequencing and the identification of clones harbouring a

607     high burden of deleterious (e.g. frameshift) mutations. **B)** Images from an In-Cell Western

608     experiment where cells were stained for a whole-cell dye and with an anti-STAT3 antibody,

609     and experiment clones were grown in the inner wells and wild-type controls in opposite

610     corners (yellow arrows), revealing clones with reduced STAT3 expression (white arrow). **C)**

611     Most samples have >1000 reads, providing good depth for assessing variant allele

612     frequency. Low quality, low abundance, or primer dimer reads are discarded from

613     subsequent analysis. **D)** Distribution of mutation frequencies and consequences across

614     targeted cell clones used for calculating gene KO scores, where variants with an allele

615   frequency of <5% (grey) do not contribute to these scores. **E,F)** Spatial heat maps of STAT3

616   protein KO score (E, compare to lower panel in B) and *STAT3* gene KO score. Note that

617   only a subset of wells from E were analysed by Illumina sequencing (blue circles). **G)** Linear

618   regression of STAT3 protein KO score and *STAT3* gene KO score shows a strong

619   correlation ($R^2$=0.79) between these independent measures of gene function. **H)** Heatmap of

620   the integrated KO score (product of protein and gene KO scores).

621

622   **Figure 3) Multiplexed identification of knockout and knock-in clones by GenEditID. A)**

623   Schematic showing the spatial location of genes in the human *FTO* locus. **B)** Parallel KO of

624   multiple genes in separate experiments, the generation of barcoded amplicons, and the

625   pooling of these barcoded amplicons for multiplexed sequencing. **C)** Mutation allele

626   frequency per targeted clone for each gene. **D)** Plate-based heat maps showing the location

627   of clones with high gene KO scores (blue). **E)** Schematic of targeted ssODN-mediate

628   mutation knock-in (red bases) in the first coding exon of *FTO*, to ablate the CRISPR PAM

629   sequence (blue) and introduce and early stop codon (red) to ablate protein production. **F)**

630   Schematic for the systematic modulation of CRISPR conditions and assessment of knock-in

631   efficiency in un-cloned cell populations. **G)** Calculated knock-in efficiencies revealed a

632   correlation with ssODN concentration that appeared to peak near 100-125 pmol. n=3 and

633   n=6 for each concentration for Study 1 and Study 2, respectively.

634

635

636   **SUPPLEMENTARY FIGURES**

637   **Supplementary Figure S1. Validation of GenEditID with the clonal selection of *STAT3***

638   **KO lines. A)** A schematic diagram depicting the site-specific CRISPR designs to target exon

639   3 or 4 of the human *STAT3* gene. **B)** STAT3 protein expression by Fluorescent Western blot

640   confirmed that STAT3 (green) was effectively knocked down in the MCF7 cells. Actin (red)

641   was used as loading control. **C)** Comparison of growth rate with loss of STAT3 protein

642   expression (ratio to total cell stain) and presence of insertion/deletion mutations (indels) in

643    Sanger sequencing data. **D)** Sanger sequencing at the target site showing sequence

644    chromatograms for clone B6 from plate 4 (right) showing an "A" insertion introduces a

645    premature STOP codon. Clone E2 from plate 1 (left) as an example of STAT3 wild type at

646    the corresponding locus.

647

648    **Supplementary Figure S2. PCR amplicon barcoding and pooling for NGS. A)** Genomic

649    DNA is extracted from cell clones 96 well plates that were duplicated from another plate

650    which is maintained or frozen. **B)** The CRISPR-targeted region is PCR amplified with locus-

651    specific primers containing universal (L1 and L2) linker sequences (orange box). **C)** In a

652    second PCR cycle (blue box), one of 96 unique Fluidigm forward barcodes are added to the

653    end of the PCR product. To further increase the multiplexing capacity, up to 96 distinct

654    reverse barcodes can be used. The resulting barcoded amplicons are then concentration

655    normalised, pooled, and sequenced on the Illumina MiSeq in the presence of PhiX to

656    increase library diversity.

657

658    **Supplementary Figure S3. Amplicon sequencing coverage and criteria used for**

659    **variant quality control.** $Log_{10}$-transformed distribution of read depth (A) and variant

660    frequency (B) at one of the two PCR amplicons generated at the human *STAT3* locus,

661    showing criteria for inclusion for downstream analysis: minimum read depth >1000 after

662    filtering (A), and minimum variant allele frequency >5% (B).

663

664    **Supplementary Figure S4. Schematic flowchart of steps involved in variant calling and**

665    **annotation by AmpliCount.** Throughout this flowchart, green indicates data outputs,

666    orange indicates tools or file inputs, and processing steps are indicated in white. FASTQ files

667    are first demultiplexed and reads are combined by merging or joining, depending on

668    amplicon length (please see Materials and Methods). The "ampli_count" tool then filters out

669    reads with low sequence quality, reads likely corresponding to primer dimers, and reads

670    corresponding to low abundance (e.g. <60 reads) sequences to retain filtered reads from

671    which variant frequencies are computed. The "variant_id" tool then removes all reads that

672    have low abundance (e.g. <5%) relative to all filtered reads to streamline downstream

673    alignment and variant classification steps. Remaining variants are pairwise-aligned,

674    classified by variant consequence, scored as indicated. Variants in the same consequence

675    categories are combined, and the combined frequencies are multiplied by the consequence

676    scores and summed to yield gene KO scores for each clone (see also Fig. 1). The resulting

677    data outputs include plots and tables, including graphical visualisation of gene KO scores in

678    plate-based heat maps.

679

680    **Supplementary Figure S5. Multiplexed analysis of CRISPR targeting of the human *FTO***

681    **locus. A)** SNPs within intron 1 of *FTO* have a strong association with obesity, implicating the

682    causal involvement of several nearby genes. **B)** Representative gel image of an *in vitro*

683    cleavage assay evaluating gRNA:Cas9-mediated target cleavage. The arrowheads indicate

684    nuclease cleaved products. **C)** Schema of the CRISPR RNA targeting early conserved

685    coding exons in four candidate genes in the *FTO* locus. The CRISPR recognition sequence

686    is shown in red and the PAM sequence is shown in blue. **D)** $Log_{10}$-transformed sequencing

687    depth per barcode of the four targeted sequenced and analysed in parallel indicating ample

688    sequence quality and depth for variant analysis.

689

690    **Supplementary Figure S6. A)** Results from an *in vitro* Cas9 cutting assay where a

691    biotinylated form of Cas9 (Bio-Cas9) was mixed with a biotinylated-ssODN and streptavidin

692    was added at increasing concentrations to physically link Cas9 and the ssODN repair oligo,

693    but this approach appeared to disrupt cutting ability, as indicated by the lower molecular

694    weight bands observed with non-biotinylated (WT) Cas9. **B)** All hESC cell lines were tested

695    for mycoplasma before and after gene editing and found to be negative.

696

697     **Supplementary Figure S7.** A screenshot of an example of the homepage of the Web App in

698     the web browser to facilitate the tracking of different projects. Some key features of the Web

699     App including "Create project" which enables the creation of a new project. "edit" allows the

700     uploading of the configuration file that are subsequently loaded into the database to facilitate

701     sample tracking. "view" displays target amplicon and a list of the samples. "ngs data" house

702     all the resulting data outputs including the csv tables and plots.

703

704     **SUPPLEMENTARY TABLES**

705     **Supplementary Table 1. Compilation of primer, barcode, CRISPR and ssODN**

706     **sequences used in this study.** The table contains all the sequence information used in all

707     the studies covered in this manuscript, including the CRISPR guide RNA sequences, target

708     primer sequences used for the *in vitro* Cas9 cutting assay and target primer sequences for

709     the first round (linker) PCR for NGS. The ssODN sequence is complementary to the *FTO*

710     target region and carries mutations to introduce a premature STOP codon and a mutation to

711     disrupt the PAM site to prevent further CRISPR-Cas9 activity.

712

713     **Supplementary Table 2. Example of a configuration file used to submit and track**

714     **samples.** An example of a configuration file that are submitted by the user when initiating a

715     project that contains all the relevant information from the user including the CRISPR

716     sequences used, primers sequences for NGS analysis and sample coordination on the plate.

717     All the information in the configuration file is later accessible via the web app allowing

718     samples to be uniquely identified for tracking, integration and visualisation.

719

720     **Supplementary Table 3. Amplicon data from the human *STAT3* locus generated with**

721     **AmpliCount.** "**Config STAT3**" contains the information including the target genomic location

722     and the two sets of primers use to identify the two amplicons of interest. "**variantid STAT3**"

723     contains the total reads after each step of the filtering process to remove reads with low

724     quality, short length (primer dimers), or low abundance, and provides information on the

725 variant frequency, type and consequence to give a variant score where the score is the

726 consequence weight x variant frequency. "**impact STAT3**" tablet the combined frequency of

727 all variants with identical consequence categories (high/medium/low impact). Finally, the

728 "**koscores_amplicon**" provides individual table for each amplicon (i.e. STAT3 exon3 &

729 exon4) on variants grouped in each consequence categories (i.e. variant with identical

730 impact weighing) and their combined frequencies to yield a gene KO score where the score

731 is the sum of impact weight x impact frequency for each allele and a score of 1 would

732 correspond to all deleterious variants.

733

734 **Supplementary Table 4. Amplicon data from the human *FTO* locus generated with**

735 **AmpliCount.** "**Config FTO plus**" contains the genomic location and the primer pairs use to

736 identify the four target gene of interest, namely *FTO, IRX3*, *IRX5* and *RPGRIP1L.* "**variantid**

737 **FTO plus**" table contains all total reads for each barcode and sequences reads after each

738 step of the filtering process such as reads with low quality, primer dimers, low abundance,

739 and lists the variant type and consequence to give a variant score. "**impact FTO plus**"

740 shows data on the combined frequency of all variants with identical consequence weighing

741 and categories the consequence into high/medium/low impact. "**koscores_amplicon**" tablet

742 each of the four amplicons included in this study and provides data on variants grouped by

743 same consequences categories and their combined frequencies to yield a gene KO score for

744 each allele where a score of 1 would correspond to all deleterious variants.

745

746 **Supplementary Table 5. Targeted editing with ssODNs to introduce targeted**

747 **mutations into *FTO*.** "**config ssODN**" lists the four types of variant sequence that were

748 tested for, namely the wild-type *FTO* sequence, the sequence with both target sites mutated,

749 the variant sequence with target site only where a STOP codon was introduced in the early

750 coding exon of FTO and sequence where only the PAM site was mutated. "**amplicount**

751 **Study 1**" shows the total and filtered (e.g. matching target sequences named above)

752 sequencing reads in eight cell populations gene edited with either WT Cas9 or Cas9-biotin

28

753 proteins (for ssODN concentration of 20 pmol, 50 pmol, 125 pmol and 312.5 pmol; n=3 for

754 each condition). "**amplicount Study 2**" shows the total and filtered sequencing reads in

755 another independent four groups of un-cloned cell populations (for ssODN concentration of 0

756 pmol, 50 pmol, 100 pmol and 200 pmol; n=6 for each condition). "**KoIN Study 1**" and "**KoIN**

757 **study 2**" shows the calculated knock-in efficiencies for Study 1 and Study 2 respectively

758 where the % of HDR efficiency were represented by the percentage of specific variant

759 filtered read over the sum of all filtered reads for each barcode.

760

## REFERENCES

761 **REFERENCES**

762 1. McCarroll SA, Kuruvilla FG, Korn JM, Cawley S, Nemesh J, Wysoker A, et al. Integrated

763 detection and population-genetic analysis of SNPs and copy number variation. Nat

764 Genet. 2008;40(10):1166-74.

765 2. Hirschhorn JN. Genetic approaches to studying common diseases and complex traits.

766 Pediatr Res. 2005;57(5 Pt 2):74R-7R.

767 3. Johnson GC, Todd JA. Strategies in complex disease mapping. Curr Opin Genet Dev.

768 2000;10(3):330-4.

769 4. Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, et al. The NHGRI

770 GWAS Catalog, a curated resource of SNP-trait associations. Nucleic Acids Res.

771 2014;42(Database issue):D1001-6.

772 5. MacArthur DG, Manolio TA, Dimmock DP, Rehm HL, Shendure J, Abecasis GR, et al.

773 Guidelines for investigating causality of sequence variants in human disease. Nature.

774 2014;508(7497):469-76.

775 6. Maurano MT, Humbert R, Rynes E, Thurman RE, Haugen E, Wang H, et al. Systematic

776 localization of common disease-associated variation in regulatory DNA. Science.

777 2012;337(6099):1190-5.

778    7.    Cecil JE, Tavendale R, Watt P, Hetherington MM, Palmer CN. An obesity-associated

779         FTO gene variant and increased energy intake in children. N Engl J Med.

780         2008;359(24):2558-66.

781    8.    Speakman JR, Rance KA, Johnstone AM. Polymorphisms of the FTO gene are

782         associated with variation in energy intake, but not energy expenditure. Obesity (Silver

783         Spring). 2008;16(8):1961-5.

784    9.    Wardle J, Carnell S, Haworth CM, Farooqi IS, O'Rahilly S, Plomin R. Obesity

785         associated genetic variation in FTO is associated with diminished satiety. J Clin

786         Endocrinol Metab. 2008;93(9):3640-3.

787    10.   Wardle J, Llewellyn C, Sanderson S, Plomin R. The FTO gene and measured food

788         intake in children. Int J Obes (Lond). 2009;33(1):42-5.

789    11.   Jacobsson JA, Schioth HB, Fredriksson R. The impact of intronic single nucleotide

790         polymorphisms and ethnic diversity for studies on the obesity gene FTO. Obesity

791         reviews : an official journal of the International Association for the Study of Obesity.

792         2012;13(12):1096-109.

793    12.   Sovio U, Mook-Kanamori DO, Warrington NM, Lawrence R, Briollais L, Palmer CN, et

794         al. Association between common variation at the FTO locus and changes in body mass

795         index from infancy to late childhood: the complex nature of genetic association through

796         growth and development. PLoS genetics. 2011;7(2):e1001307.

797    13.   Stratigopoulos G, Burnett LC, Rausch R, Gill R, Penn DB, Skowronski AA, et al.

798         Hypomorphism of Fto and Rpgrip1l causes obesity in mice. J Clin Invest.

799         2016;126(5):1897-910.

800    14.   Stratigopoulos G, Martin Carli JF, O'Day DR, Wang L, Leduc CA, Lanzano P, et al.

801         Hypomorphism for RPGRIP1L, a ciliary gene vicinal to the FTO locus, causes increased

802         adiposity in mice. Cell metabolism. 2014;19(5):767-79.

803    15. Wang L, De Solis AJ, Goffer Y, Birkenbach KE, Engle SE, Tanis R, et al. Ciliary gene

804        RPGRIP1L is required for hypothalamic arcuate neuron development. JCI Insight.

805        2019;4(3).

806    16. Claussnitzer M, Dankel SN, Kim KH, Quon G, Meuleman W, Haugen C, et al. FTO

807        Obesity Variant Circuitry and Adipocyte Browning in Humans. N Engl J Med.

808        2015;373(10):895-907.

809    17. Smemo S, Tena JJ, Kim KH, Gamazon ER, Sakabe NJ, Gomez-Marin C, et al. Obesity-

810        associated variants within FTO form long-range functional connections with IRX3.

811        Nature. 2014;507(7492):371-5.

812    18. Merkle FT, Eggan K. Modeling human disease with pluripotent stem cells: from genome

813        association to function. Cell stem cell. 2013;12(6):656-68.

814    19. Sandoe J, Eggan K. Opportunities and challenges of pluripotent stem cell

815        neurodegenerative disease models. Nature neuroscience. 2013;16(7):780-9.

816    20. Cohen DE, Melton D. Turning straw into gold: directing cell fate for regenerative

817        medicine. Nat Rev Genet. 2011;12(4):243-52.

818    21. Murry CE, Keller G. Differentiation of embryonic stem cells to clinically relevant

819        populations: lessons from embryonic development. Cell. 2008;132(4):661-80.

820    22. Kirwan P, Jura M, Merkle FT. Generation and Characterization of Functional Human

821        Hypothalamic Neurons. Current protocols in neuroscience. 2017;81:3.33.1-3..24.

822    23. Merkle FT, Maroof A, Wataya T, Sasai Y, Studer L, Eggan K, et al. Generation of

823        neuropeptidergic hypothalamic neurons from human pluripotent stem cells.

824        Development (Cambridge, England). 2015;142(4):633-43.

825    24. Cong L, Ran FA, Cox D, Lin S, Barretto R, Habib N, et al. Multiplex genome engineering

826        using CRISPR/Cas systems. Science. 2013;339(6121):819-23.

827    25. Mali P, Yang L, Esvelt KM, Aach J, Guell M, DiCarlo JE, et al. RNA-guided human
828        genome engineering via Cas9. Science. 2013;339(6121):823-6.

829    26. Zetsche B, Gootenberg JS, Abudayyeh OO, Slaymaker IM, Makarova KS,
830        Essletzbichler P, et al. Cpf1 is a single RNA-guided endonuclease of a class 2 CRISPR-
831        Cas system. Cell. 2015;163(3):759-71.

832    27. Esvelt KM, Mali P, Braff JL, Moosburner M, Yaung SJ, Church GM. Orthogonal Cas9
833        proteins for RNA-guided gene regulation and editing. Nature methods.
834        2013;10(11):1116-21.

835    28. Jinek M, Chylinski K, Fonfara I, Hauer M, Doudna JA, Charpentier E. A programmable
836        dual-RNA-guided DNA endonuclease in adaptive bacterial immunity. Science.
837        2012;337(6096):816-21.

838    29. Wang T, Wei JJ, Sabatini DM, Lander ES. Genetic screens in human cells using the
839        CRISPR-Cas9 system. Science. 2014;343(6166):80-4.

840    30. Pastwa E, Blasiak J. Non-homologous DNA end joining. Acta biochimica Polonica.
841        2003;50(4):891-908.

842    31. Aizawa E, Hirabayashi Y, Iwanaga Y, Suzuki K, Sakurai K, Shimoji M, et al. Efficient
843        and accurate homologous recombination in hESCs and hiPSCs using helper-dependent
844        adenoviral vectors. Mol Ther. 2012;20(2):424-31.

845    32. Canver MC, Haeussler M, Bauer DE, Orkin SH, Sanjana NE, Shalem O, et al.
846        Integrated design, execution, and analysis of arrayed and pooled CRISPR genome-
847        editing experiments. Nature protocols. 2018;13(5):946-86.

848    33. Lindsay H, Burger A, Biyong B, Felker A, Hess C, Zaugg J, et al. CrispRVariants charts
849        the mutation spectrum of genome engineering experiments. Nature biotechnology.
850        2016;34(7):701-2.

851   34. Pinello L, Canver MC, Hoban MD, Orkin SH, Kohn DB, Bauer DE, et al. Analyzing
852       CRISPR genome-editing experiments with CRISPResso. Nature biotechnology.
853       2016;34(7):695-7.

854   35. Labun K, Guo X, Chavez A, Church G, Gagnon JA, Valen E. Accurate analysis of
855       genuine CRISPR editing events with ampliCan. Genome research. 2019.

856   36. Harrod A, Fulton J, Nguyen VTM, Periyasamy M, Ramos-Garcia L, Lai CF, et al.
857       Genomic modelling of the ESR1 Y537S mutation for evaluating function and new
858       therapeutic approaches for metastatic breast cancer. Oncogene. 2017;36(16):2286-96.

859   37. Yuan J, Zhang F, Niu R. Multiple regulation pathways and pivotal biological functions of
860       STAT3 in cancer. Scientific reports. 2015;5:17663.

861   38. Rattanapornsompong K, Ngamkham J, Chavalit T, Jitrapakdee S. Generation of Human
862       Pyruvate Carboxylase Knockout Cell Lines Using Retrovirus Expressing Short Hairpin
863       RNA and CRISPR-Cas9 as Models to Study Its Metabolic Role in Cancer Research.
864       Methods in molecular biology (Clifton, NJ). 2019;1916:273-88.

865   39. Narayan M, Peralta DA, Gibson C, Zitnyar A, Jinwal UK. An optimized InCell Western
866       screening technique identifies hexachlorophene as a novel potent TDP43 targeting
867       drug. Journal of biotechnology. 2015;207:34-8.

868   40. Fu Y, Foden JA, Khayter C, Maeder ML, Reyon D, Joung JK, et al. High-frequency off-
869       target mutagenesis induced by CRISPR-Cas nucleases in human cells. Nature
870       biotechnology. 2013;31(9):822-6.

871   41. Merkle FT, Neuhausser WM, Santos D, Valen E, Gagnon JA, Maas K, et al. Efficient
872       CRISPR-Cas9-mediated generation of knockin human pluripotent stem cells lacking
873       undesired mutations at the targeted locus. Cell reports. 2015;11(6):875-83.

874   42. Liu Y, Rao M. Gene targeting in human pluripotent stem cells. Methods in molecular
875        biology (Clifton, NJ). 2011;767:355-67.

876   43. Zwaka TP, Thomson JA. Homologous recombination in human embryonic stem cells.
877        Nature biotechnology. 2003;21(3):319-21.

878   44. Haapaniemi E, Botla S, Persson J, Schmierer B, Taipale J. CRISPR-Cas9 genome
879        editing induces a p53-mediated DNA damage response. Nature medicine.
880        2018;24(7):927-30.

881   45. Ihry RJ, Worringer KA, Salick MR, Frias E, Ho D, Theriault K, et al. p53 inhibits
882        CRISPR-Cas9 engineering in human pluripotent stem cells. Nature medicine.
883        2018;24(7):939-46.

884   46. Tu Z, Yang W, Yan S, Yin A, Gao J, Liu X, et al. Promoting Cas9 degradation reduces
885        mosaic mutations in non-human primate embryos. Scientific reports. 2017;7:42081.

886   47. Santos DP, Kiskinis E, Eggan K, Merkle FT. Comprehensive Protocols for
887        CRISPR/Cas9-based Gene Editing in Human Pluripotent Stem Cells. Current protocols
888        in stem cell biology. 2016;38:5b.6.1-5b.6.60.

889   48. Raj B, Gagnon JA, Schier AF. Large-scale reconstruction of cell lineages using single-
890        cell readout of transcriptomes and CRISPR-Cas9 barcodes by scGESTALT. Nature
891        protocols. 2018;13(11):2685-713.

892   49. Eiges R, Schuldiner M, Drukker M, Yanuka O, Itskovitz-Eldor J, Benvenisty N.
893        Establishment of human embryonic stem cell-transfected clones carrying a marker for
894        undifferentiated cells. Curr Biol. 2001;11(7):514-8.

895   50. Savic N, Ringnalda FC, Lindsay H, Berk C, Bargsten K, Li Y, et al. Covalent linkage of
896        the DNA repair template to the CRISPR-Cas9 nuclease enhances homology-directed
897        repair. eLife. 2018;7.

898    51.  Thorvaldsdottir H, Robinson JT, Mesirov JP. Integrative Genomics Viewer (IGV): high-
899          performance genomics data visualization and exploration. Brief Bioinform.
900          2013;14(2):178-92.

901    52.  Goodwin S, McPherson JD, McCombie WR. Coming of age: ten years of next-
902          generation sequencing technologies. Nat Rev Genet. 2016;17(6):333-51.

903    53.  Frayling TM, Timpson NJ, Weedon MN, Zeggini E, Freathy RM, Lindgren CM, et al. A
904          common variant in the FTO gene is associated with body mass index and predisposes
905          to childhood and adult obesity. Science. 2007;316(5826):889-94.

906    54.  Scuteri A, Sanna S, Chen WM, Uda M, Albai G, Strait J, et al. Genome-wide association
907          scan shows genetic variants in the FTO gene are associated with obesity-related traits.
908          PLoS genetics. 2007;3(7):e115.

909    55.  Yeo GS, Heisler LK. Unraveling the brain regulation of appetite: lessons from genetics.
910          Nature neuroscience. 2012;15(10):1343-9.

911    56.  Gagnon JA, Valen E, Thyme SB, Huang P, Akhmetova L, Pauli A, et al. Efficient
912          mutagenesis by Cas9 protein-mediated oligonucleotide insertion and large-scale
913          assessment of single-guide RNAs. PLoS One. 2014;9(5):e98186.

914

**A** CRISPR targeting of *STAT3*

CRISPR

*STAT3* locus

cloning by FACS, plate duplication

test for loss of STAT3 protein expression

PCR amplification of targeted region

barcode amplicons, pool, and deeply sequence

CCAGATAGGTACTCAGTTAAG
CCAGAT---------AGTTAAG

test for burden of deleterious mutations

**B**

whole-cell stain (CellTag 700)

STAT3 antibody

merge

**C** Read depth

Sample

filtered reads

low quality reads

primer dimer reads

low abundance reads

**D** Allele fraction by consequence class

Sample

high impact (1.0)

low impact (0.7)

medium impact (0.9)

WT (0.0)

low freq. allele (excl.)

**E** Protein KO score

0 — 1

**F** Gene KO score

0 — 1

**G** $R^2 = 0.79$

protein KO score

gene KO score

**H** Integrated KO score

0 — 1