

Whole-genome epigenetic function annotation through sgRNA libraries synthesized by controlled template-dependent elongation

Junling Jia (✉ junling.jia@alumni.bcm.edu)

The First Affiliated Hospital, Zhejiang University

Chen Pan

ZheJiang Univeristy

Ran Li

life sciences institute zhejiang university

Liyang Shui

Collaborative Innovation Center for Diagnosis and Treatment of Infectious Diseases, State Key Laboratory for Diagnosis and Treatment of Infectious Diseases, The First Affiliated Hospital

Yali Wang

life sciences institute zhejiang university

Zhengyun Xiao

zhejiang university

Jing Zhu

Mingtian Genetics

Chao Wu

Zhe Jiang Univeristy <https://orcid.org/0000-0002-6193-4398>

Min Zheng

ZheJiang Univeristy

Article

Keywords: Genetic screen, CRISPR/Cas9, sgRNA library, H3K4me3, CTCF, mESCs, Liver cancer, LincRNA, CDC42

Posted Date: November 5th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-100657/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Whole-genome epigenetic function annotation through**
2 **sgRNA libraries synthesized by controlled template-**
3 **dependent elongation**

4

5 Chen Pan^{1,2,#}, Ran Li^{1,2,#}, Liyan Shui^{2,#}, Yali Wang¹, Zhengyun Xiao¹, Jing Zhu³,

6 Chao Wu^{2,*}, Min Zheng^{2,*} and Junling Jia^{2,4,*}

7

8 1. Life Sciences Institute, Zhejiang University, Hangzhou, Zhejiang, 310058, PRC

9 2. Collaborative Innovation Center for Diagnosis and Treatment of Infectious Diseases,
10 State Key Laboratory for Diagnosis and Treatment of Infectious Diseases, The First
11 Affiliated Hospital, Zhejiang University, Hangzhou, Zhejiang 310003, PRC Zhejiang
12 University, Hangzhou, Zhejiang 310003, PRC

13 3. Department of Anesthesiology & Center for Shock, Trauma and Anesthesiology
14 Research, University of Maryland School of Medicine, Baltimore, MD, 21201, USA

15 4. Innovation Center for Precision Medicine, Zhongtong-Lanbo Diagnostic LTD,
16 Beijing, 100166, PRC

17 *. Correspondence: minzheng@zju.edu.cn (M.Z.), wuchao1984@zju.edu.cn (C.W.),
18 Junling.jia@alumni.bcm.edu (J.J.)

19 # These authors contributed equally

20

21 Keywords: Genetic screen, CRISPR/Cas9, sgRNA library, H3K4me3, CTCF, mESCs,

22 Liver cancer, LincRNA, CDC42

23

24

25

26 **Abstract**

27

28 Epigenome is the set of DNA-associated proteins or chemical modifications to DNA,
29 which regulates gene expression in processes of development and disease. While
30 current advances have allowed researchers to routinely profile epigenomes from given
31 samples, our understandings of the functions of epigenetic hallmarks are nonspecific at
32 best. Applying CRISPR-screening to genome-widely interrogate the function of
33 individual epigenetic hallmarks demands massive sgRNA libraries which are
34 unaffordable via commercial syntheses. Our development consists of a high throughput
35 and cost-effective controlled template-dependent elongation (CTDE) approach which
36 converts source DNA to sgRNA templates. Affiliated screenings encompass 3.8M
37 sgRNAs generated by CTDE targeting all major H3K4me3 and CTCF hallmarks in
38 mESCs and HepG2 and identified 20K essential epigenetic hallmarks, which render the
39 first batch of functional epigenome annotation of H3K4me3 and CTCF hallmarks in
40 mammals. As an application example, we show that a H3K4me3 hallmark orchestrates
41 CDC42 level and cell-cycle progression through promoting LINC00339 expression in
42 HepG2.

43

44 **Introduction**

45 Humans have over 20,000 protein-coding genes, which account for about 2% of
46 the overall genome DNA¹. The expression of those genes are modulated by elements
47 in the remaining portions of the genome, which are regarded as the regulome - an
48 assortment of chromatin components, such as promoters, transcriptional regulatory
49 regions, and high dimensional chromatin structures^{1,2}. Being a major component of
50 regulome, epigenome contains the set of DNA-associated proteins or chemical
51 modifications to DNA³. Systematic discovery and location annotation of epigenome are
52 accomplished primarily by Consortiums such as ENCODE and Roadmap Epigenomics
53 Mapping through assays such as DNase hypersensitivity assays, DNA methylation
54 assays and chromatin immunoprecipitation sequencing (ChIP-seq) assays⁴⁻⁸. These
55 refined techniques have allowed most labs to routinely profile their interested
56 epigenomic information under deliberate experimental conditions.

57 The function of many epigenetic modifications is generally known. For example,
58 H3K4me3 (tri-methylation at the 4th lysine residue of the DNA packaging protein
59 Histone H3) is involved in the positive regulation of the nearby gene transcription^{9,10}.
60 H3K4me3 plays a significant role regarding the regulation of stem cell lineage
61 potency¹¹ as well. CTCF hallmarks primarily have a CCCTC-like motif and are bound
62 by CTCF protein, a 11-zinc finger protein¹². Moreover, CTCF hallmarks are involved
63 in various cellular processes such as transcriptional regulation, insulator activity, V(D)J
64 recombination and the regulation of chromatin architecture¹³⁻¹⁷.

65 Mouse embryonic stem cells (mESCs) have around 42K (23.8Mb) H3K4me3
66 hallmarks and 37K (10.6Mb) CTCF hallmarks. Human liver cancer cells (HepG2) have
67 around 41K (46.9Mb) H3K4me3 hallmarks and 28K (5.2 Mb) CTCF hallmarks.
68 However, there is still a scarcity of a functional epigenome annotation of these
69 hallmarks considering genetic screens have prioritized protein-coding genes or
70 expressed non-coding loci^{18,19}.

71 Numerous reports have depicted pooled CRISPR-screenings that interrogate the
72 function of regulatory elements using designed dense tiling sgRNA libraries targeting
73 limited genomic loci^{8,20,21}. While these studies have provided proof of concept for the
74 application of pooled CRISPR-screening in the functional characterization of
75 regulatory elements, the high-cost of synthesizing a dense tiling sgRNA library to cover
76 an epigenetic hallmark genome-widely (around half-million USD for H3K4me3 or
77 CTCF) is unrealistic and a major hurdle to further functional epigenome studies.

78 We have developed a simple- and cost-effective controlled template-dependent
79 elongation (CTDE) approach that can convert any DNA sample to a sgRNAs library,
80 which covers 98.47% of the effective CRISPR/Cas9 targeting sites within the source
81 DNA. Significantly, over 99% CTDE-sgRNAs targeting sequences have a protospacer
82 adjacent motif (PAM)²²⁻²⁴. We have generated sgRNA libraries targeting all H3K4me3
83 and CTCF hallmarks in mESCs and HepG2. In total, we have screened 3.8 M sgRNAs
84 and identified 14K (14265) H3K4me3 and 6K (6235) CTCF essential hallmarks for the
85 proliferation of mESCs and HepG2. mESCs CTCF dataset shows that mESCs

86 maintains a high proportion of non-essential cell-type specific CTCF hallmarks, which
87 may be important for the implement of pluripotency. Importantly, the HepG2
88 H3K4me3 dataset helps confirm that an essential H3K4me3 hallmark inside the intron
89 of LINC00339 orchestrates the cell-cycle progression and the expression of CDC42, a
90 pivotal factor for the proliferation and invasion of cancer cells, through promoting
91 LINC00339 expression.

92 Our studies have brought us to develop an efficient and budget-friendly approach
93 (CTDE) to convert DNA to a sgRNA library. Several breakthroughs have ensued, one
94 of them being the first batch of functional epigenome annotation of H3K4me3 and
95 CTCF hallmarks in mammals through CTDE library coupled CRISPR-screening.
96 Various other findings through characterizing an essential H3K4me3 hallmark in
97 HepG2 have shown the significance of functional epigenome annotation in cancer
98 research.

99

100

101

102

103 **Result**

104

105 **Convert DNA to a sgRNA library through synthesized by controlled template-**
106 **dependent elongation (CTDE)**

107 The CRISPR/Cas9 system has been developed into genome mutating tools with
108 wide-ranging applications²⁵. The single guide RNA (sgRNA) binds to the Cas9 enzyme
109 and guides the complex to a target via pairing the complementary DNA sequences
110 followed by a protospacer adjacent motif (PAM), where Cas9 performs its
111 endonuclease activity immediately^{22,25}. Currently, sgRNA templates are designed as
112 17-22bp DNA sequences and commercially synthesized *in vitro* as ready-to-clone
113 fragments²⁵⁻²⁷.

114 To directly generate large-scale sgRNA templates from source DNA, we
115 fragmented the DNA template within a 1kb length (Figure 1a). We then ligated the
116 DNA fragments with A1 adaptors (red), immobilizing them on streptavidin beads and
117 washing away their positive strands under a denaturing condition (Figure 1a). Next, we
118 annealed the priming primer (the positive strand of A1 adaptor) onto the immobilized
119 minus strand, then extending the primer using DNA polymerase coupled with reversible
120 terminator (RT) nucleotides (3'-O-N₃-dNTP) that allow singular nucleotide
121 incorporation before the restoration of their 3'-hydroxy groups (Figure 1a)²⁸. We
122 restored the 3'-hydroxyl group via tris (2-carboxyethyl) phosphine (TCEP) treatment
123 and repeated another round of nucleotide incorporation following the previous
124 extension (Figure 1a)²⁸. After 23 rounds of cycling, we blunt the 3' terminus to get 23bp
125 DNA fragments (not including the adapter) (Figure 1a). The most used Cas9 from

126 *Streptococcus pyogenes* recognizes the 5'-NGG-3' (where "N" can be any nucleotide
127 base) PAM sequence. If the last two nucleotides of 3'end of the 23bp DNA are GG, it
128 will compose an AscI cutting site after the A2 adaptor ligation, which is used to select
129 the DNA fragment with 5'-NGG-3' PAM (Figure 1a). Following PAM selection, we
130 remove the NGG triplet using a type II restrict endonuclease (BbsI) and put an A4
131 adaptor onto the 3' terminus for further Gibson assembly into a sgRNA expressing
132 vector (Figure 1a)²⁹. Since only two rare-cutting endonucleases are employed, the
133 dropout rate of sgRNA template caused by endonuclease cutting is very low (0.86%
134 per mouse genome and 0.82% per human genome).

135 To test the efficiency of this technique, we used a 14.9 kb plasmid (lentiCRISPR-
136 v2) modeling a DNA template. After implementing the CTDE steps described above
137 (Figure 1a and S1a), we generated sgRNAs targeting 98.47% of the sites with 5'-NGG-
138 3' PAM sequence (Figure 1b and S1b). As expected, few sgRNAs can be generated
139 from AT rich region because of the low complimentary binding affinity between AT
140 rich sequences (Figure S1b). We also checked the capability of CTDE to enrichment
141 the DNA fragments with 3' NGG triplet. There are 6.23% input DNA fragments that
142 are adjacent to a 3' NGG triplet, and after enrichment, the rate is raised to 99.7% (15.98
143 folds) (Figure 1c). As expected, the lengths of the sgRNA templates are predominantly
144 20bp (83.93%) and the functional sgRNA templates (17-22bp) are 98.66% (Figure
145 1d)^{26,27}.

146 An evaluation of the faithful of the CTDE library ensued the CTDE library was
147 compared with the popular sgRNA libraries (GeCKO 2.0 human/mouse) generated by
148 Chip Based DNA Synthesis (CBDS). Typically, 96.4% sgRNAs of CTDE library can
149 perfectly match a position in their template, and 1.6% sgRNAs of CTDE library carry
150 one mismatch (Figure 1e). The error rate of the CTDE procedure is around 0.97 bases
151 per 1000 bases. Around 90% sgRNAs of CBDS can perfectly match their targeting
152 positions and around 7.5%/9.5% sgRNAs carry one mismatch (Figure 1e). The CBDS
153 library, having an error rate of around 4.5-5.7 bases per 1000 bases, assimilates itself
154 with the CTDE library. We also compared the amplification bias between CTDE library
155 and CBDS library. The sgRNA abundance of CBDS library and CTDE library are both
156 similarly low (Figure 1f). Treating of the abovesaid data, CTDE, consistent in the
157 conversion of source DNA to sgRNA library, produces the same qualitative results as
158 the library generated via Chip Based DNA Synthesis (CBDS).

159 **Validations of CTDE sgRNA library screening**

160 To implement mega-level screenings, we planned to infect cells with the lenti-viral
161 CTDE library expressing sgRNAs along with Cas9 protein (around 40K sgRNAs per
162 batch) (Figure 2a). Then, we sequence the abundance of each sgRNA template at two
163 time points (3 days selected in puromycin media as P1 and 20 days expanded after
164 puromycin selection as P10) to calculate the abundance change from P1 to P10 of each
165 sgRNA template (Figure 2a; details in method).

166 Current bioinformatics efforts for the analysis of pooled CRISPR screens are
167 devoted to identifying genes rather than non-coding genomic loci³⁰. A coding gene is
168 targeted usually by multiple sgRNAs with similar mutational abilities in traditional
169 CRISPR screens; and the essentiality of genes are evaluated through the integrative
170 analysis abundance change of these sgRNAs in frequented algorithms, such as
171 MAGeCK³¹. For the purpose, however, of screening non-coding regulatory loci inside
172 epigenetic hallmarks using CTDE library, the objective becomes the identification of
173 narrow essential genomic sites, the majority of said can only be efficiently targeted by
174 one sgRNA. Traditional calling algorithms are unreliable in this scenario as it will
175 report numerous false positive significantly changed sgRNA (ssgRNA) (Figure 2b).

176 During the screening, many minor uncertain factors bias the abundance of each
177 sgRNA randomly and cause the abundance change of the sgRNAs following a normal
178 distribution in the scenario of non-selection pressure³². Some sgRNA disruptions will
179 cause a negative selection pressure, bias their abundance in P10 systematically (<20%),
180 and result in their abundance change following another normal distribution. While the
181 distance of above two normal distribution is large enough (Figure S2a), we can
182 efficiently identify the sgRNAs (FDR<0.1) that leads to negative selections using a self-
183 developed straight-forward approach (NSgRNAShot; details in method), which has
184 96.42% precision rate and 96.27% recall rate on a simulation dataset (Figure S2a; Table
185 S1). Most importantly, on identical testing datasets, NSgRNAShot has significantly
186 lower false positive rates than MAGeCK (Figure 2b). Additionally, we compared the

187 ability of NSgRNAShot to call true positive ssgRNA with MAGeCK using two
188 published essential gene screen datasets. Due to significant false positive rates,
189 MAGeCK reports a noticeably higher amount of ssgRNAs than NSgRNAShot, most of
190 which target non-essential genes (Figure 2b-c), whereas major portions of ssgRNA
191 identified by NSgRNAShot target essential genes and overlap (100%) with ssgRNAs
192 targeting essential genes reported by MAGeCK (Figure 2c).

193 Following this, through targeting a drug resistance gene (Neo) in mESCs (Figure
194 S2b), we compared the efficiency of CTDE library with the library designed via
195 standard CBDS methods (details in method). The CBDS library has 115 sgRNAs
196 targeting all possible sites (NGG PAM) in Neo gene, along with 20 non-targeting
197 control gRNAs. The CDTE library was generated from the abovesaid Neo gene
198 fragments with additions of the same 20 non-targeting control gRNAs (Figure S2b).
199 The mESCs (expressing Neo) were infected with above lenti-viral libraries and were
200 cultured in medium with Neomycin, and the screening was performed according to the
201 abovesaid procedure. As expected, the abundance of Neo targeting sgRNAs of both
202 libraries exhibit significant decrease and can be identified as ssgRNAs (Figure 2d).
203 Most importantly, the ssgRNAs from both libraries are highly overlapped (97%)
204 (Figure 2d), exhibiting implications that sgRNA library generated by CTDE performs
205 similarly, if not better as the library produced by standard method such as CBDS. Taken
206 together, multiple validations have shown that our CTDE sgRNA library screening
207 approach is both practicable and efficient.

208

209 Annotation of the essential H3K4me3 hallmarks for mESCs self-renewal

210 Tri-methylation of Histone H3 at Lysine 4 (H3K4me3) hallmarks functions in
211 transcriptional activation and maintaining bivalent chromatin of mESCs¹¹. Previous
212 small-scale CRISPR/Cas9-based screens of *Pou5f1* promoter in human ES cells show
213 that numerous critical regulatory elements exist inside H3K4me3 hallmarks and that
214 they are of great important regarding our understanding of gene regulation mechanisms
215 under a biological context^{21,33}.

216 To genome-widely interrogate essential regulatory elements within H3K4me3
217 hallmarks during mESCs self-renewal, we first acquired the H3K4me3 labeled DNA
218 fragments through chromatin immunoprecipitation (ChIP) and converted them to
219 sgRNA libraries through CTDE (Figure 1a). To implement the screening, we infected
220 the mESCs with the lenti-viral library expressing sgRNAs along with Cas9 protein
221 (around 40K sgRNAs per batch). Then, we performed the screening as described above
222 (Figure 2a).

223 In total, we screened 926K sgRNAs targeting 82.99% of the H3K4me3 enriched
224 regions (Figure 3a-d; Table S2-3). In H3K4me3 highly enriched regions (top 100), the
225 sgRNA density is high and reaches to 59 sgRNA per kb. The abundance distribution of
226 sgRNAs inside H3K4me3 hallmarks observes the pattern of their template DNAs
227 (Figure 3a-b and S3a). We identified 24189 sgRNAs causing significant negative
228 selection (called ssgRNA from here), which indicates that their targeting sites in

229 H3K4me3 hallmarks are essential regulatory elements for mESCs self-renewal (Figure
230 3a-d; Table S4). We verified three randomly picked ssgRNAs in non-coding regions.
231 As expected, all three ssgRNAs can significantly inhibit the proliferation of mESCs
232 (Figure 3e and S3b-c).

233 H3K4me3 ssgRNAs distribute evenly on most chromosomes, regions on
234 chromosome 6, 9 and X being an exception (Figure 3a). Given these biases are not the
235 results of lacking matrix H3K4me3 hallmarks or sgRNAs (Figure 3a), we reason that
236 these regions contain fewer essential genes regulated by H3K4me3 hallmarks.
237 ssgRNAs appear in regions, arrayed from weak to strong H3K4me3 elements, the
238 majority of which are among strong elements regions (Figure 3b). Among all essential
239 H3K4me3 elements, 63.24% of the elements are targeted by 1 ssgRNA, and the
240 remaining 37.76% are targeted by multiple ssgRNAs (Table S3). Detailed positions of
241 ssgRNAs inside their H3K4me3 elements can be found in the Supplemental table 4.
242 Given that the H3K4me3 elements are generally wide, the location of these ssgRNAs
243 should reflect the core regulatory sites within the elements.

244 Many H3K4me3 hallmarks locate on exons (Figure 3c; Table S4). Exon mutations
245 will inactivate their genes and cause stronger phenotypes rather than disrupting
246 H3K4me3's regulatory elements. Consistent with this supposition, the percentage of
247 ssgRNAs on exons is significantly greater than the percentage of their matrix sgRNAs,
248 whereas the percentage of the ssgRNAs on other regions are similar as the percentage
249 of their matrix sgRNAs (Figure 3c). Because the sgRNAs targeting known essential

250 genes should be efficiently identified as ssgRNAs, the exons targeting sgRNAs can be
251 used as spike-ins, providing further validation to the CTDE approach under an authentic
252 mega-library screening condition. In sum, exons of 8165 genes are targeted, while 659
253 known essential genes are included (Figure 3f). As expected, our studies exhibit the
254 ssgRNAs targeting 534 essential genes from the said 659, strongly accrediting the
255 success of our screening (Figure 3f).

256 The enriched Gene Ontology (GO) terms of ssgRNAs targeting exons are related
257 to essential biological processes such as ncRNA metabolic and ribonucleoprotein
258 biogenesis (Figure S3d; Table S5). The corresponding genes of ssgRNAs targeting
259 proximal promoters and UTRs can also be confidently identified (Details in method),
260 and their enriched GO terms are related to the essential biological processes of cell
261 survival as well (Figure S3e-f; Table S5).

262 Thus, we have successfully performed a genome-wide CRISPR-screening to
263 interrogate the essential H3K4me3 regulatory elements for mESCs self-renewal.

264

265 **Annotation of the essential CTCF hallmarks for mESCs self-renewal**

266 Gene expressions are orchestrated by regulatory elements at local, long-range
267 and high-dimensional levels³⁴⁻³⁶. CTCF stabilizes chromosomal architecture and
268 coordinates the genome spatial positioning, which functions as a transcriptional
269 activator or repressor¹⁶. In addition to local level regulatory elements (H3K4me3),

270 interrogating CTCF hallmarks provide an important understanding for gene regulation
271 mechanisms in a biological context.

272 Mouse epigenome has around 55,000-65,000 CTCF hallmarks³⁷. ~50% of them
273 are intergenic and ~35% of them are intragenic³⁷. To genome-widely interrogate
274 essential CTCF hallmarks during mESCs self-renewal, we converted CTCF ChIPed
275 DNA fragments to sgRNA libraries via CTDE (Figure 1a) and performed the CRISPR-
276 screening as described above (Figure 2a). In total, we screened 848K sgRNAs which
277 targets 64.12% of the CTCF hallmarks in mESCs (Figure 4a-c and S4a-b; Table S3 and
278 S6). The sgRNA density is 24 sgRNAs per kb in CTCF strong binding sites (top 100).
279 As CTCF elements display a consistent size, whilst maintaining a moderate diversity
280 of the input DNA amount, that of which is significantly smaller than the input amount
281 of H3K4me3, the abundance distribution of sgRNAs inside CTCF hallmarks is notably
282 more even than that of H3K4me3, although the pattern of their template DNAs is still
283 observed (Figure 4b and S4a). We identified 3038 CTCF ssgRNAs, which indicates
284 that the corresponding CTCF hallmarks (47.02 % in intergenic regions) are essential
285 for mESCs self-renewal (Figure 4a-c and S4b; Table S4). ssgRNAs appear in regions,
286 arrayed from weak to strong CTCF elements. Among all essential CTCF elements,
287 87.07% elements are targeted by 1 ssgRNA whilst 12.93% elements are targeted by
288 multiple ssgRNAs (Table S3). Detailed positions of ssgRNAs inside CTCF elements
289 can be found in the Supplemental table 4. CTCF elements are narrow (83% < 400 bp)
290 and have clear binding motif, and the ssgRNAs are rather important in indication of the

291 essentiality of their belonging CTCF elements than their targeting location inside. We
292 verified three randomly picked intergenic CTCF ssgRNAs. As expected, all can
293 significantly compromise mESCs proliferation (Figure 4d and S4c-d).

294 Because CTCF hallmarks stabilize high-dimensional architecture of chromosome
295 at multiple levels, disrupting CTCF hallmarks will disrupt the topological structure of
296 genomes at different levels and can cause different level cell stresses. Thus, the
297 distribution of CTCF ssgRNA on most chromosomes is not even as the distribution of
298 their matrix sgRNAs (Figure 4a-b).

299 Distal promoters, introns and proximal promoters are major parts on which CTCF
300 hallmarks locate (Figure S4b; Table S4). The GO terms of genes close to ssgRNAs
301 targeting introns and proximal promoters are essential biological processes and tissue
302 developments (Figure 4e and S4e; Table S7). Because the expression of differentiation
303 and development related genes generally antagonizes mESCs pluripotency and self-
304 renewal, we believe that these essential CTCF hallmarks should inhibit their expression.
305 Unlike H3K4me3 ssgRNAs, only a minor part of CTCF ssgRNAs target exons (Figure
306 S4b; Table S4). Although most GO terms of expressed genes of these exons are also
307 related to essential biological processes, their significance is much lower than that of
308 H3K4me3 ssgRNA (Figure S4e and S3d; Table S5 and S7). We reason that major
309 functions of these essential CTCF hallmarks are beyond promoting the expression of
310 their sitting genes.

311 mESCs differentiate into various cell types during embryonic development.
312 Previous studies have shown that the chromosome spatial structure will rearrange
313 accordingly to fit the change of gene expression patterns during differentiation³⁸. We
314 compared the CTCF hallmarks with 16 mouse cell types/tissues and found that 59.63%
315 CTCF hallmarks in mESCs are cell-type specific and 40.37% are common (Figure 4f;
316 Table S8)^{37,39}. The common CTCF hallmarks should help maintain the universal spatial
317 structure of chromosome, while the cell-type specific CTCF hallmarks should be either
318 mESCs specific or pre-loaded hallmarks for further differentiated cells. Consistent with
319 this supposition, the percentage of the cell-type specific essential CTCF hallmarks
320 (28.85%) of mESCs is significantly smaller than the percentage of the cell-type specific
321 CTCF hallmarks (59.63%) (Figure 4f; Table S8).

322

323 **Annotation of the essential H3K4me3 hallmarks in human liver cancer cells**

324 Whole-genome sequencing has surveyed large sets of cancer genomes and studied
325 the role and extent of single-nucleotide variants (SNVs), small insertions/deletions
326 (indels) and larger structural variants in cancers^{40,41}. While the initial focus on the
327 genetic variations in protein-coding regions has dramatically expanded our knowledge
328 of cancer genetics, the remaining (>90%) non-coding part of the genetic variations are
329 much more difficult to understand and have remained largely unexplored⁴², which is
330 due to a lack of functional annotation of regulatory elements inside.

331 To genome-widely interrogate essential activating regulatory elements in human
332 liver cancer cells (HepG2) (Figure S5a-b), we performed a H3K4me3 CRISPR-
333 screening as described above (Figure 1a and 2a). In total, we screened 1.19M sgRNAs
334 targeting 80.91% of the H3K4me3 hallmarks in HepG2 (Figure 5a-c and S5c-d; Table
335 S3 and S9). In H3K4me3 highly enriched regions (top100), the sgRNA density is
336 43sgRNAs per kb. The abundance distribution of sgRNAs inside H3K4me3 hallmarks
337 observes the pattern of their template DNAs (Figure 5b and S5c). We have identified
338 14540 ssgRNAs (75.82% are inside non-coding regions), which represent 6475
339 essential regulatory elements in HepG2 (Figure 5a-c and S5d; Table S4). ssgRNAs
340 appear in regions from weak to strong H3K4me3 elements, and the majority are in
341 strong elements regions (Figure 5b). Among all essential H3K4me3 elements 62.89%
342 elements are targeted by 1 ssgRNA while 37.11% elements are targeted by multiple
343 ssgRNAs (Table S3). Detailed positions of ssgRNAs inside H3K4me3 elements can be
344 found in the Supplemental table 4, indicating the core regulatory sites of these essential
345 H3K4me3 elements. We also verified three randomly picked ssgRNAs targeting non-
346 coding regions, and all can significantly inhibit HepG2 growth (Figure 5d and S5e-f).

347 The H3K4me3 ssgRNAs evenly distribute on most chromosomes, with regions on
348 chromosome 5, 8 and 13 being exceptions, in which there is no shortage of H3K4me3
349 hallmarks and sgRNAs (Figure 5a). This indicates that fewer essential genes exist in
350 these regions.

351 Most H3K4me3 ssgRNAs locate on exons, proximal promoters, and introns
352 (Figure S5d; Table S4). Normally, genes are more efficiently inactivated by mutations
353 on their exons than their regulatory regions. Hence, the percentage of the ssgRNAs on
354 exons is significantly greater than the percentage of the ssgRNAs on other regions
355 (Figure S5d). The enriched GO terms of ssgRNA targeting exons are related to essential
356 biological processes such as ncRNA metabolic and mitochondrial function, which play
357 central roles in malignancy through macromolecular synthesis and energy
358 production^{43,44} (Figure S5g; Table S10), while the enriched GO terms of ssgRNAs
359 targeting proximal promoters and UTRs are also related to the essential biological
360 processes of cell survival (Figure 5e and S5g; Table S10).

361

362 **Annotation of the essential CTCF hallmarks in human liver cancer cells**

363 CTCF hallmarks stabilize mammalian genomes into discrete structural and
364 regulatory domains, those of which can either prevent or facilitate the interactions of
365 promoters and enhancers across their boundaries¹⁵. Although previous works
366 substantiate that CTCF hallmarks can evolve in human cancers⁴⁵, their underlying
367 mechanisms are principally more difficult to define as a result of missing functional
368 annotation of CTCF hallmarks in human.

369 To interrogate essential CTCF hallmarks in HepG2, we generated CTCF sgRNA
370 libraries and performed the CRISPR-screening as described above (Figure 1a and 2a).
371 Overall, we screened 1.06M sgRNAs, targeting 79.75% of the CTCF hallmarks in

372 HepG2 (Figure 6a-c and S6a-b; Table S3 and S11). In strong CTCF binding regions
373 (top100), the sgRNA density is 63sgRNAs per kb. As expected, the abundance
374 distribution of sgRNAs inside CTCF hallmarks is notably more even than that of
375 H3K4me3, and the patterns of their template DNAs are observed all the while (Figure
376 6b and 5b). We identified 4628 CTCF ssgRNAs which represent 3583 (44.63 % inside
377 intergenic regions) essential CTCF hallmarks for HepG2 growth (Figure S6b; Table
378 S4). ssgRNAs appear in regions from weak to strong CTCF elements (Figure 6b).
379 Among all HepG2 essential CTCF elements, 78.76% elements are targeted by 1
380 ssgRNA while 21.24% elements are targeted by multiple ssgRNAs (Table S3). Detailed
381 positions of ssgRNAs inside CTCF elements can be found in the Supplemental table 4.
382 We verified three randomly picked intergenic CTCF ssgRNAs. As expected, all can
383 significantly inhibit HepG2 growth (Figure 6d and S6c-d).

384 CTCF ssgRNAs on chromosomes roughly follow the distribution of their matrix
385 hallmarks and sgRNAs (Figure 6a). As we known that CTCF hallmarks stabilize
386 genome from smaller loops into huge megabase-sized loops called topologically
387 associated domains (TADs)^{46,47}, CTCF ssgRNAs will disrupt the topological structure
388 at different levels and result in their uneven distribution on many regions (Figure 6a).

389 Major CTCF hallmarks exist on distal promoters, introns, and proximal promoters
390 in HepG2 (Figure S6b; Table S4). CTCF ssgRNA targeting introns and proximal
391 promoters have significant GO terms including signaling transduction and cell
392 morphogenesis, all of which are crucial for cancers (Figure 6e and S6e; Table S12).

393 Only a minor portion of CTCF ssgRNAs target exons (Figure S6b). The majority of
394 GO terms of expressed genes whose exons are targeted by ssgRNAs are related to
395 essential biological processes (Figure S6f; Table S12). Because major functions of
396 CTCF hallmarks on exons are beyond promoting their sitting genes expression¹⁵, the
397 GO term significance is lesser than that of H3K4me3 ssgRNA targeting exons (Figure
398 S6f and S5g; Table S12 and S10).

399 The majority of CTCF hallmarks in HepG2 are common (79.79%) among 55
400 human cell types (Figure 6f; Table S13). As a tissue-specific cell line, the spatial
401 chromosome structure of HepG2 has been adapted to the requirements of liver functions,
402 and it is not necessary to keep so many spatial chromosome structures specific for other
403 cell types. Unlike mESCs, the percentage of the cell-type specific essential CTCF
404 hallmarks corresponds with the percentage of the cell-type specific CTCF hallmarks in
405 HepG2 (Figure 6f and 4f).

406

407 **H3K4me3 hallmark-LINC00339-CDC42 axis maintaining HepG2 growth**

408 We believe that functional epigenome annotation can facilitate uncovering novel
409 regulation mechanisms and biomarkers of cancer cells. We focused on a H3K4me3
410 ssgRNA (chr1-22352881), which targets the intron of LINC00339 that is highly
411 expressed in multiple cancer cell lines (Figure 7a). Because H3K4me3 activates local
412 gene expression, we checked if ssgRNA (chr1-22352881) would disrupt LINC00339
413 expression. As predicted, the LINC00339 level is significantly decreased after ssgRNA

414 (chr1-22352881) disruption (Figure 7b). Ensuing, we knocked down the expression of
415 LINC00339, and discovered that the proliferation of HepG2 is significantly inhibited
416 into the same level as ssgRNA (chr1-22352881) disruption (Figure 7c-e). Cell-cycle
417 analysis shows that both ssgRNA (chr1-22352881) disruption and knocking down
418 LINC00339 blocks the S-phase entry of HepG2 but not hESCs (Figure 7f). Altogether,
419 the data suggests that ssgRNA (chr1-22352881) disruption compromises HepG2
420 proliferation through downregulating the expression of LINC00339.

421 It has recently become apparent that long non-coding RNAs (lncRNAs) can
422 function as transcriptional activators⁴⁸. Through binding to histone-modifying
423 complexes, transcription factors and RNA polymerase II, lncRNAs can promote gene
424 expression in *cis* or in *trans*⁴⁸. Located at the immediate downstream of LINC00339,
425 CDC42 functions in cell-cycle and anchorage-independent growth (Figure S5f)^{49,50}.
426 CDC42 also transduces growth and adhesion signals to drives cell-cycle progress from
427 G1 to S phase^{49,50}. Therefore, LINC00339 may promote CDC42 expression in HepG2.
428 As expected, we found that both ssgRNA (chr1-22352881) disruption and knocking-
429 down LINC00339 significantly decreases the expression of CDC42 (Figure 7g-h). Our
430 data establishes that the H3K4me3 hallmark targeted by ssgRNA (chr1-22352881)
431 orchestrates the activity of LINC00339-CDC42 axis to promote HepG2 growth.

432

433 **Discussions**

434 A vast proportion of mammalian genomes are non-coding and contain vital
435 regulatory roles. Unlike coding-regions, functional analyses of non-coding regions
436 through CRISPR-screening requires high density sgRNA coverage. The commercial
437 synthesis of a sgRNA library covering an epigenetic hallmark (H3K4me3 for example:
438 200-300K USD) is unaffordable for most labs. As a result, molecular biology
439 techniques that directly convert the source DNA into sgRNA library are urgently
440 needed to unleash the full capabilities of CRISPR-screening for decoding non-coding-
441 genomes. To this end, several approaches have been developed. One of which, named
442 CORALINA (comprehensive gRNA library generation through controlled nuclease
443 activity), employs MNase to break down the source DNA⁵¹. CORALINA incorporates
444 fragmented DNA into library vector without a step to enrich the fragment flanking a
445 Protospacer Adjacent Motif (PAM)⁵¹. Without PAM, around 70-77% sgRNAs
446 generated by CORALINA are ineffective. Hiroshi Arakawa developed another
447 approach that employs six type II or type III restriction enzymes (EcoP15I, Bgl II, Acu
448 I, XbaI, Bsm BI and AatI) digestions and PAGE-Gel purifications to convert mRNA
449 into sgRNAs with flanking PAM⁵². Because six restriction enzymes digestions will
450 destroy 6.2% potential sgRNAs and multiple PAGE-Gel purification significantly
451 increases DNA loss, Arakawa's approach is not practicable for large-scale CRISPR-
452 screening. Our CTDE approach overcomes pervious limitations by generating sgRNA
453 libraries from source DNAs with a simple- and cost-effective procedure.

454 Because the CTDE sgRNA library is directly from a DNA source, the abundance
455 distribution of sgRNAs follows the abundance distribution of their input DNA.
456 Benefited by this characteristic, our screen can focus more on the highly enriched
457 H3K4me3 and CTCF regions which normally would have more biological significance.
458 However, the biological significance of RNAs do not always correlate with their
459 abundance, which normally varies up to hundreds of folds inside cells. Thus, CTDE
460 sgRNA library is not practicable for whole transcriptome screening.

461 The approach utilizing paired-guide RNA CRISPR/Cas9 library(CREST-seq) has
462 been shown efficient discovery and functional characterization of regulatory elements²¹.
463 This type of approach needs more than three paired-guide RNAs for each hallmark, and
464 the library cost to cover all CTCF and H3K4me3 hallmarks (79K in mESCs and 69K
465 in HepG2) proves to be too costly. In addition, once asynchronous cutting happens, the
466 sequence of the first cutting site will change⁵³ and the hallmark cannot be deleted by
467 the paired-guide RNAs anymore. So, the paired-guide RNA approach needs
468 CRISPR/Cas9 cut both sites simultaneously, which leads to a lower efficiency than
469 single sgRNA system. Thus, the paired-guide RNA library-based screening requires a
470 more sensitive readout than the survival readout used in this work.

471 The essential hallmarks datasets of H3K4me3 and CTCF in mESCs and HepG2
472 can provide potential benefits for multiple fields in the future. For instance, Induced
473 Pluripotent Stem Cells (iPSCs) researchers should evaluate the somatic mutations
474 inside the essential hallmarks of mESCs, which can significantly compromise the

475 reprogramming efficiency. Additionally, the H3K4me3 essential hallmarks related to
476 core transcriptional factors can be utilized to boost the reprogramming efficiency
477 through CRISPR mediated activation⁵⁴. The HepG2 datasets permits cancer researchers
478 to direct their focus on H3K4me3 or/and CTCF essential hallmarks, which are specific
479 in liver cancer cells. Through reversible or permanent inactivation of these hallmarks
480 via CRISPR mediated approaches, potential novel liver cancer treatment strategies can
481 be developed⁵⁵.

482 CTDE provides a simple, time- and cost-effective procedure to convert source
483 DNA to sgRNA library, which could be more broadly applied to functional epigenome
484 annotation in a biological context.

485

486

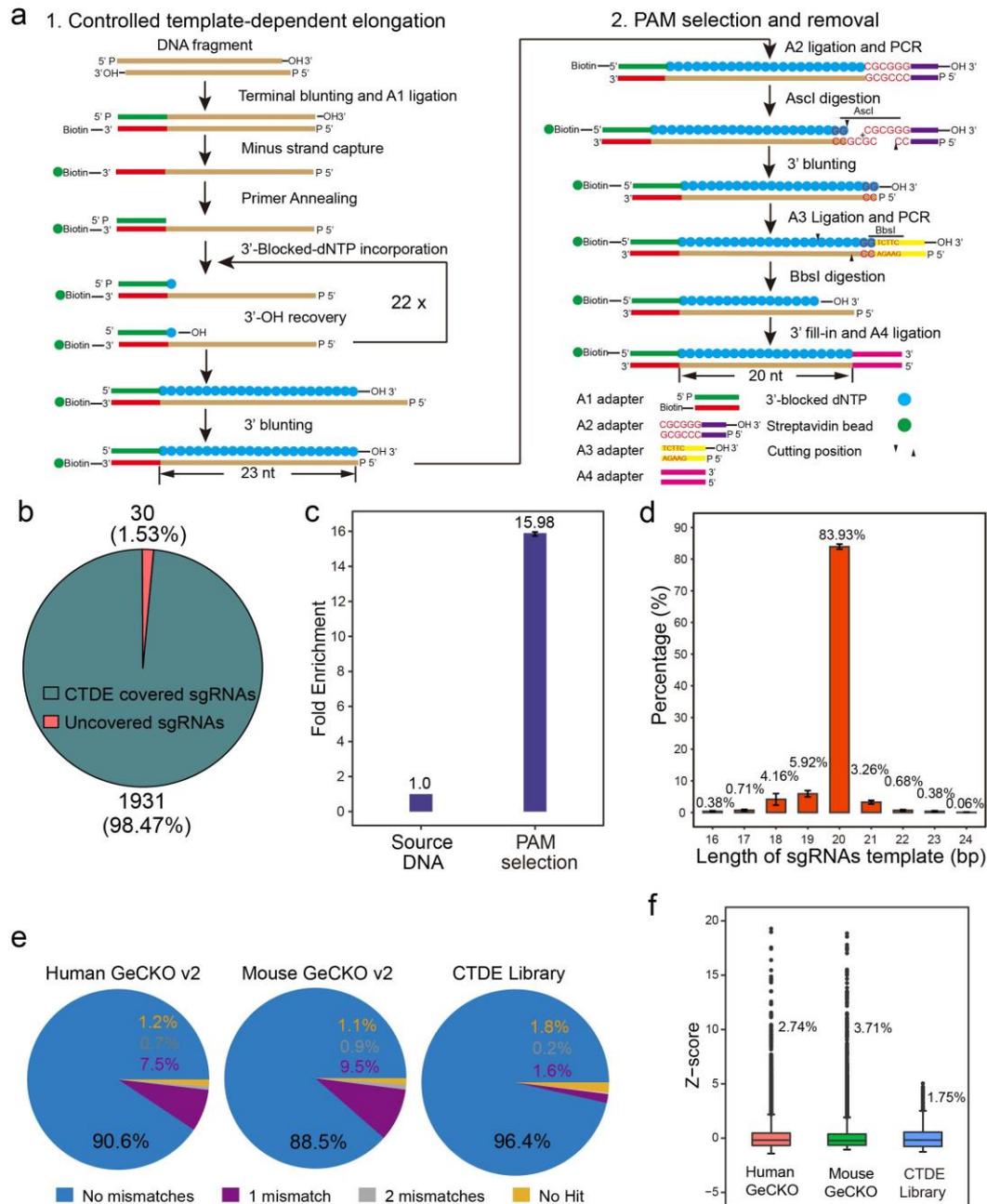
487 **References**

- 488 1. TA., B. Genomes. 2nd edition.
- 489 2. Bonev, B. & Cavalli, G. Organization and function of the 3D genome. *Nat Rev Genet* **17**, 772
490 (2016).
- 491 3. Bernstein, B.E., Meissner, A. & Lander, E.S. The mammalian epigenome. *Cell* **128**, 669-681
492 (2007).
- 493 4. Abbott, A. Project set to map marks on genome. *Nature* **463**, 596-597 (2010).
- 494 5. Consortium, E.P. et al. Identification and analysis of functional elements in 1% of the human
495 genome by the ENCODE pilot project. *Nature* **447**, 799-816 (2007).
- 496 6. Boyle, A.P. et al. High-resolution mapping and characterization of open chromatin across the
497 genome. *Cell* **132**, 311-322 (2008).
- 498 7. Kurdyukov, S. & Bullock, M. DNA Methylation Analysis: Choosing the Right Method. *Biology*
499 (*Basel*) **5** (2016).
- 500 8. Barski, A. et al. High-resolution profiling of histone methylations in the human genome. *Cell*
501 **129**, 823-837 (2007).
- 502 9. Sims, R.J., 3rd, Nishioka, K. & Reinberg, D. Histone lysine methylation: a signature for chromatin
503 function. *Trends Genet* **19**, 629-639 (2003).
- 504 10. Wysocka, J. et al. A PHD finger of NURF couples histone H3 lysine 4 trimethylation with
505 chromatin remodelling. *Nature* **442**, 86-90 (2006).
- 506 11. Bernstein, B.E. et al. A bivalent chromatin structure marks key developmental genes in
507 embryonic stem cells. *Cell* **125**, 315-326 (2006).
- 508 12. Phillips, J.E. & Corces, V.G. CTCF: master weaver of the genome. *Cell* **137**, 1194-1211 (2009).
- 509 13. Chaumeil, J. & Skok, J.A. The role of CTCF in regulating V(D)J recombination. *Curr Opin Immunol*
510 **24**, 153-159 (2012).
- 511 14. Guelen, L. et al. Domain organization of human chromosomes revealed by mapping of nuclear
512 lamina interactions. *Nature* **453**, 948-951 (2008).
- 513 15. Kim, S., Yu, N.K. & Kaang, B.K. CTCF as a multifunctional protein in genome regulation and gene
514 expression. *Exp Mol Med* **47**, e166 (2015).
- 515 16. Ong, C.T. & Corces, V.G. CTCF: an architectural protein bridging genome topology and function.
516 *Nat Rev Genet* **15**, 234-246 (2014).
- 517 17. Khoury, A. et al. Constitutively bound CTCF sites maintain 3D chromatin architecture and long-
518 range epigenetically regulated domains. *Nat Commun* **11**, 54 (2020).
- 519 18. Joung, J. et al. Genome-scale CRISPR-Cas9 knockout and transcriptional activation screening.
520 *Nat Protoc* **12**, 828-863 (2017).
- 521 19. Shalem, O., Sanjana, N.E. & Zhang, F. High-throughput functional genomics using CRISPR-Cas9.
522 *Nat Rev Genet* **16**, 299-311 (2015).
- 523 20. Fulco, C.P. et al. Systematic mapping of functional enhancer-promoter connections with CRISPR
524 interference. *Science* **354**, 769-773 (2016).
- 525 21. Diao, Y. et al. A tiling-deletion-based genetic screen for cis-regulatory element identification in
526 mammalian cells. *Nat Methods* **14**, 629-635 (2017).
- 527 22. Shah, S.A., Erdmann, S., Mojica, F.J. & Garrett, R.A. Protospacer recognition motifs: mixed
528 identities and functional diversity. *RNA Biol* **10**, 891-899 (2013).

- 529 23. Jinek, M. et al. A programmable dual-RNA-guided DNA endonuclease in adaptive bacterial
530 immunity. *Science* **337**, 816-821 (2012).
- 531 24. Sternberg, S.H., Redding, S., Jinek, M., Greene, E.C. & Doudna, J.A. DNA interrogation by the
532 CRISPR RNA-guided endonuclease Cas9. *Nature* **507**, 62-67 (2014).
- 533 25. Hsu, P.D., Lander, E.S. & Zhang, F. Development and applications of CRISPR-Cas9 for genome
534 engineering. *Cell* **157**, 1262-1278 (2014).
- 535 26. Fu, Y., Sander, J.D., Reyon, D., Cascio, V.M. & Joung, J.K. Improving CRISPR-Cas nuclease
536 specificity using truncated guide RNAs. *Nat Biotechnol* **32**, 279-284 (2014).
- 537 27. Ran, F.A. et al. Double nicking by RNA-guided CRISPR Cas9 for enhanced genome editing
538 specificity. *Cell* **154**, 1380-1389 (2013).
- 539 28. Metzker, M.L. et al. Termination of DNA synthesis by novel 3'-modified-deoxyribonucleoside
540 5'-triphosphates. *Nucleic Acids Res* **22**, 4259-4267 (1994).
- 541 29. Gibson, D.G. et al. Enzymatic assembly of DNA molecules up to several hundred kilobases. *Nat*
542 *Methods* **6**, 343-345 (2009).
- 543 30. Bodapati, S., Daley, T.P., Lin, X., Zou, J. & Qi, L.S. A benchmark of algorithms for the analysis of
544 pooled CRISPR screens. *Genome Biol* **21**, 62 (2020).
- 545 31. Li, W. et al. MAGeCK enables robust identification of essential genes from genome-scale
546 CRISPR/Cas9 knockout screens. *Genome Biol* **15**, 554 (2014).
- 547 32. Bracewell, R. The Fourier Transform & Its Applications 3rd Edition. *Book*.
- 548 33. Diao, Y. et al. A new class of temporarily phenotypic enhancers identified by CRISPR/Cas9-
549 mediated genetic screening. *Genome Res* **26**, 397-405 (2016).
- 550 34. Weake, V.M. & Workman, J.L. Inducible gene expression: diverse regulatory mechanisms. *Nat*
551 *Rev Genet* **11**, 426-437 (2010).
- 552 35. Malik, S. & Roeder, R.G. The metazoan Mediator co-activator complex as an integrative hub for
553 transcriptional regulation. *Nat Rev Genet* **11**, 761-772 (2010).
- 554 36. Ong, C.T. & Corces, V.G. Enhancer function: new insights into the regulation of tissue-specific
555 gene expression. *Nat Rev Genet* **12**, 283-293 (2011).
- 556 37. Kim, T.H. et al. Analysis of the vertebrate insulator protein CTCF-binding sites in the human
557 genome. *Cell* **128**, 1231-1245 (2007).
- 558 38. Dixon, J.R. et al. Chromatin architecture reorganization during stem cell differentiation. *Nature*
559 **518**, 331-336 (2015).
- 560 39. Cuddapah, S. et al. Global analysis of the insulator binding protein CTCF in chromatin barrier
561 regions reveals demarcation of active and repressive domains. *Genome Res* **19**, 24-32 (2009).
- 562 40. Blum, A., Wang, P. & Zenklusen, J.C. SnapShot: TCGA-Analyzed Tumors. *Cell* **173**, 530 (2018).
- 563 41. Consortium, I.T.P.-C.A.o.W.G. Pan-cancer analysis of whole genomes. *Nature* **578**, 82-93 (2020).
- 564 42. Rheinbay, E. et al. Analyses of non-coding somatic drivers in 2,658 cancer whole genomes.
565 *Nature* **578**, 102-111 (2020).
- 566 43. Zong, W.X., Rabinowitz, J.D. & White, E. Mitochondria and Cancer. *Mol Cell* **61**, 667-676 (2016).
- 567 44. Anastasiadou, E., Jacob, L.S. & Slack, F.J. Non-coding RNA networks in cancer. *Nat Rev Cancer*
568 **18**, 5-18 (2018).
- 569 45. Song, S.H. & Kim, T.Y. CTCF, Cohesin, and Chromatin in Human Cancer. *Genomics Inform* **15**,
570 114-122 (2017).

- 571 46. Yu, M. & Ren, B. The Three-Dimensional Organization of Mammalian Genomes. *Annu Rev Cell*
572 *Dev Biol* **33**, 265-289 (2017).
- 573 47. Pombo, A. & Dillon, N. Three-dimensional genome architecture: players and mechanisms. *Nat*
574 *Rev Mol Cell Biol* **16**, 245-257 (2015).
- 575 48. Long, Y., Wang, X., Youmans, D.T. & Cech, T.R. How do lncRNAs regulate transcription? *Sci Adv*
576 **3**, eaao2110 (2017).
- 577 49. Chou, M.M., Masuda-Robens, J.M. & Gupta, M.L. Cdc42 promotes G1 progression through p70
578 S6 kinase-mediated induction of cyclin E expression. *J Biol Chem* **278**, 35241-35247 (2003).
- 579 50. Qadir, M.I., Parveen, A. & Ali, M. Cdc42: Role in Cancer Management. *Chem Biol Drug Des* **86**,
580 432-439 (2015).
- 581 51. Kofler, A. et al. CORALINA: a universal method for the generation of gRNA libraries for CRISPR-
582 based screening. *BMC Genomics* **17**, 917 (2016).
- 583 52. Arakawa, H. A method to convert mRNA into a gRNA library for CRISPR/Cas9 editing of any
584 organism. *Sci Adv* **2**, e1600699 (2016).
- 585 53. Sander, J.D. & Joung, J.K. CRISPR-Cas systems for editing, regulating and targeting genomes.
586 *Nat Biotechnol* **32**, 347-355 (2014).
- 587 54. Konermann, S. et al. Genome-scale transcriptional activation by an engineered CRISPR-Cas9
588 complex. *Nature* **517**, 583-588 (2015).
- 589 55. Qi, L.S. et al. Repurposing CRISPR as an RNA-guided platform for sequence-specific control of
590 gene expression. *Cell* **152**, 1173-1183 (2013).
- 591
592
593

Figure 1



594

595 **Figure 1 Controlled template-dependent elongation (CTDE) converts source DNA to a sgRNA**
 596 **library**

597 (a) Detailed schematics of the CTDE procedure. Briefly, sgRNA templates (23bp) are generated
 598 from fragmented input DNA by controlled template-dependent elongation; sgRNA templates with
 599 NGG triplet in 3' terminus are enriched and NGG triplets are removed before sgRNA templates are
 600 cloned into a vector.

601 (b) The coverage rate of the potential sgRNA templates on source DNA (lentiCRISPR-V2) by CTDE.

602 (c) Fold enrichment of the sgRNA templates after NGG-PAM selection step. Error bars SD of
 603 triplicates.

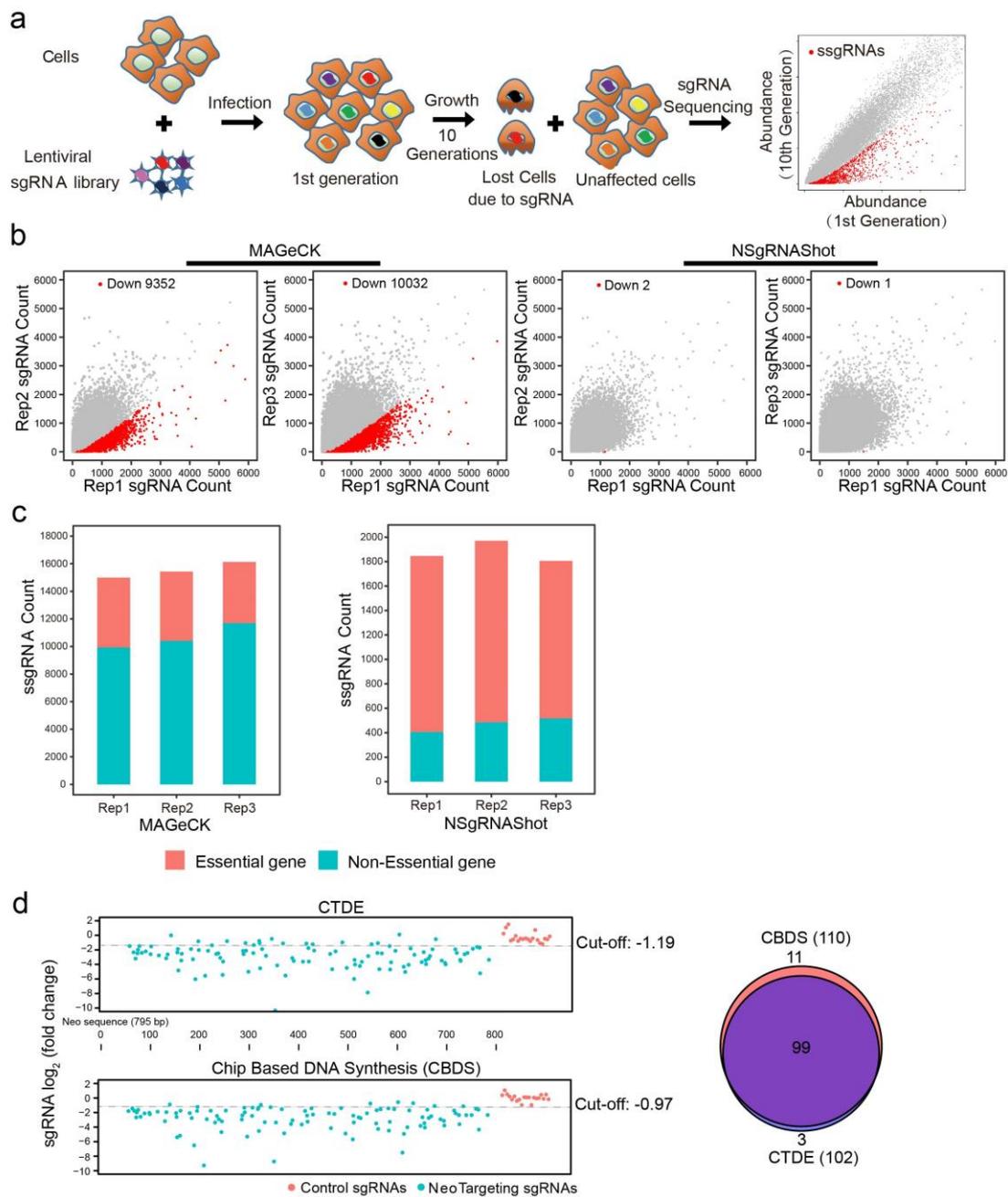
604 (d) Length distribution of the sgRNA templates from source DNA (lentiCRISPR-V2). Error bars
605 SD of triplicates.

606 (e) The fidelity of sgRNAs generated by CBDS (GeCKO libraries) and CTDE approaches; the
607 mapping quality of the sgRNAs against their targeting sites is used to show the library fidelity.

608 (f) The abundance bias of sgRNAs of CBDS (GeCKO libraries) and CTDE approaches; distribution
609 of sgRNA counts z-score of two GeCKO libraries (red and green) and CTDE library (blue) are
610 displayed.

611

Figure 2



612

613

614 **Figure 2 Validations of CTDE sgRNA library screening**

615 (a) The schematic of CRISPR/Cas9 dropout screen: Cells are infected by library virus at low MOI;
 616 infected cells are cultured 10 generations; by comparing the abundance of sgRNA templates
 617 between the 1st and 10th generation, the sgRNAs which limit cell proliferation are identified by
 618 NSgRNAShot as significant sgRNAs(ssgRNAs).

619 (b) Comparison of false-positive rate of significantly changed sgRNA (ssgRNA) calling between
 620 MAGeCK and NSgRNAShot. Datasets of three biological replicates from an CRISPR dropout
 621 screen project (Tzelepis et al., 2016) are used (details in method). Scatter-plots exhibit sgRNA
 622 abundance at same time point between two biological replicate datasets. MAGeCK identifies

623 noticeably more false-positive ssgRNAs (red dots); NSgRNAShot identifies few false-positive
624 ssgRNAs.

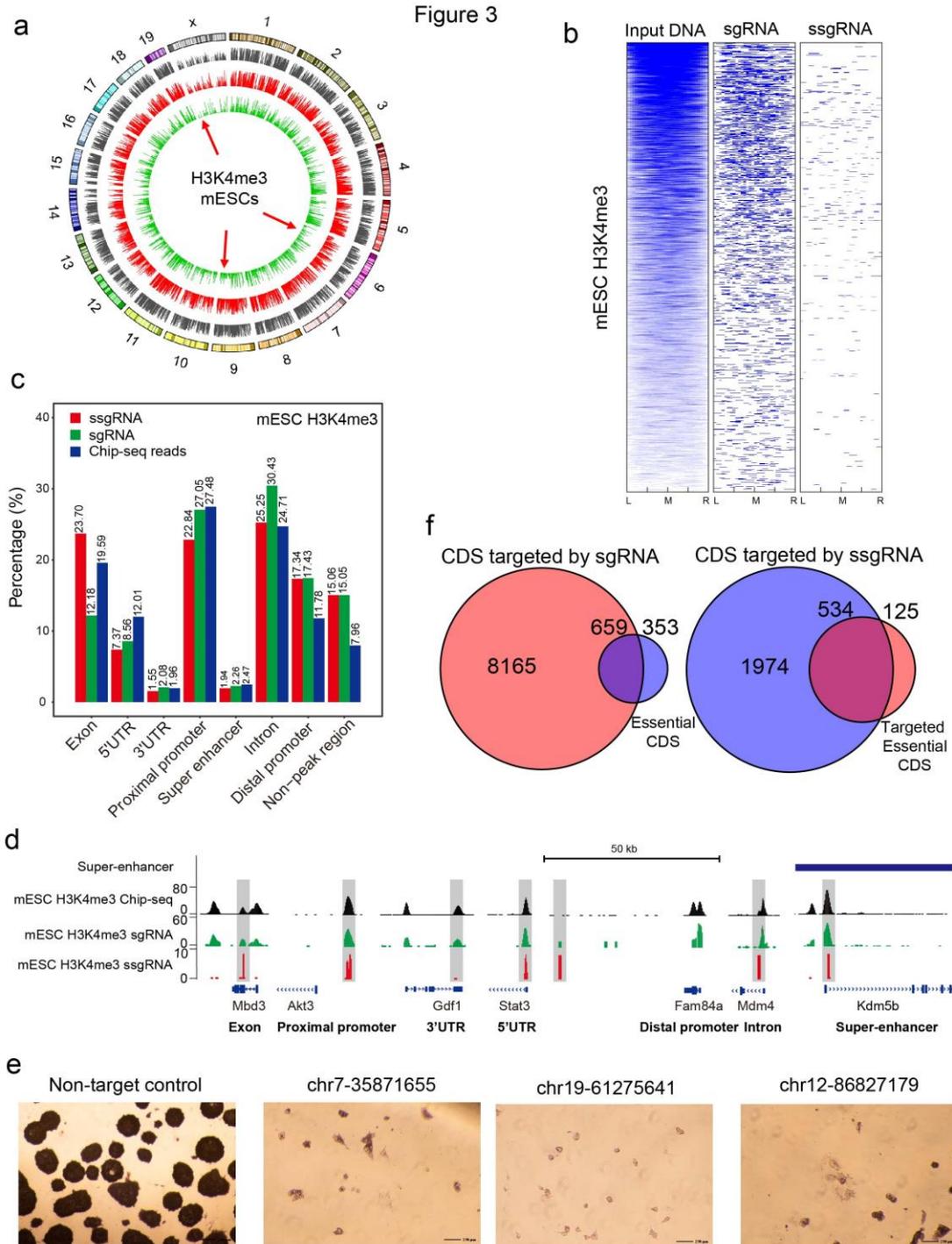
625 (c) Bar plots display the number and distribution (essential genes and non-essential genes) of
626 ssgRNAs identified by MAGeCK and NSgRNAShot from above (b) CRISPR dropout screen
627 project.

628 (d) Left panel: \log_2 fold change of the sgRNA abundance of Neo gene targeting libraries generated
629 by CBDS and CTDE during validation screening (details in method). Red dots are non-targeting
630 control sgRNAs; green dots are Neo targeting sgRNAs. Right panel: the overlap of significantly
631 changed sgRNA (ssgRNA) from CBDS library and CTDE library.

632

633

634

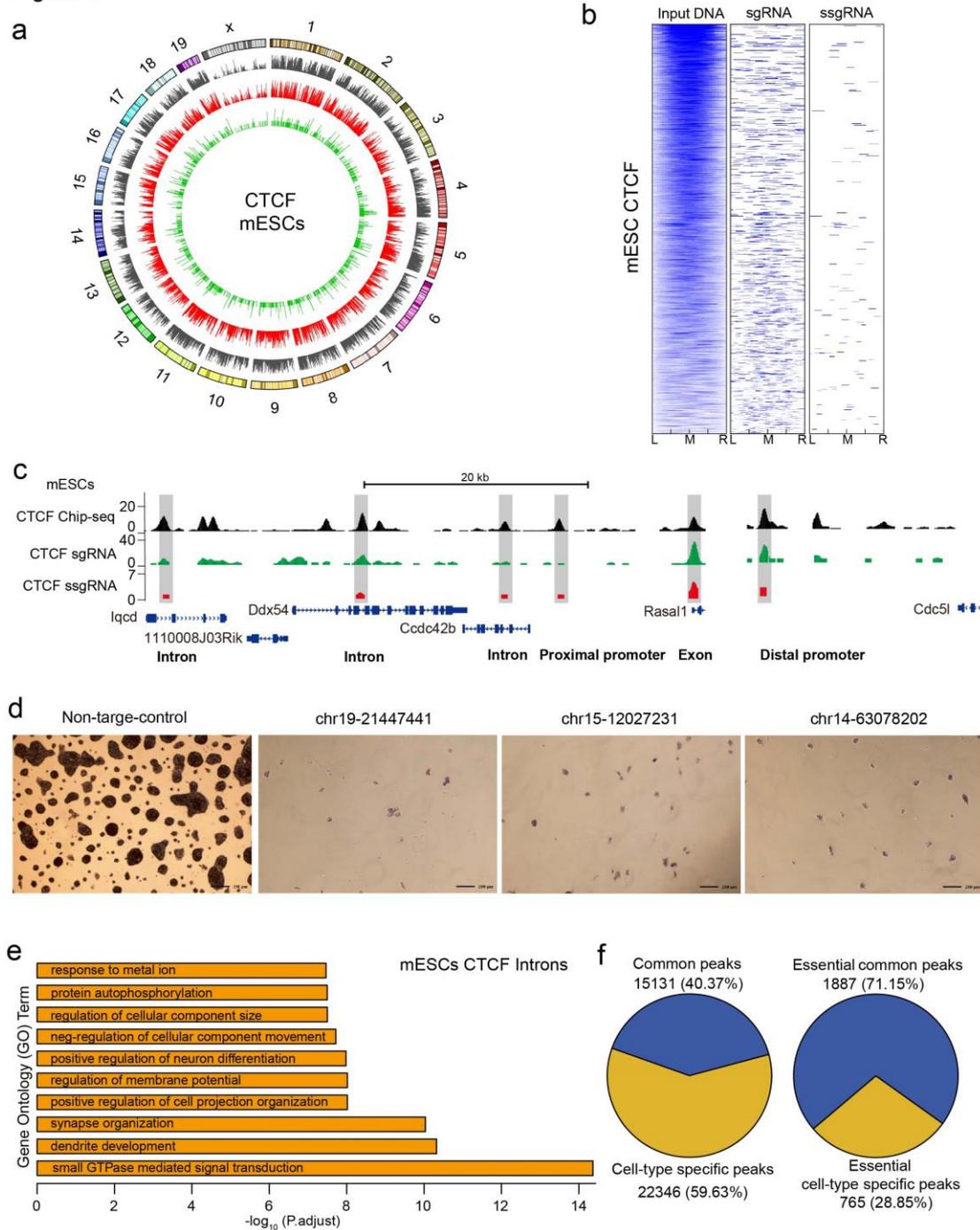


635
636

637 Figure 3 Annotation of the essential H3K4me3 hallmarks in mESCs self-renewal
638 (a) The circos plot of H3K4me3 ChIP-seq, H3K4me3 sgRNAs and H3K4me3 ssgRNAs profiles on
639 mouse genome. The circles, from exterior to interior, illustrate mouse genome (mm9), ChIP-Seq
640 reads density (black), H3K4me3 sgRNAs density (red) and H3K4me3 ssgRNAs density (green).
641 Arrows show regions with biased ssgRNAs distribution.
642 (b) Heat-map of input DNA, sgRNA and ssgRNA enrichment inside mESCs H3K4me3 elements.
643 All H3K4me3 elements are displayed in relative scale. L: Left boundary; M: Middle; R: Right
644 boundary. The heat-maps are rank-ordered according to the enrichment of input DNA (blue,

645 enriched; white, not enrichment).
646 (c) The distribution of H3K4me3 ChIP-seq reads, H3K4me3 sgRNAs and H3K4me3 ssgRNAs on
647 major regulatory regions of mouse genome.
648 (d) Plots of H3K4me3 ChIP-seq reads, H3K4me3 sgRNAs and H3K4me3 ssgRNAs at seven
649 genomic loci. Y-axes, RPKM. Genomic regions with enriched H3K4me3 ssgRNA are shaded grey.
650 (e) Validation of three H3K4me3 ssgRNA in mESCs self-renewal. Alkaline phosphatase (AP)
651 staining illustrates that three randomly selected H3K4me3 ssgRNAs significantly inhibit mESCs
652 self-renewal. ssgRNA is named by three factors (chromosome number; targeting strain + or -;
653 mapping position). Scale bar, 250 μ m.
654 (f) Left panel: CDS targeted by mESCs H3K4me3 sgRNAs includes 659 well-recognized essential
655 genes (Shohat and Shifman, 2019; Tzelepis et al., 2016). Right panel: CDS targeted by mESCs
656 H3K4me3 ssgRNAs includes 534 genes from above 659 essential genes.
657
658

Figure 4



659
660

661 **Figure 4 Annotation of the essential CTCF hallmarks in mESCs self-renewal**

662 (a) The circos plot of CTCF ChIP-seq, CTCF sgRNAs and CTCF ssgRNAs profiles on mouse
663 genome. The circles, from exterior to interior, illustrate mouse genome (mm9), CTCF ChIP-Seq
664 reads density (black), CTCF sgRNAs density (red) and CTCF ssgRNAs density (green).

665 (b) Heat-map of input DNA, sgRNA and ssgRNA enrichment inside mESCs CTCF elements. All
666 CTCF elements are displayed in relative scale. L: Left boundary; M: Middle; R: Right boundary.
667 The heat-maps are rank-ordered based on the enrichment of input DNA as follows: blue, enriched;
668 white, not enrichment.

669 (c) Plots of CTCF ChIP-seq reads, CTCF sgRNAs and CTCF ssgRNAs at six genomic loci. Y-axes,
670 RPKM. Genomic regions with CTCF ssgRNA are shaded grey.

671 (d) Validation of three CTCF ssgRNA in mESCs self-renewal. Alkaline phosphatase (AP) staining
672 shows that three randomly selected CTCF ssgRNAs significantly inhibit mESCs self-renewal.
673 ssgRNA is named by three factors (chromosome number; targeting strain + or -; mapping position).

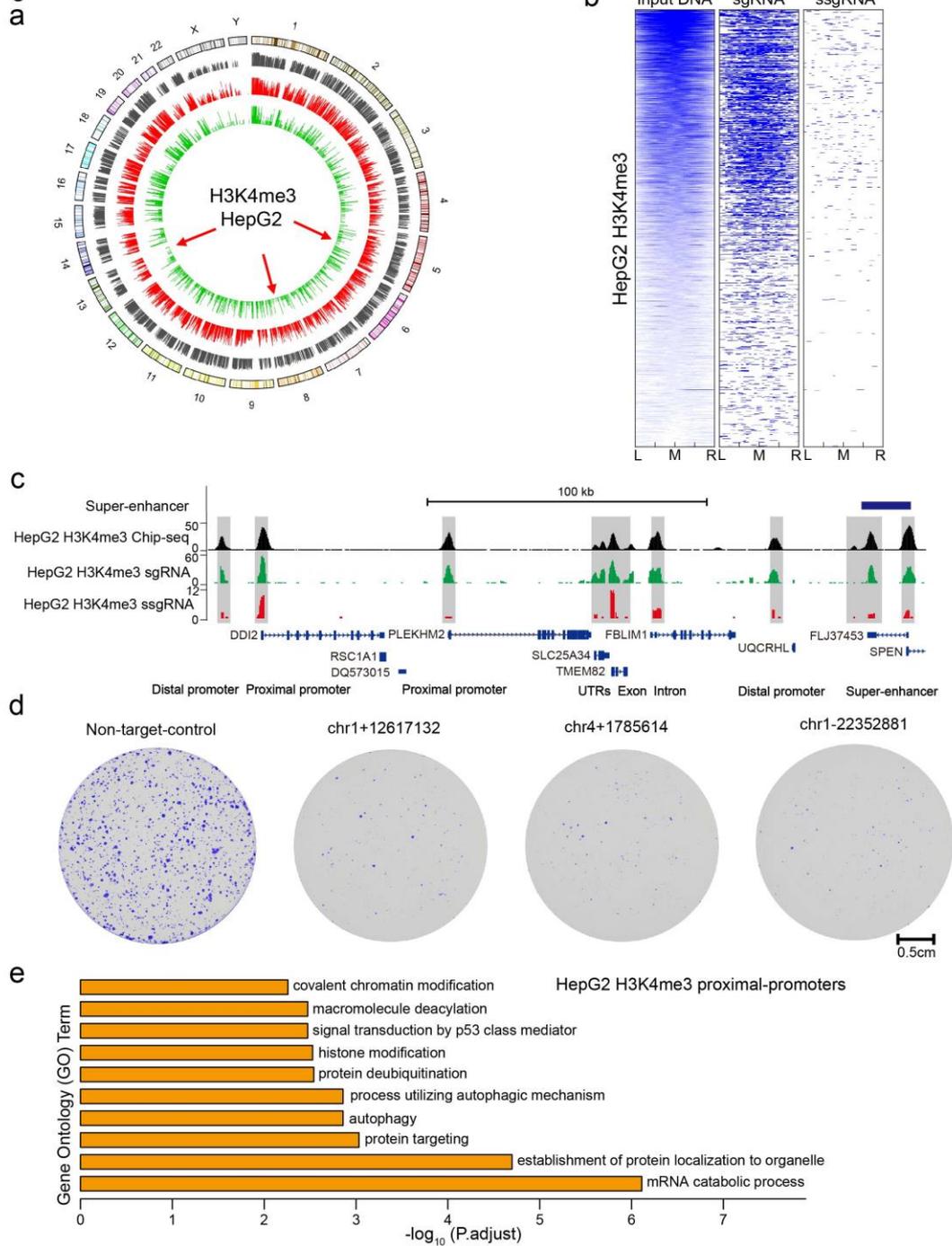
674 Scale bar, 250 μ m.

675 (e) Top 10 GO terms of the genes whose introns are targeted by CTCF ssgRNAs.

676 (f) Cell-type specific analysis of essential CTCF hallmark (targeted by ssgRNA) in mESCs self-
677 renewal. Major part (71.15%) of essential CTCF hallmark in mESCs self-renewal are common VS
678 major part (59.63%) of CTCF hallmarks in mESCs are cell-type specific.

679

Figure 5



680

681

682 **Figure 5 Annotation of the essential H3K4me3 hallmarks in HepG2**

683 (a) The circos plot of H3K4me3 ChIP-seq, H3K4me3 sgRNAs and H3K4me3 ssgRNAs profiles on
 684 human genome. The circles, from exterior to interior, illustrate human genome (hg19), ChIP-Seq
 685 reads density (black), H3K4me3 sgRNAs density (red) and H3K4me3 ssgRNAs density (green).
 686 Arrows show regions with biased ssgRNAs distribution.

687 (b) Heat-map of input DNA, sgRNA and ssgRNA enrichment inside HepG2 H3K4me3 elements.
 688 All H3K4me3 elements are displayed in relative scale. L: Left boundary; M: Middle; R: Right
 689 boundary. The heat-maps are rank-ordered based on the enrichment of input DNA as follows: blue,

690 enriched; white, not enrichment.

691 (c) Plots of H3K4me3 ChIP-seq reads, H3K4me3 sgRNAs and H3K4me3 ssgRNAs at eight
692 genomic loci. Y-axes, RPKM. Genomic regions with enriched H3K4me3 ssgRNA are shaded grey.

693 (d) Validation of three H3K4me3 ssgRNAs in HepG2. Crystal violet staining shows that three
694 randomly selected H3K4me3 ssgRNAs significantly inhibit HepG2 proliferation. ssgRNA is named
695 by three factors (chromosome number; targeting strain + or -; mapping position). Scale bar, 500 μ m

696 (e) Top 10 GO terms of the genes whose proximal promoters are targeted by H3K4me3 ssgRNAs.

697

698

699

700

Figure 6

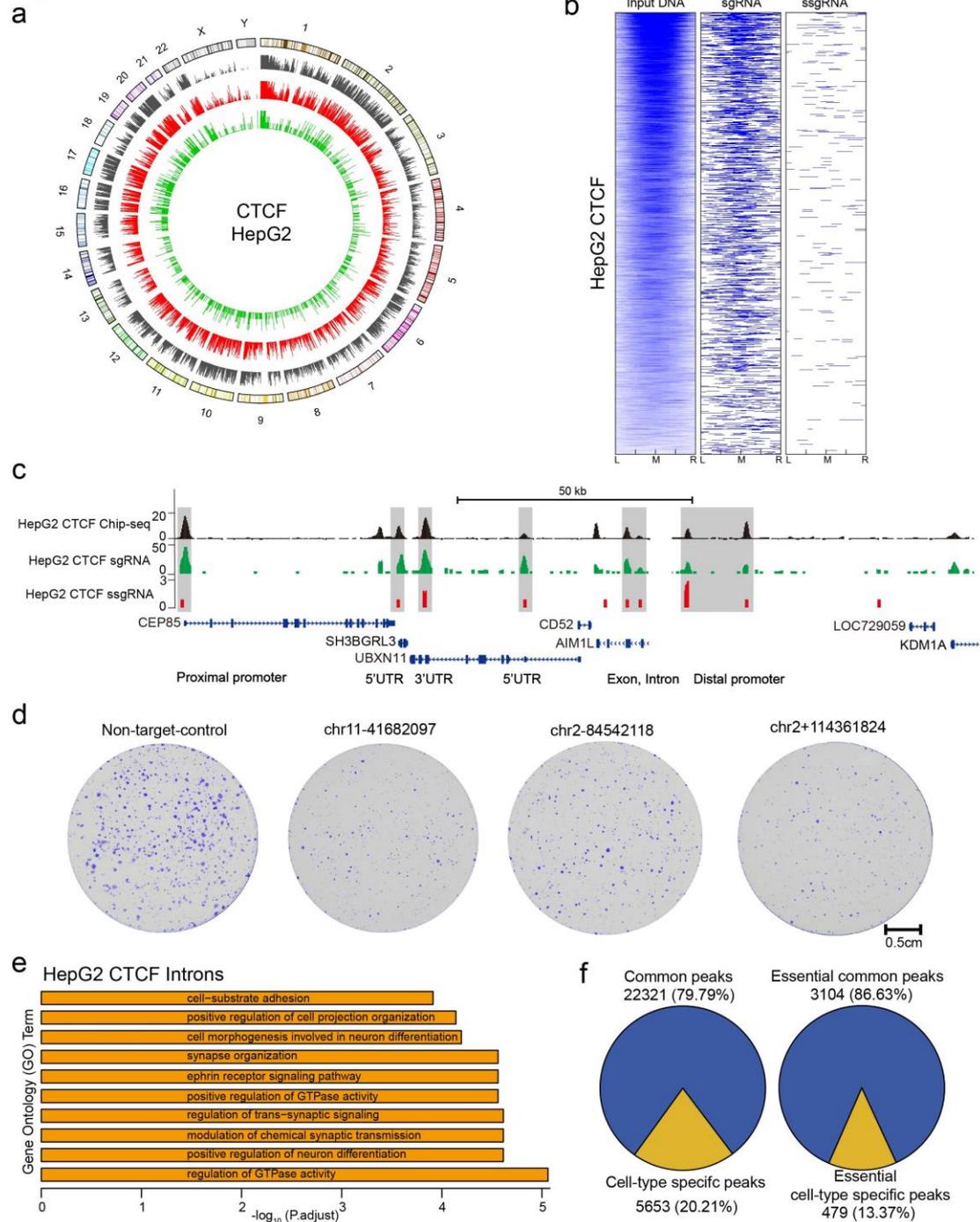


Figure 6 Annotation of the essential CTCF hallmarks in HepG2

(a) The circos plot of CTCF ChIP-seq, CTCF sgRNAs and CTCF ssgRNAs profiles on human genome. The circles, from outside to inside, illustrate human genome (hg19), CTCF ChIP-Seq reads density (black), CTCF sgRNAs density (red) and CTCF ssgRNAs density (green).

(b) Heat-map of input DNA, sgRNA and ssgRNA enrichment inside HepG2 CTCF elements. All CTCF elements are displayed in relative scale. L: Left boundary; M: Middle; R: Right boundary. The heat-maps are rank-ordered based on the enrichment of input DNA as follows: blue, enriched; white, not enrichment.

701
702

711 (c) Plots of CTCF ChIP-seq reads, CTCF sgRNAs and CTCF ssgRNAs at seven genomic loci. Y-
712 axes, RPKM. Genomic regions with CTCF ssgRNA are shaded grey.

713 (d) Validation of three CTCF ssgRNAs in HepG2 proliferation. Crystal violet staining illustrates
714 that three randomly selected CTCF ssgRNAs significantly inhibit HepG2 proliferation. ssgRNA is
715 named by three factors (chromosome number; targeting strain + or -; mapping position). Scale bar,
716 500 μ m

717 (e) Top 10 GO terms of the genes whose intron are targeted by CTCF ssgRNAs.

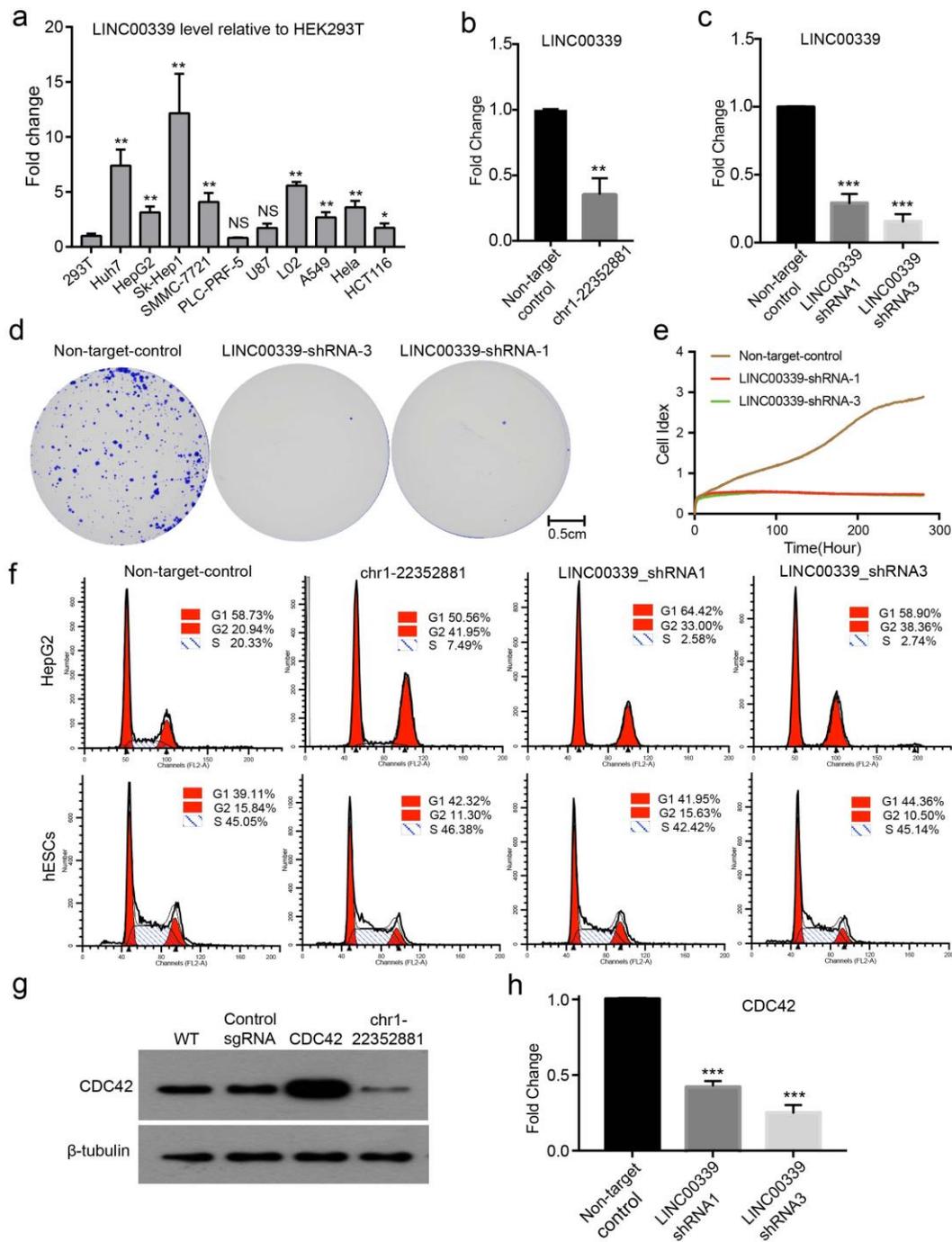
718 (f) Cell-type specific analysis of essential CTCF hallmark (targeted by ssgRNA) in HepG2
719 proliferation. Major part of both CTCF hallmarks (79.79%) and essential CTCF hallmarks (86.63%)
720 are common.

721

722

723

Figure 7



724
725

726 **Figure 7 H3K4me3 hallmark-LINC00339-cdc42 axis maintaining HepG2 growth**

727 (a) LINC00339 expression in ten human cancer cell lines by RT-PCR. Most human cancer cell lines
728 have more LINC00339 than HEK293T. Student's T-test, *** $p < 0.001$, ** $p < 0.01$. Error bars SD of
729 triplicates.

730 (b) Quantification of LINC00339 RNA level by RT-PCR after ssgRNA (chr1-22352881) disruption.
731 ssgRNA (chr1-22352881) disruption significantly decreases LINC00339 level in HepG2. Student's
732 T-test, *** $p < 0.01$. Error bars SD of triplicates.

733 (c) Quantification of LINC00339 RNA level by RT-PCR after shRNA knockdown. Two shRNAs

734 can significantly decrease LINC00339 level in HepG2. Student's T-test, *** $p < 0.001$. Error bars SD
735 of triplicates.

736 (d) Proliferation assay of HepG2 after knockdown LINC00339. Knockdown LINC00339
737 significantly inhibits HepG2 proliferation. 14 days culture after LINC00339 knockdown, cells were
738 stained by crystal violet.

739 (e) Knockdown LINC00339 significantly inhibits HepG2 proliferation. Real-time cell numbers are
740 plotted.

741 (f) Flow cytometry analysis of HepG2 and hESCs after ssgRNA (chr1-22352881) disruptio and
742 LINC00339 knockdown. ssgRNA (chr1-22352881) disruption and LINC00339 knockdown only
743 block S-phase entry in HepG2.

744 (g) Western blot analysis of CDC42 after ssgRNA (chr1-22352881) disruption in HepG2. CDC42
745 over-expression and a scramble sgRNA as controls. Loading control, tubulin.

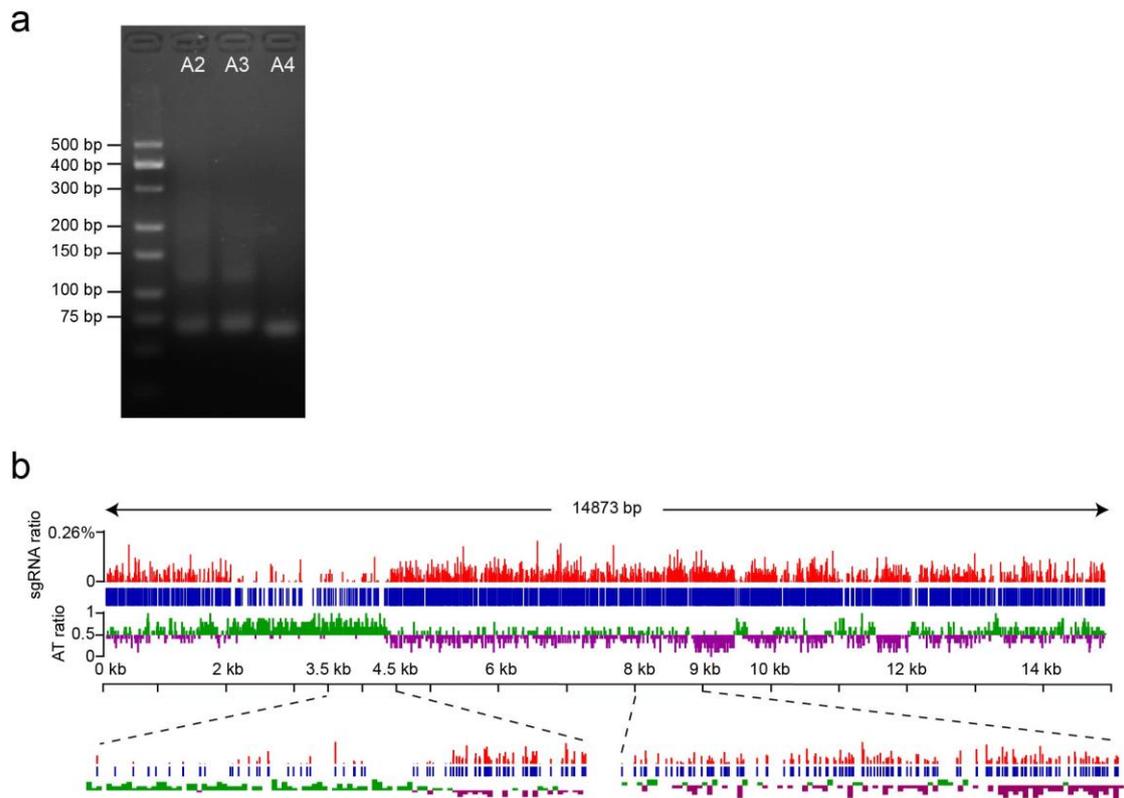
746 (h) Quantification of CDC42 mRNA by RT-PCR after LINC00339 knockdown. Student's T-test,
747 *** $p < 0.001$. Error bars SD of triplicates.

748

749

750

Figure S1

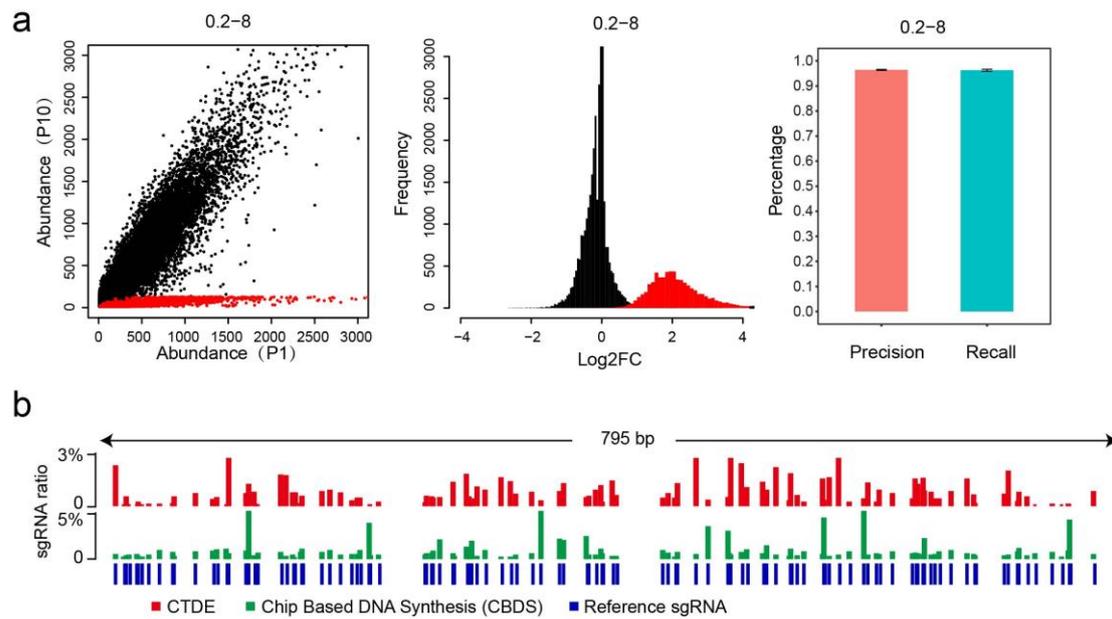


751
752
753
754
755
756
757
758
759
760
761

Figure S1

(a) Agarose gel (3.5%) image of PCR product after adapter (A2, A3 and A4) ligation.
 (b) The abundance distribution of sgRNA templates generated by CTDE on source DNA (LentiCRISPR-V2). Up panel: Red displays the relative abundance of sgRNA templates on each PAM (NGG) site. Middle: blue represents is potential sgRNA sites (NGG PAM) on source DNA. Down panel: The green and purple are the relative abundance of A/T in 10bp bins (details in method). A A/T-rich region and a balanced nucleotide region are zoomed in.

Figure S2



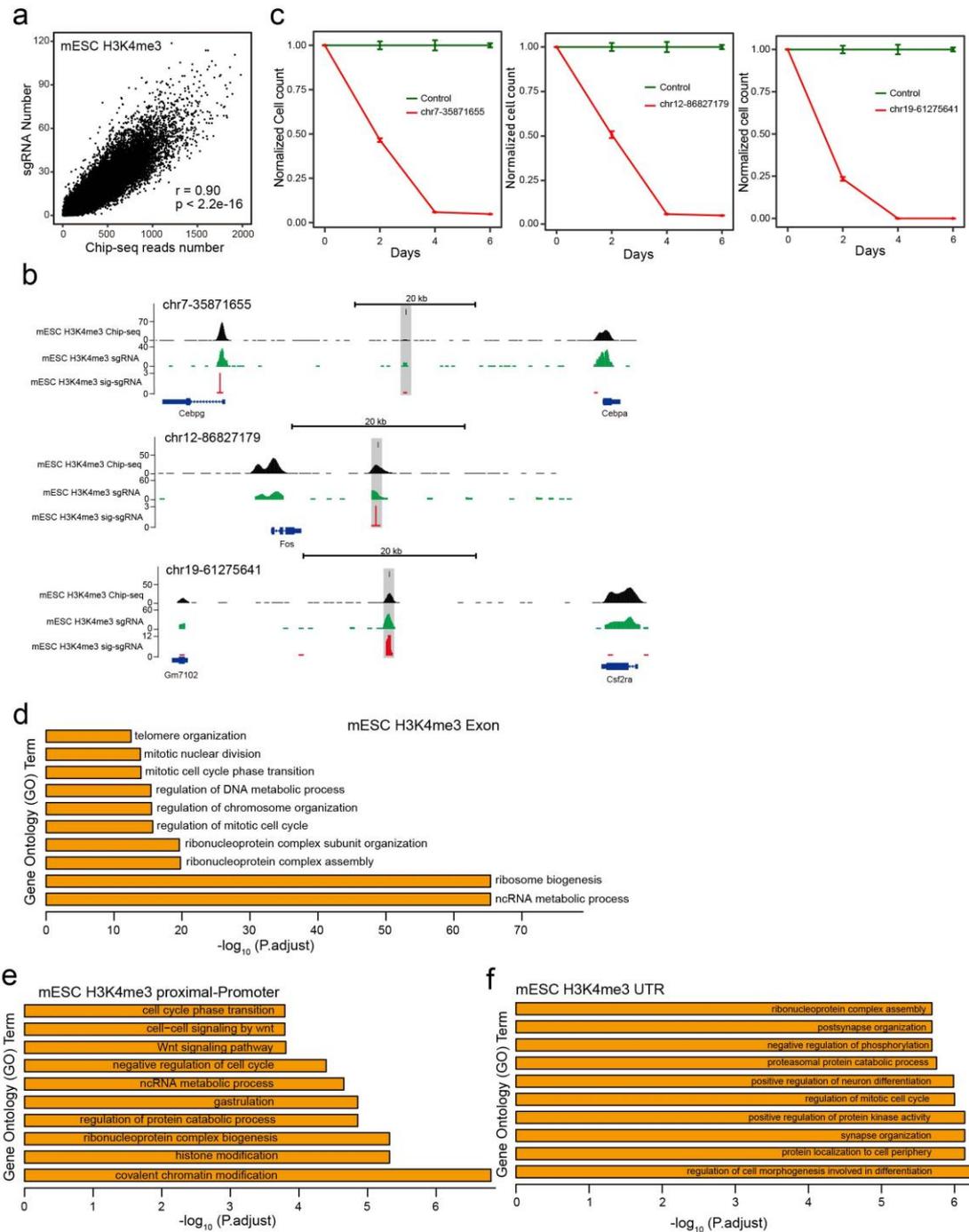
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776

Figure S2

(a) NSgRNAShot can efficiently identify ssgRNA from a typical ($\alpha=0.2$; $\beta=8$; details in method) simulated dataset. Left picture: Scatter-plot of normalized read counts of sgRNAs at P1 and P10 of the simulated dataset. Middle picture: the abundance fold-change (FC) distribution of the sgRNAs. Right picture: the precision rate and the recall rate of NSgRNAShot in simulated dataset. The black represents the sgRNAs without suffering dropout. The red represents the sgRNAs suffering dropout.

(b) The abundance distribution of sgRNA templates generated by CTDE and CBDS on Neo gene. Red displays the relative abundance of sgRNA templates (CTDE) on each PAM (NGG) site; green is the relative abundance of sgRNA templates (CBDS) on each PAM (NGG) site; blue shows potential sgRNA sites (NGG PAM) on Neo gene.

Figure S3



777

778 **Figure S3**

779 (a) Scatter-plot of the correlation between ChIP-seq reads number and sgRNA number in mESCs
 780 H3K4me3 elements. Spearson’s correlation coefficient with two-tailed test was calculated.

781 (b) Plots of H3K4me3 ChIP-seq reads, H3K4me3 sgRNAs and H3K4me3 ssgRNAs at three
 782 genomic loci. Y-axes, RPKM. The genomic regions with H3K4me3 ssgRNAs validated in Figure

783 3e are shaded grey.

784 (c) Three H3K4me3 ssgRNAs validated in Figure 3e significantly inhibit mESCs proliferation. Cell
785 numbers of three time points are plotted. Error bars the SD of triplicates

786 (d) Top 10 GO terms of the genes whose exons are targeted by H3K4me3 ssgRNAs.

787 (e) Top 10 GO terms of the genes whose proximal promoters are targeted by H3K4me3 ssgRNAs.

788 (f) Top 10 GO terms of the genes whose UTRs are targeted by H3K4me3 ssgRNAs.

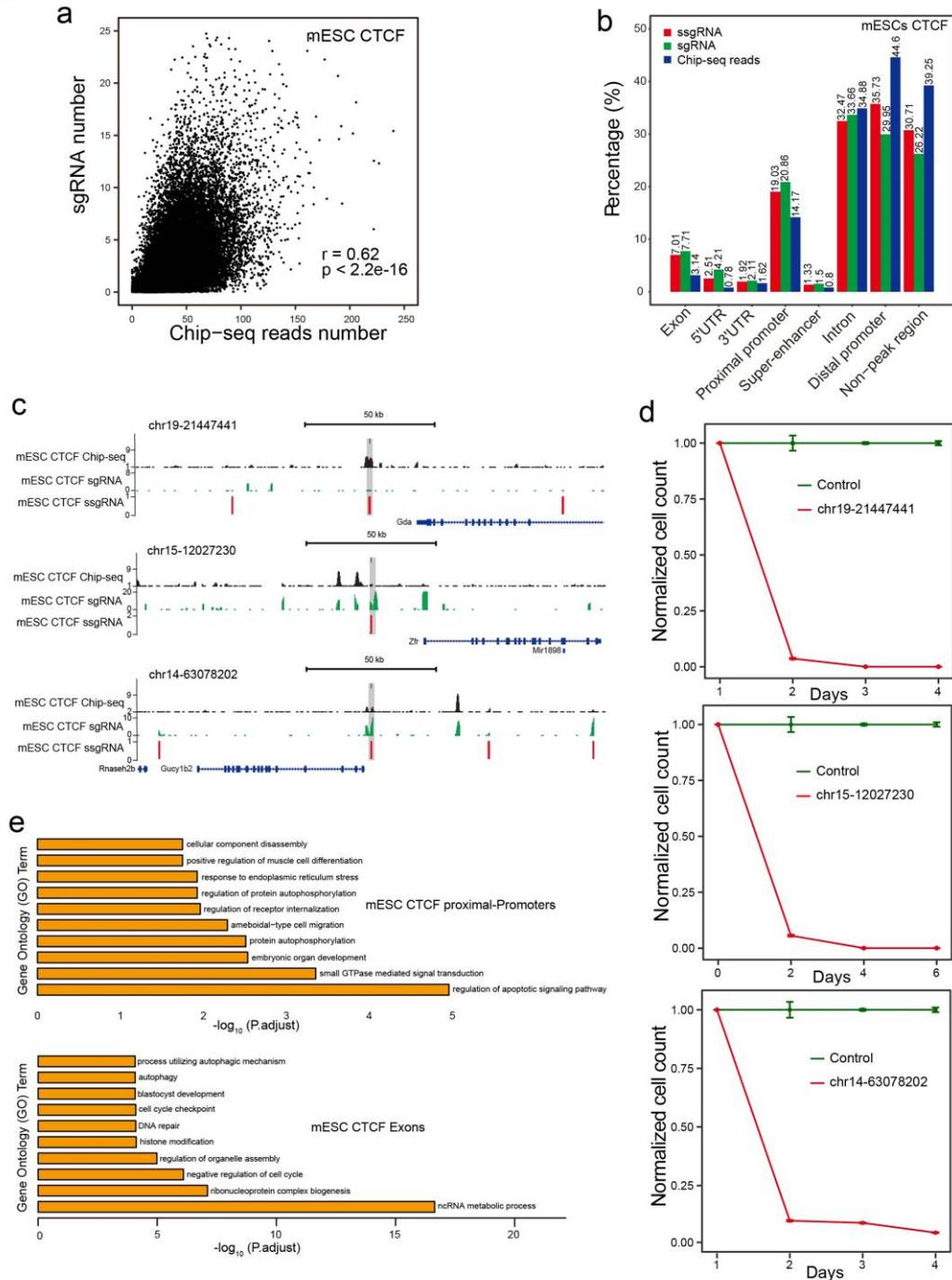
789

790

791

792

Figure S4



793

794

795 **Figure S4**

796 (a) Scatter-plot of the correlation between ChIP-seq reads number and sgRNA number in mESCs
797 CTCF elements. Spearson's correlation coefficient with two-tailed test was calculated.

798 (b) The distribution of CTCF ChIP-seq reads, CTCF sgRNAs and CTCF ssgRNAs on major
799 regulatory regions of mouse genome.

800 (c) Plots of CTCF ChIP-seq reads, CTCF sgRNAs and CTCF ssgRNAs at three genomic loci. Y-
801 axes, RPKM. The genomic regions with CTCF ssgRNAs (validated in Figure 4d) are shaded grey.

802 (d) Three CTCF ssgRNA (validated in Figure 4d) significantly inhibits mESCs proliferation. Cell

803 numbers of three time-points are plotted. Error bars the SD of triplicates.

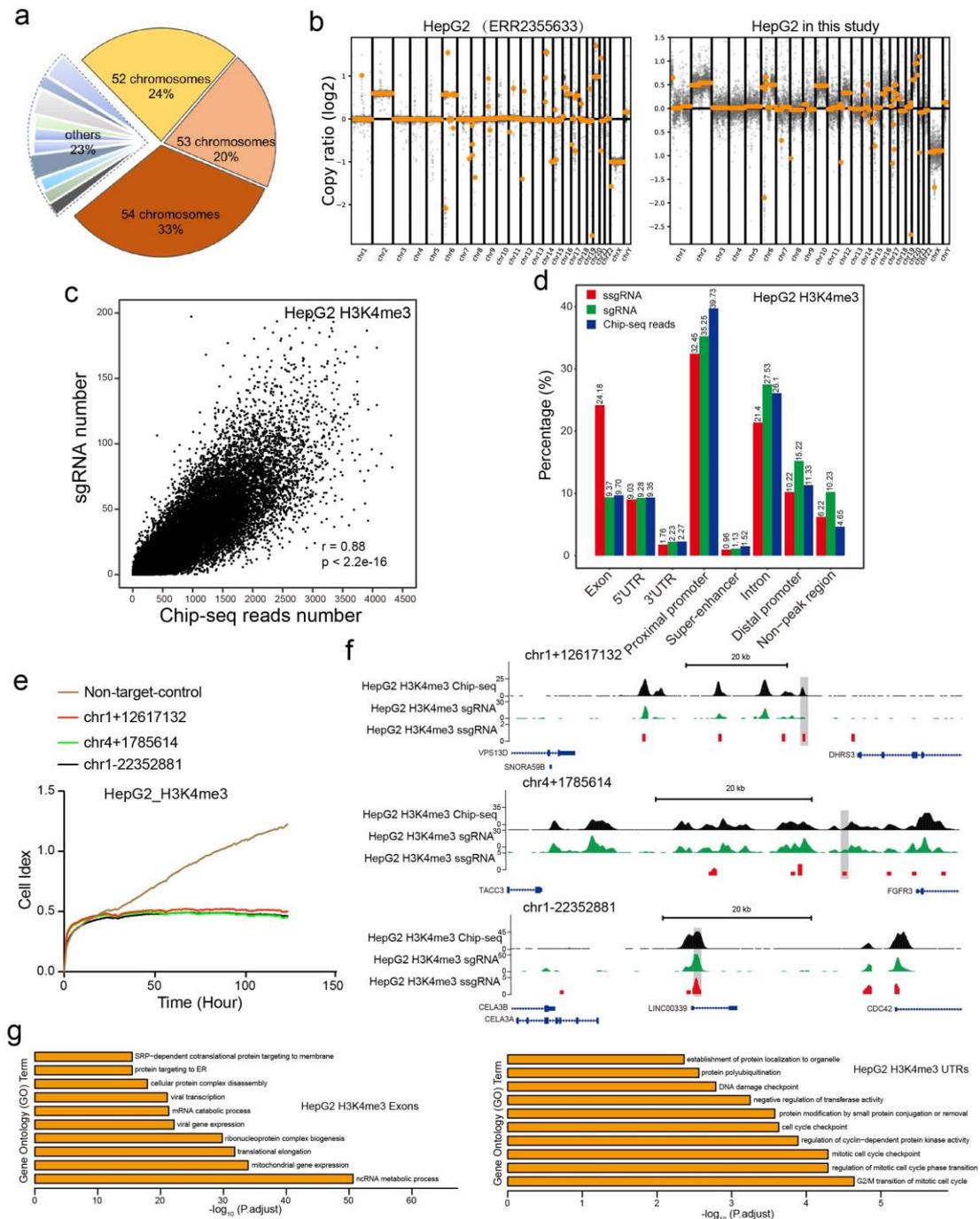
804 (e) Top 10 GO terms of the genes whose proximal promoters and exons are targeted by CTCF
805 ssgRNAs.

806

807

808

Figure S5



809
810

811 **Figure S5**

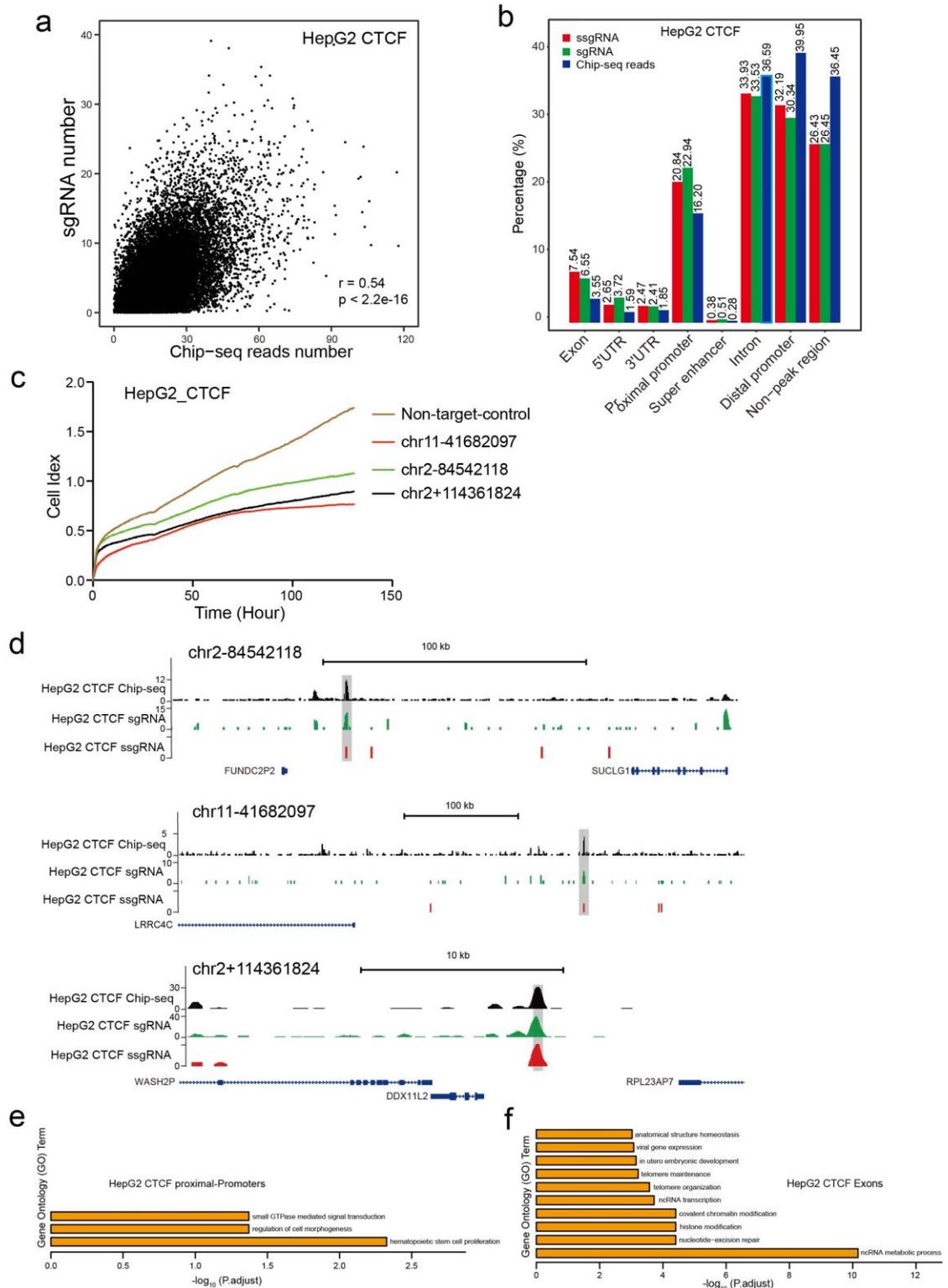
812 (a) Chromosome number of HepG2 used in this study based on 100 Giemsa-stained metaphases,
813 which illustrates that the modal number was 52-54 chromosomes.

814 (b) Copy Number Variation (CNV) analysis of HepG2 in this study through low depth whole
815 genome sequencing. The CNV distribution of HepG2 used in this study is similar with that of a
816 standard HepG2 cell from Sequence Red Archive (SRA) database (ERR2355633).

817 (c) Scatter-plot of the correlation between ChIP-seq reads number and sgRNA number in HepG2
818 H3K4me3 elements. Spearson's correlation coefficient with two-tailed test was calculated.

- 819 (d) The distribution of H3K4me3 ChIP-seq reads, H3K4me3 sgRNAs and H3K4me3 ssgRNAs on
820 major regulatory regions of human genome.
- 821 (e) All three H3K4me3 ssgRNAs (validated in Figure 5d) can significantly inhibit HepG2
822 proliferation. Real-time cell numbers are plotted.
- 823 (f) Plots of H3K4me3 ChIP-seq reads, H3K4me3 sgRNAs and H3K4me3 ssgRNAs at three
824 genomic loci. Y-axes, RPKM. The genomic regions with H3K4me3 ssgRNAs (validated in Figure
825 5d) are shaded grey.
- 826 (g) Top 10 GO terms of the genes whose exons and UTRs are targeted by H3K4me3 ssgRNAs.
827

Figure S6



828

829

830 **Figure S6**

831 (a) Scatter-plot of the correlation between ChIP-seq reads number and sgRNA number in HepG2
832 CTCF elements. Spearson's correlation coefficient with two-tailed test was calculated.

833 (b) The distribution of CTCF ChIP-seq reads, CTCF sgRNAs and CTCF ssgRNAs on major

834 regulatory regions of human genome.

835 (c) All three CTCF ssgRNAs (validated in Figure 6d) can significantly inhibit HepG2 proliferation.

836 Real-time cell numbers are plotted.

837 (d) Plots of CTCF ChIP-seq reads, CTCF sgRNAs and CTCF ssgRNAs at three genomic loci. Y-
838 axes, RPKM. The genomic regions with CTCF ssgRNAs (validated in Figure 6d) are shaded grey.

839 (e) All three GO terms of the genes whose proximal promoters are targeted by CTCF ssgRNAs.

840 (f) Top 10 GO terms of the genes whose exons are targeted by CTCF ssgRNAs.

841

842 **Methods and Materials**

843

844 **Part 1: Wet lab experiments**

845 **sgRNA library construction through CTDE (Controlled Template-dependent**
846 **elongation)**

847 1. 23bp DNA fragment generation

848 Source DNA is fragmented. We ligate them to the A1 adapter and capture them onto
849 T1 streptavidin magnetic beads (Thermo Fisher, 65601). Then we denature DNA with
850 0.1M NaOH (Sigma, 79724) for 10min, and anneal primer for chain extension. We
851 apply 2U Therminator™ DNA Polymerase (NEB, M0261) to incorporate one 3'
852 hydroxyl-reversible dNTP (Jena Bioscience, 3'-O-N₃-dNTPs). We restore the 3'
853 hydroxyl group by 100mM TCEP (Sigma, 646547) treatment and repeat the
854 incorporation-and-reversion for 22 cycles. We blunt DNA with mung bean nuclease
855 (NEB, M0250) for 30min, and then apply T4 PNK (NEB, M0201) to phosphorylate it
856 for 30min.

857 2. NGG PAM selection

858 We ligate the A2 adapter and amplify the library for ten cycles. After the gel extraction
859 of the amplified library, we digest the DNA with AscI (NEB, R0558). Then we capture
860 the library onto streptavidin beads.

861 3. NGG PAM removal

862 We ligate the A3 adapter and amplify the library for ten cycles, and then digest the
863 library with BbsI (NEB, R3539). We fill in the gap with T4 DNA polymerase (NEB,
864 M0203), ligate A4 adapter, and amplify the library with the KAPA HiFi polymerase

865 mix (Roche, KK2631) for ten cycles. We apply 20% TBE-PAGE to select the size of
866 the library (61nt). After releasing DNA from PAGE, we amplify the library by PCR
867 with KAPA HiFi polymerase and primer (the sequence is below) and size-selected via
868 2% agarose gel¹.

869 ArrayF

870 TAACTTGAAAGTATTTTCGATTTCTTGGCTTTATATATCTTGTGGAAAGGAC

871 GAAACACCG

872 ArrayR

873 ACTTTTTCAAGTTGATAACGGACTAGCCTTATTTTAACTTGCTATTTCT

874 AGCTCTAAAAC

875 A1-F GAAAGGACGAAACACCG_T

876 A1-R cGGTGTTCGTCCTTTCCACA_{agat}AGATCGGAAGAGCGTC-Biotin

877 Anneal Primer

878 tgGACGCTCTTCCGATCTATCTTGTGGAAAGGACGAAACACCGT

879 A2-F

880 CGCGCCACACGTCTGAACTCCAGTCACAGTCAACAATCTCGTATGCCGT

881 CTTCTGCT

882 A2-R

883 CAAGCAGAAGACGGCATAACGAGATTGTTGACTGTGACTGGAGTTCAGACG

884 TGTGGGCGCG

885 A3-F

886 TCTTCAAGCTTGGCGTAATCATGGTCATAGCTGTTTCCTGTGTGAAATTGT

887 TATCCGC

888 A3-R

889 AGCGGATAACAATTTACACAGGAAACAGCTATGACCATGATTACGCCAA

890 GCTTGAAGA

891 A4-F GTTTTAGAGCTAGAAATAGCAAGTTAAA

892

893 **Chromatin Immunoprecipitation (ChIP)**

894 The ChIP assay was performed according to manufacture of Millipore ChIP kit

895 (Millipore, 20-153). Briefly, formalin was added to cells to crosslink protein and DNA

896 and then cells (10^6 - 10^7) were lysed in SDS lysis buffer. Chromatin was fragmented via

897 sonication to the size of 200bp-500bp and then immunoprecipitated with antibodies (5-

898 15 μ g each sample). The enriched DNA was purified by QIAquick PCR purification kit

899 (QIAGEN). These DNA samples were applied to CTDE protocol or illumina

900 sequencing library construction. The antibodies are anti-H3K4me3 (Cell Signaling,

901 9751) and anti-CTCF (Millipore, 07-729).

902

903 **Cell culturing**

904 V6.5 ESCs were maintained in DMEM supplemented with 15% fetal bovine serum

905 (Gibco), 0.1mM β -mercaptoethanol, 2mM L-glutamine, 0.1mM nonessential amino

906 acid, 1000U/ml recombinant leukemia inhibitory factor (Millipore LIF, ESG1107), and

907 30U/ml penicillin/streptomycin. For feeder-free culturing, ESCs were grown on plates
908 coated with 0.1% gelatin (Millipore, ES-006-B). HepG2 cells were cultured in DMEM
909 supplemented with 10% FBS (Gibco) and 1% penicillin/streptomycin.

910

911 **sgRNA library cloning into lentiCRISPR v2**

912 The lentiCRISPR v2 vector was digested with BsmBI (Fermentas, FD0454), treated
913 with alkaline phosphatase (Fermentas, EF0654) at 37°C for 2 hours and gel-purified on
914 a 1% TAE-Agarose gel. A 20µl Gibson ligation reaction (NEB, E2611) was performed
915 using 30ng of inserts and 200ng of digested vector. After ligation, 10µl of the reaction
916 was transformed into 100µl of T1 Resistant Chemically competent cells (TransGen
917 Biotech, CD501-02) according to the manufacturer's protocol. The sgRNA diversity of
918 one reaction will be 40-60 thousand. For example one library, whose template is mESC
919 H3K4me3 ChIPed DNA, contains 0.3-0.4 million kinds of sgRNAs, we took about 20
920 assembly reactions so that we can get representative sub-pool. For template is not a
921 linear DNA, increasing assembly reaction will not significantly enhance representation
922 of this library. Plasmid DNA was extracted using Plasmid Preparation Kit (Axygen,
923 AP-MN-P-250G).

924

925 **Lentivirus production and transduction**

926 Lentivirus was produced through the co-transfection of the lentiviral vectors with
927 psPAX2 (Addgene, #12260) and pMD2.G (Addgene, #12259) into HEK293T cells

928 using PEI (Polysciences, 24765). Virus-containing supernatant was collected and
929 filtered through a 0.45µm low protein-binding membrane (Millipore, SLHV033RB) 48
930 hours after transfection. We performed spin-infection in medium containing polybrene
931 (Sigma-Aldrich, H9268) at 1800 rpm for 45min at room temperature. According to
932 Poisson distribution², we ensured a multiplicity of infection (MOI) of less than 0.3 to
933 get single-infected percentage of over 80%. For a 30000-40000 sgRNA diversity sub-
934 pool, we used 1.2×10^7 cells for each transduction. Transduced cells were selected 24
935 hours under puromycin (Gibco, A1113803) for 3 days. 1µg/ml for mESCs and 3µg/ml
936 for HepG2. After puromycin selection (3 days for mESC, 5 days for HepG2), we
937 applied these cells for dropout screen.

938

939 **CRISPR/Cas9 based dropout screen**

940 After puromycin selection, 4 million cells were seeded into a 10-cm dish every
941 generation. 10 million cells from the 1st generation and 10th generation were used for
942 genomic DNA extraction and sgRNA-template targeted sequencing. For DNA
943 extraction, we used ethanol-precipitated method to avoid DNA loss. Cell pellet was
944 lysed in lysis buffer (1% SDS, 50mM Tris, 1mM EDTA, 20µg Protease K, 10µg RNase
945 A) in 65°C water bath overnight. Cell lysate was denatured by
946 phenol:chloroform:isopentanol (25:24:1), and thoroughly mixed by shaking. After
947 centrifugation, DNA was precipitated from supernatant by 2 volume 100% ethanol. We

948 washed DNA pellet twice with chilled 70% ethanol. After air drying, we dissolved
949 DNA in ddH₂O.

950

951 **sgRNA-template targeted sequencing**

952 Two PCR steps were performed:

953 1. PCR from enough genomic DNA to preserve library complexity.

954 sgRNA template containing cassette was amplified using primers specific to the
955 lentiCRISPR v2 vector (Round1-F and Round1-R)

956 Round1-F AATGGACTATCATATGCTTACCGTAACTTGAAAGTATTTTCG

957 Round1-R CCAACTTCTCGGGGACTGTGGGCGATGTGCGCTCTGCCCACTGA

958 24 × 50µl PCR reactions were performed with 2µg genomic DNA in each tube to
959 achieve over 100 × coverage. 20 cycles were used for minus bias.

960 2. Add illumina sequencing adapters.

961 10µl DNA from first round PCR product (combined all tubes together) was used as
962 template for the second round PCR. Second round PCR product was purified by
963 QIAquick PCR clean-up Kit and quantified by Agilent 2200 Bio-analyzer before
964 illumina sequencing.

965 Round 2-F

966 AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTT

967 CCGATCTtcttgtgaaaggacgaaacaccg

968 Round2-R

969 CAAGCAGAAGACGGCATAACGAGATxxxxxxGTGACTGGAGTTCAGACGTGT

970 GCTCTTCCGATCTcgctctgcccactgacgggcaccgg (xxxxxx is a barcode)

971

972 **Anti-neomycin cell-line preparation**

973 Anti-Neomycin mESCs cells were performed as described with modifications³. Briefly,
974 the pEF1a-BirA-V5-neo (#100548, Addgene) plasmid was stably transfected into V6.5
975 ESCs via electroporation. Cells were selected by neomycin and were applied to screen
976 experiment.

977

978 **Commercial synthesis Neo screen library construction**

979 According to 795bp neo fragment sequence, we pick all gRNA with NGG PAM, totally
980 115 gRNA. Also, we chose 20 gRNA from Gecko negative control (they can't target
981 hg19, mm9 and neo reference sequence). We order these gRNA from oligo synthesis
982 company (Shangya), and these oligoes were in uniform format: 5'-
983 GTGGAAAGGACGAAACACCGNNNNNNNNNNNNNNNNNNNNNGTTTTAGA
984 GCTAGAAATAGC-3', where N20 represent neo gRNA or negative control sequence.

985 We combined 115 neo gRNA and 20 negative control equally together, and amplified
986 5µl library(10µM) with 2X HiFi DNA polymerase, Array-F and Array-R primer for 5
987 cycles. 140bp library was gel extracted and assembled into LentiCRISPR V2 plasmid.
988 20 negative control were separately prepared as above and assembled into
989 LentiCRISPR v2 plasmid, and this library was spike-in of CTDE-neo library.

990

991 **Neo fragment obtaining**

992 We obtained neo fragment directly by PCR using pEF1a-BirA-V5-neo (Addgene,
993 #100548) as template, and gel extract 795bp DNA as source of CTDE method to
994 construct CTDE neo library.

995

996 **Neo library Screen experiment**

997 After puromycin selection, 1 million cells were seeded into a 10-cm dish, and 1 million
998 cells were collected and labelled with Prev-selection. Anti-Neomycin mESCs infected
999 by CTDE-Neo and Commercial synthesis Neo library were cultured with 150 μ g/ml
1000 Neomycin (Sangon) for 7 days. 7 days later, cells were collected and labelled with
1001 Post-selection. Prev- and Post-selection sample were extracted genomic DNA and
1002 illumina sequencing library were construct as above.

1003

1004 **Plasmid construction**

1005 Oligoes (for sgRNA or shRNA) were annealed by gradient cooling process and were
1006 ligated to lentiCRISPR v2 (BsmBI digested) or PLKO.1 (digested by AgeI/EcoRI)
1007 using T4 DNA ligase. After transformation, colonies were picked. The confirmed
1008 colonies by sanger sequencing were preserved in glycerol and plasmid were applied to
1009 transfection. The CDC42 cDNA was obtained from ORF collection (Ultimate™ ORFs,

1010 GeneID: 998). The entire open-reading frame was amplified by PCR and cloned into
1011 PLVX vector at BamHI and XhoI sites.

1012

1013 **AP staining and survival assay**

1014 mESCs transduced with lenti-sgRNA were first selected under 1µg/ml puromycin for
1015 3 days. AP staining was performed according to the manufacturer's instructions
1016 (Beyotime Biotechnology, C3206). HepG2 cells transduced with sgRNA or
1017 LINC00339-shRNA were plated 5 days after transduction at 1.8×10^5 cells per well in
1018 a 24-well plate. Then cells were plated in 6-well plates at a density of 3000 cells per
1019 well. The culture media with 1µg/ml puromycin was refreshed every 3 days for ~14
1020 days. After PBS washing, colonies were stained with 0.1% crystal violet (Beyotime,
1021 C0121) for 10min at room temperature. Then the images of each well were recorded by
1022 a digital camera and the number of colonies was counted using ImageJ software.

1023

1024 **Karyotype analysis**

1025 HepG2 (10^5) were trypsinized and resuspended in 5ml 0.56% KCl in water. After
1026 incubation at 37°C for 10min, 1ml 25% freshly made fixative solution (methanol:glacial
1027 acetic acid = 3:1 by volume) in water was added. Cells were pelleted and resuspended
1028 in 1ml 100% fixative solution three times. To make cells spread, one drop of the cell
1029 suspension was dropped (from a height of 0.5m) onto a glass slide and allowed to air

1030 dry before DAPI staining. The mitotic chromosome number was counted under the
1031 fluorescence microscopy.

1032

1033 **Cell cycle analysis**

1034 Cells were transduced with lentivirus via spinfection in 6-well plates. After puromycin
1035 selection, cells were harvested and fixed in chilled 70% ethanol at -20°C overnight, and
1036 stained with 40µg/ml propidium iodide (Sigma-Aldrich, P4170) containing 50µg/ml
1037 RNase (Sigma-Aldrich, R6513) at 37°C for 15min in the dark, then analyzed cell cycle
1038 by FACSCalibur flow cytometer (BD Biosciences). Data were analyzed with ModFit
1039 LT software.

1040

1041 **qPCR quantification**

1042 Total RNA was harvested using the RNeasy Plus Mini Kit (Qiagen, 74134) and 1µg
1043 RNA was used for reverse transcription with the PrimeScript™ RT Master Mix (Takara,
1044 RR036A). After reverse transcription, TB Green® Fast qPCR (Takara, RR430A) was
1045 performed with related primers (Table S14).

1046

1047 **Western blot**

1048 Cells were harvested and lysed in SDS-PAGE sample buffer. Equal amounts of total
1049 protein were loaded in each lane. Proteins were resolved by SDS-PAGE, transferred to
1050 0.45µm PVDF membranes, and probed with the indicated antibodies. The antibodies

1051 used for western blotting are: anti-CDC42 (Cell Signaling Technologies, 2466), anti- β -
1052 tubulin (Cell Signaling Technologies, 2128).

1053

1054 **Real-Time cell number analysis.**

1055 Experiments were performed using Real-Time Cell Analyzer (iCELLigence Analyzer,
1056 ACEA Biosciences Inc.) at 37 °C with 5% CO₂. To measure the background, we placed
1057 200 μ l DMEM medium in the E-plate L8 (ACEA Biosciences, 00300600840). Then, we
1058 added 400 μ l HepG2 (transduced with sgRNA or shRNA) cell suspension (8000
1059 cells/well) with 10% FBS DMEM medium. The impedance (cell index), which was
1060 detected by sensor electrodes in E-plate L8 every 15min for 120 hours. The plot of the
1061 cell index was calculated automatically by the RTCA software package 1.1.1 (ACEA,
1062 Biosciences Inc.).

1063

1064 **Part 2: Dry lab experiments and data analysis**

1065 **Sequencing data pre-processing**

1066 The designed lentiCRISPR v2 sgRNA fragments were sequenced. The sequenced read
1067 comprises a 5'-adapter, a sgRNA, and a 3'-adapter, one by one. We trimmed the read,
1068 remove the pre-designed 5'-adapter and 3'-adapter (Table S14) to get the sequence of
1069 the sgRNA. Then we adhered "NGG" to the end of each sgRNA to recover its PAM
1070 sequence and prepared it for mapping on the reference genome.

1071

1072 **The assessment of sgRNA library generated by CTDE from plasmid lentiCRISPR**

1073 **v2**

1074 1. We evaluated the detected sgRNAs from the designed reference of plasmid
1075 lentiCRISPR v2 (Addgene #52961, 14873bp) with the following steps:

1076 Step 1, we built a mapping reference set from the plasmid lentiCRISPR v2 sequences
1077 using Bowtie build-index function.

1078 Step 2, we prepared three sgRNA libraries with the CTDE approach. For each of the
1079 libraries, we sampled one million reads and then extracted 19-21 bp length sgRNAs
1080 from the “Sequencing data pre-processing” analysis. We merged three sgRNA libraries
1081 to get a sgRNA set.

1082 Step 3, we aligned the sgRNAs set from step2 to the reference set by the bowtie aligner,
1083 allowing one base mismatch and reporting all alignments (-v 1 -a).

1084 Step 4, we counted the number of mapped sgRNAs as x . We identified the “NGG”
1085 locus on either sense or antisense strands of the plasmid lentiCRISPR v2 sequences.

1086 For locus i , we counted the number of sgRNAs whose PAM tails mapped to the locus
1087 as y . Then we calculated the ratio of sgRNAs for locus i (Ratio_{locus_i}) as follows:

1088
$$\text{Ratio}_{locus_i} = \frac{y}{x} \times 100$$

1089 We repeatedly calculated the number and ratio of sgRNAs for each locus. Next we
1090 divided lentiCRISPR v2 sequences into equal-sized 10 bp bins and calculated the ratio
1091 of AT for each bin. We finally used Integrative Genomics Viewer (IGV) (v2.8.0) to

1092 display the ratio of sgRNAs in each “NGG” locus and AT in each bin of the
1093 lentiCRISPR v2 sequences.

1094 2. We calculated the coverage of the detected sgRNAs on the “NGG” locus in the
1095 designed reference with the following steps:

1096 Step 1, we count all possible the “NGG” loci in the plasmid lentiCRISPR v2 sequences
1097 as N_{total} .

1098 Step 2, we count the number of “NGG” loci covered by sgRNAs from the sgRNA set
1099 described above as $N_{covered}$. Then we calculated the coverage of the detected sgRNAs
1100 on the “NGG” locus in the designed reference ($Cover_{SCTE}$) as follows:

$$1101 \quad Cover_{SCTE} = \frac{N_{covered}}{N_{total}}$$

1102 3. We next analyzed the enrichment of the detected sgRNAs on the “NGG” locus with
1103 the following steps:

1104 Step 1, we sequenced the plasmid lentiCRISPR v2, then we trimmed the sequenced
1105 reads to 23bp and aligned them to the reference set using bowtie, not allowing base
1106 mismatch (-v 0). The reads mapped on the plasmid lentiCRISPR v2 sequences and
1107 mapped on the “NGG” locus of the sequences were counted. Next, we calculated the
1108 percentage of reads on the “NGG” locus from the mapped reads as $P_{NGG-PAM}$.

1109 Step 2, we sampled one million reads from each of three sgRNA libraries, then we got
1110 the sgRNAs from the “Sequencing data pre-processing” analysis and extracted the
1111 sgRNAs of 20bp length.

1112 Step 3, we processed the sgRNAs in two different ways. In one way, we removed their
 1113 “NGG” tails and mapped them to the reference set by bowtie, not allowing base
 1114 mismatch (-v 0). Simultaneously, we mapped them to the reference set by bowtie,
 1115 allowing one base mismatch (-v 1). After alignment, we counted the number of mapped
 1116 sgRNAs as $N_{NGG-PAM}$. Then the fold enrichment of the detected sgRNAs on the “NGG”
 1117 locus ($Enrich_{SCTE}$) was estimated as follows:

$$1118 \quad Enrich_{SCTE} = \frac{N_{NGG-PAM}}{N_{total} \times P_{NGG-PAM}}$$

1119 4. We grouped the detected sgRNAs by their length and calculated the proportion of
 1120 sgRNAs of different length with the following steps:

1121 Step 1, for each of the three sgRNA libraries, we got sgRNAs from the “Sequencing
 1122 data pre-processing” analysis and mapped them to reference set by bowtie, allowing
 1123 one base mismatch (-v 1).

1124 Step 2, we grouped sgRNAs by their length. For the sgRNAs of length i , we counted
 1125 the number as N_i , and calculated their proportion as follows:

$$1126 \quad Proportion_i = \frac{N_i}{N_{all}}$$

1127 Here N_{all} is the number of all mapped sgRNAs. We repeatedly calculated the
 1128 proportion for the sgRNAs of 16-24bp from the three sgRNA libraries.

1129 5. To evaluate the how faithful sgRNA synthesis was, we mapped 16-24bp sgRNA to
 1130 reference set by the following priority order: with no mismatch, with one mismatch,
 1131 with two mismatches (not include mismatch in PAM site). For PAM site, we only
 1132 tolerate mismatch at the first base. The error rate was calculated as total number of

1133 mismatch base divided by total number of base of mapped sgRNA (not include PAM
1134 site).

1135

1136 **Comparison between CTDE and commercial method library in Neo resistant**
1137 **screen**

1138 To demonstrate the performance of CTDE, we generated two libraries based on the Neo
1139 sequence using CTDE and commercial methods (Table S14). Twenty non-targeting
1140 control sgRNAs randomly selected from GeCKO v2 library were added to two libraries.
1141 After illumina sequencing and initial data processing, the number of reads were counted
1142 with an in-house script for all libraries. Then library sizes were normalized using total
1143 reads count of control sgRNA in pre- and post-selection library. After normalization,
1144 the fold change of each sgRNA between post- and pre-selection library was calculated.
1145 Targeting sgRNAs whose fold change less than the minimum fold change of non-
1146 targeting control sgRNAs were recognized as ssgRNA.

1147

1148 **Comparison between NSgRNAShot and MAGeCK in identify ssgRNA**

1149 We downloaded one genome-wide CRISPR-Cas9 dropout screen data of mESCs as the
1150 benchmark data⁴. MAGeCK is a widely used method to identify of essential genes from
1151 genome-scale CRISPR/Cas9 knockout screens⁵. As described in the previous paper,
1152 since the CRISPR/Cas9 knockout system should show no difference in selection
1153 preference between control samples or between replicated treatment samples, a good

1154 method should not detect many significantly selected sgRNAs and genes between these
1155 samples⁵. So we took the strategy to detect the possible false-positive ssgRNA with
1156 FDR < 0.1 between two replicates of ESC treatment samples with NSgRNAShot and
1157 MAGeCK based on benchmark data. When multiple sgRNAs are available for one gene,
1158 MAGeCK demonstrated the best performance to detect essential genes⁶. Therefore we
1159 evaluated the sensitivity of ssgRNA identified by our method with the following
1160 strategy: we took the essential genes reported by MAGeCK as the gold standard set, a
1161 good method should largely report the ssgRNA from essential genes but not from other
1162 genes.

1163

1164 **ChIP-seq data analysis**

1165 For ChIP-seq of H3K4me3 and CTCF in mESC and HepG2 cell lines, libraries were
1166 sequenced using Illumina HiSeq X Ten and paired-end 150 bp long reads were obtained.

1167 Chip-seq reads were analyzed with the following steps:

1168 Step 1, we trimmed Chip-seq reads to 100 bp (from 5' to 3'), and for redundant reads
1169 which have the same sequence, only one was retained. Then reads were mapped to the
1170 reference genome (mm9 for mESC or hg19 for HepG2) by bowtie, only uniquely
1171 mapped reads were reported (-m 1).

1172 Step 2, Chip-seq peaks for H3K4me3 and CTCF were detected by MACS2 with default
1173 settings. For each mouse ES CTCF peak, we calculated its length (*peak length*) and

1174 counted the number of reads mapped on the peak (*number of mapped reads*), then
1175 we calculated its average read depth (*average depth per peak*) as follows:

$$1176 \quad \text{average depth per peak} = \frac{\text{reads length} \times \text{number of mapped reads}}{\text{peak length}}$$

1177 We filtered peaks whose average read depth less than ten.

1178 Step 3, we manually identified regions which consistently contain significantly
1179 enriched reads in multiple input ChIP-Seq data provided by Encode project (Table S15)
1180 and regarded them as false positive peak regions. We filtered the peaks overlapped with
1181 these regions.

1182

1183 **RNA-seq data analysis**

1184 1. We downloaded two total RNA-seq datasets of HepG2 from Gene Expression
1185 Omnibus (GEO) database (GEO:GSE88089), and two total RNA-Seq datasets of
1186 mESC from Beijing Institute of Genomics (BIG) Data Center (CRA001133).

1187 2. For each RNA-seq dataset, raw reads were mapped to reference genome (mm9 for
1188 mESC or hg19 for HepG2) by bowtie, allowing one base mismatch and only uniquely
1189 mapped reads were retained (-v 1 -m 1). For redundant reads that fall on the same
1190 position on the genome, only one was retained.

1191 3. Based on the UCSC knownGene annotations, we counted mapped reads in the exon
1192 regions of transcripts, and then we calculated transcript expression levels by
1193 normalizing the number of reads for each transcript to total mapped reads and mRNA
1194 length, namely reads per kilo-bases per million reads (RPKM).

1195 4. We got the average RPKMs of transcripts for mESC and HepG2, respectively, and
1196 the largest RPKM of transcripts for each gene represents gene expression level. We
1197 selected the genes with RPKM higher than one as expressed genes.

1198

1199 **Whole-genome sequencing data analysis**

1200 SNP (single nucleotide polymorphism) existed in the human and mouse population. To
1201 identified the sgRNAs from genomic regions with SNP, we built additional reference
1202 from the regions with the following steps:

1203 1. We downloaded one whole-genome sequencing (WGS) dataset of HepG2 from the
1204 Sequence Read Archive (SRA) database (ERR2355633). Moreover, the WGS library
1205 of mESC was sequenced using Illumina HiSeq X Ten, and paired-end 150 bp long reads
1206 were obtained.

1207 2. For each WGS dataset, paired-end reads were merged together. Then we trimmed
1208 adapter residues from reads by AdapterRemoval (Version 2.1.7), and reads shorter than
1209 150 bp after trimming were filtered out.

1210 3. We performed step-wise mapping process, step 1, reads were mapped to the reference
1211 genome (mm9 for mESC and hg19 for Hepg2) using bowtie2 with default settings. Step
1212 2, mapped reads from step 1 were mapped to the reference genome using bowtie, not
1213 allowing base mismatch (-v 0). Unmapped reads from step 2 were regarded as derived
1214 from regions of mESC or HepG2 with SNP, we used these uniquely unmapped reads
1215 to build a mapping reference set using the Bowtie build-index function.

1216

1217 **Identification of negatively selected sgRNAs**

1218 1. The sgRNA libraries generated by H3K4me3 or CTCF labeled DNA sequences were
1219 sequenced using Illumina HiSeq X Ten. For each sgRNA library, we got sgRNAs from
1220 “Sequencing data pre-processing” analysis.

1221 2. We selected 18-21 bp length sgRNAs which have ability to direct Cas9-induced
1222 indels at target sites for further analysis. Then we performed step-wise mapping process.
1223 Step 1, sgRNAs were mapped to reference genome (mm9 for mESC and hg19 for
1224 HepG2) using bowtie, allowing one base mismatch (-v 1). Step 2, unmapped sgRNAs
1225 from step 1 were mapped to reference set from “Whole-genome sequencing data
1226 analysis” using bowtie, allowing one base mismatch (-v 1). We counted the number of
1227 each mapped sgRNA from above two steps.

1228 3. The sgRNAs whose abundance had significantly reduced from P1 to P10 were
1229 identified using our designed NSgRNAShot algorithm. We used the term ‘ssgRNA’ to
1230 refer to negatively selected sgRNA.

1231

1232 **Genome-wide distribution of chip-seq reads, sgRNAs and ssgRNAs**

1233 We analyzed the genomic distribution of chip-seq reads, sgRNAs, and ssgRNAs with
1234 following steps:

1235 1. We got sgRNAs from “Identification of negatively selected sgRNAs”. Then we
1236 mapped sgRNAs to reference genome (mm9 for mESC and hg19 for HepG2) using
1237 bowtie, allowing one base mismatch (-v 1).

1238 2. We got ssgRNAs from “Identification of negatively selected sgRNAs” and mapped
1239 them to reference genome using bowtie, allowing one base mismatch and refraining
1240 from reporting any alignments have more than three reportable alignments (-v 1 -a -m
1241 3). For ssgRNAs having more than one reportable alignment, we only retained
1242 alignments located in peak region, otherwise, all alignments were retained.

1243 3. We calculated the chip-seq reads, sgRNAs, and ssgRNAs density on each 100 bp
1244 window and sliding along the chromosomes with a step length of 20 bp. We used Circos
1245 (v0.69-6) tool to display the genome-wide density distribution of chip-seq reads,
1246 sgRNAs and ssgRNAs, and local density distribution was visualized by the UCSC
1247 Genome Browser.

1248 4. To display the distribution of chip-seq reads, sgRNAs and ssgRNAs in functional
1249 elements of the genome, we assigned annotation for each nucleotide in the reference
1250 genome by using the following priority order: coding-exon > 5'-UTR > 3'-UTR >
1251 proximal promoter (10kb upstream of a TSS) > super-enhancer > intron > distal
1252 promoter (>10kb upstream of a TSS) = non-peak region. We assigned chip-seq read,
1253 sgRNAs and ssgRNAs to these categories based on location.

1254 5. To demonstrate what the distribution of chip-seq reads, sgRNAs and ssgRNAs per
1255 peak look like, we split each peak into non-overlapping 30 bins. Then for each of these

1256 non-overlapping bins, RPKM of chip-seq reads, sgRNAs and ssgRNAs was calculated
1257 for visualization separately. Peaks were ranked from highest to lowest based on median
1258 RPKM of chip-seq reads of bins.

1259

1260 **Recall rate of essential genes from CTDE**

1261 We downloaded two gene sets described as essential for survival and proliferation of
1262 mESCs^{4, 7}. The overlap of two gene sets was used as reference essential genes for
1263 mESCs. We mapped sgRNAs and ssgRNAs to reference genome (mm9) using bowtie
1264 (-v 1 -m 1), essential genes whose exons were targeted by more than 2 sgRNA were
1265 regarded as covered by CTDE. The ratio of essential genes, which are covered by
1266 CTDE, targeted by ssgRNA is the recall rate of essential genes.

1267

1268 **Gene Ontology (GO) analysis**

1269 We assessed whether related genes of exon, intron, utr, and proximal promoters covered
1270 by ssRNAs were enriched for particular GO categories by calculating adjusted P-value
1271 using the Benjamini & Hochberg method. We performed GO analyses using R package
1272 clusterProfiler with default parameters.

1273 1. For uniquely mapped ssgRNAs generated by H3K4me3 labeled DNA sequences,
1274 only related expressed genes (FPKM > 1) of an exon, utr, intron, and proximal
1275 promoters with ssgRNAs located were used to perform GO analysis.

1276 2. For uniquely mapped ssgRNAs generated by CTCF labeled DNA sequences, all
1277 related genes of utr, intron and proximal promoters with ssgRNAs located were used
1278 for GO analysis, but only related expressed genes of exons with ssgRNAs located were
1279 used for GO analysis.

1280 3. We sorted all GO categories according to adjusted P-values in an ascending order.

1281

1282 **Analysis of CTCF hallmarks shared occupation in different cell types**

1283 1. We downloaded CTCF hallmarks of 16 mouse cell lines and tissues from ENCODE,
1284 the CTCF hallmarks of mESC were occupied in higher than 9 of 16 mouse cell lines
1285 and tissues were defined as common hallmarks.

1286 2. We downloaded CTCF hallmarks of 55 human cell lines from ENCODE, the CTCF
1287 hallmarks of HepG2 were occupied in higher than 33 of 55 human cell lines were
1288 defined as common hallmarks.

1289 3. We mapped ssgRNAs from “Identification of negatively selected sgRNAs” to
1290 reference genome (mm9 for mESC and hg19 for HepG2) using bowtie, allowing one
1291 base mismatch and reported uniquely mapped reads (-v 1 -m 1). Based on the position
1292 of CTCF hallmarks and ssgRNAs in the genome, we calculated the distribution of
1293 ssgRNAs in common and cell-type specific CTCF hallmarks. We used the term
1294 ‘essential common peaks’ to refer to common CTCF hallmarks with ssgRNAs located.

1295

1296

1297 **Negative-SgRNA-Shot (NSgRNAShot) method**

1298 We assumed that in the dish, during culture, 1) most sgRNAs do not change their
 1299 abundance, which suggests they suffer no selection pressure; 2) no sgRNAs is
 1300 undergoing positive selection and increases the abundance; 3) a few sgRNAs are
 1301 suffering negative selection and decrease their abundance. Then we developed a
 1302 method named Negative-SgRNA-Shot (NSgRNAShot) to detect significant-
 1303 negatively-selected sgRNAs (ssgRNA). Here are the details.

1304 Step 1. We got reads counts (r) profiles of sgRNAs at 1st generation (p_1) and 10th
 1305 generation (p_{10}) in the dish from “Identification of negatively selected sgRNAs”
 1306 analysis. Suppose we have detected n and m sgRNAs from p_1 and p_{10} , then we
 1307 calculated the sequencing depth of sgRNAs profiles at p_1 (Sum_{p_1}) and p_{10} ($\text{Sum}_{p_{10}}$)
 1308 with the following formula,

$$1309 \quad \text{Sum}_{p_1} = \sum_{i=1}^n r_i, \quad \text{Sum}_{p_{10}} = \sum_{j=1}^m r_j;$$

1310 If $\text{Sum}_{p_1} \geq \text{Sum}_{p_{10}}$, we calculated the normalized reads counts (nr) of sgRNAs at p_1
 1311 with the formula: $nr_i = r_i$, where $1 \leq i \leq n$; and we calculated the normalized reads
 1312 counts of sgRNAs at p_{10} with the formula: $nr_j = r_j * (\text{Sum}_{p_1}/\text{Sum}_{p_{10}})$, where $1 \leq$
 1313 $j \leq m$.

1314 If $\text{Sum}_{p_1} < \text{Sum}_{p_{10}}$, we calculated the normalized reads counts of sgRNAs at p_1 with
 1315 the formula: $nr_i = r_i * (\text{Sum}_{p_{10}}/\text{Sum}_{p_1})$, where $1 \leq i \leq n$; and we calculated the
 1316 normalized reads counts of sgRNAs at p_{10} with the following formula: $nr_j = r_j$,
 1317 where $1 \leq j \leq m$.

1318 Step 2. We calculated log2 fold change (log2fc) value of abundance of a sgRNA from
 1319 p1 and p10 with the following formula:

$$1320 \quad \log_2 \text{fc} = \log_2 \left(\frac{\max\{nr_{p1}, 10\}}{\max\{nr_{p10}, 10\}} \right),$$

1321 where nr_{p1} and nr_{p10} are the normalized reads count of the sgRNA at p1 and p10.

1322 However, if the sgRNA is not detected at either p1 or p10, we set $nr_{p1} = 0$ or

1323 $nr_{p10} = 0$, correspondingly.

1324 With that, we calculated log2fc of all sgRNAs in the dish.

1325 Step 3. At first, we made some deductions. In the dish, we assumed the number of the

1326 negatively-selected sgRNAs is no more than 20% of the non-selected sgRNAs. We

1327 assumed the log2fc of the former ones follows a normal distribution $N(u, \sigma^2)$, and

1328 the log2fc of the latter ones follows a normal distribution $N(0, \sigma^2)$. Here $u > 0$ and

1329 $u > 4\sigma^2$. Then the log2fc of sgRNAs in the dish follows a mixture-normal

1330 distribution. The mode, which is the value occurring a maximum number of times in

1331 the distribution, will indicate the center of the normal distribution of the log2fc from

1332 the non-selected sgRNAs and should be zero.

1333 Then, we conducted the analysis. We calculated the distribution of log2fc of all

1334 sgRNAs in the dish, find mode (m) of the distribution and correct log2fc of the

1335 sgRNAs with the following formula:

$$1336 \quad \log_2 \text{fc}_m = \log_2 \text{fc} - m.$$

1337 With that, the $\log_2\text{fc}_m$ of sgRNAs in the dish follows a mixture-normal distribution.
 1338 The mode of the distribution indicates the center of the normal distribution of the log
 1339 fold change from the non-selected sgRNAs and is zero.

1340 Step 4. Again, we made some deductions at first. The previous deductions in step 3
 1341 demonstrated the $\log_2\text{fc}_m$ of the negatively-selected sgRNAs follows the normal
 1342 distribution $N(u, \sigma^2)$. We have assumed that $u > 4\sigma^2$, thus the negatively-selected
 1343 sgRNAs whose $\log_2\text{fc}_m < 0$ should occupy 4.5% of all negatively-selected sgRNAs.
 1344 However, the percent of negatively-selected sgRNAs is less than 20% of non-selected
 1345 sgRNAs in the dish. Thus the negatively selected sgRNAs whose $\log_2\text{fc}_m < 0$
 1346 should occupy less than 1% ($0.2/(1+0.2) \times 0.045 = 0.0075$) sgRNAs in the dish. Since this
 1347 percent is relatively small, we assumed that the sgRNAs with $\log_2\text{fc}_m < 0$ are from
 1348 the non-selected sgRNAs but not from the negatively-selected sgRNAs.

1349 Then, we conducted the analysis. We randomly picked a sgRNA k with the positive
 1350 $\log_2\text{fc}_m$ ($\log_2\text{fc}_{m_k}$) from the dish. We took the sgRNAs in the dish whose
 1351 $\log_2\text{fc}_m < -\log_2\text{fc}_{m_k}$ being from the non-selected sgRNAs, and recorded their
 1352 number as x_1 . We counted the ones whose $\log_2\text{fc}_m > \log_2\text{fc}_{m_k}$ in the dish as x_2 .
 1353 If $x_1 > 0$ and $x_2 > x_1$, then we can deduce that x_1 sgRNAs are from the non-
 1354 selected sgRNAs and $x_2 - x_1$ sgRNAs are from the negatively-selected ones in all the
 1355 x_2 sgRNAs. Then, taking $\log_2\text{fc}_{m_k}$ as the threshold, we can calculate the false
 1356 discovery rate (fdr_k) of the negatively-selected sgRNAs as:

$$1357 \quad \text{fdr}_k = x_1/x_2.$$

1358 Step 5. We went through any sgRNA l in the dish with positive $\log_2\text{fc}_m$
 1359 ($\log_2\text{fc}_m_l$), took the $\log_2\text{fc}_m_l$ as a threshold, recorded the number of sgRNAs
 1360 whose $\log_2\text{fc}_m > \log_2\text{fc}_m_l$ in the dish as $x_{2,l}$, and calculated the fdr (fdr_l) of the
 1361 negatively-selected sgRNAs. Then we paired $\log_2\text{fc}_m_l, \text{fdr}_l, x_{2,l}$, recorded the pair of
 1362 any sgRNA to get the pairs set. Next, we screened the pairs set and selected the pairs
 1363 whose $\text{fdr}_l < 0.1$, and further identified the one with the largest $x_{2,l}$ in the selected
 1364 pairs and recorded its $\log_2\text{fc}_m_l$. Finally, taking the $\log_2\text{fc}_m_l$ as the threshold, we
 1365 identified the sgRNAs whose $\log_2\text{fc}_m > \log_2\text{fc}_m_l$ in the dish as the ssgRNA.
 1366 However, if we cannot find any pair whose $\text{fdr} < 0.1$ in pairs set, we claimed no
 1367 ssgRNA in the dish.

1368 Step 6. We can get the ssgRNAs by conducting the steps 1-5. However, some ssgRNAs
 1369 may not be saturated sequenced and have small reads counts at either p1 or p10. The
 1370 change of their abundance from p1 to p10 was suspected. With that, we modified step
 1371 2 to give weights to the \log_2 fold change ($\log_2\text{fc}$) of the sgRNAs in the dish. At first,
 1372 we selected the larger one of its normalized reads counts (nr) at p1 and p10 for each
 1373 sgRNA. Then we ranked all sgRNAs by their selected nr in descending order. Next,
 1374 we calculated its percentile (perct) for each sgRNA in the rank list. At last, picking a
 1375 sgRNA z with the percentile value perct_z , we calculated its $\log_2\text{fc}$ ($\log_2\text{fc}_z$) with
 1376 the following formula,

$$1377 \quad \log_2\text{fc}_z = \text{perct}_z * \log_2\left(\frac{\max\{\text{nr}_{p1}, 10\}}{\max\{\text{nr}_{p10}, 10\}}\right),$$

1378 where nr_{p1} and nr_{p10} are the normalized reads count of the sgRNA at p1 and p10.

1379 However, if the sgRNA is not detected at either p1 or p10, we set $nr_{p1} = 0$ or

1380 $nr_{p10} = 0$, correspondingly.

1381 We modified the \log_2fc of all sgRNAs in the dish by this approach.

1382 Finally, we conducted step 1,6,3,4 and 5 to detect ssgRNAs in the dish.

1383

1384 **Evaluation of Negative-SgRNA-Shot method**

1385 We designed a simulation experiment to test the performance of NSgRNAShot.

1386 At first, we generated a series of simulated datasets.

1387 Step 1. We picked a dish. We conducted steps 1 to 3 of NSgRNAShot to get the

1388 \log_2fc_m of all the sgRNAs in the dish. We also recorded the normalized reads count

1389 of the sgRNAs at p10 (nr_{p10}). For the sgRNAs who got no reads count at p10, we set

1390 their nr_{p10} as zero.

1391 Step 2. We counted sgRNAs whose $\log_2fc_m < 0$ in the dish and recorded the

1392 number as r . We randomly picked a sgRNA p from the r sgRNAs and recorded its

1393 \log_2fc_m as $\log_2fc_{m_p}$. Then we tried to find a sgRNA q and recorded its \log_2fc_m

1394 as $\log_2fc_{m_q}$ where $|\log_2fc_{m_p} + \log_2fc_{m_q}| < 0.05$, thus we got a sgRNAs pair

1395 whose \log_2fc_m values are symmetrical about zero.

1396 We repeated the process $r/2$ times, got $r/2$ sgRNAs pairs to form a sgRNAs set

1397 whose \log_2fc_m values are symmetrical about zero. Since we had recorded the nr_{p10}

1398 of the sgRNAs in the set, we inferred their normalized reads count of the sgRNAs at p1
 1399 (nr_{p1}) with the following formula,

$$1400 \quad nr_{p1} = nr_{p10} * 2^{\log2fc_m} .$$

1401 Here we got a sgRNAs set with the normalized reads count at both p1 and p10. We saw
 1402 the set as the simulated non-selected sgRNAs set.

1403 Step 3. We calculated variation (σ^2) of the simulated non-selected sgRNAs set. We
 1404 randomly picked α ($0 < \alpha < 0.2$) percent of the non-selected sgRNAs, calculated their
 1405 $\log2fc_m$ and modified their $\log2fc_m$ values to $\log2fc_m + \beta\sigma^2$, where ($\beta \geq 4$).

1406 We had recorded the nr_{p10} of these sgRNAs, we inferred their normalized reads count
 1407 of the sgRNAs at p1 (nr_{p1}) with the following formula,

$$1408 \quad nr_{p1} = nr_{p10} * 2^{\log2fc_m + \beta\sigma^2} .$$

1409 Thus we got a sgRNAs set with the normalized reads count at both p1 and p10. We saw
 1410 the set as the simulated negatively-selected sgRNAs set. Finally, we combined the
 1411 simulated non-selected and negatively-selected sgRNAs sets to get a simulated sgRNAs
 1412 dataset.

1413 Step 4. We took (α, β) in step 3 as the parameter pair. We set the parameter pairs as
 1414 (0.01,6), (0.01,8), (0.05,4), (0.05,6), (0.05,8), (0.1,4), (0.1,6), (0.1,8), (0.2,4), (0.2,6),
 1415 and (0.2,8). We repeated steps 3 and 4 one hundred times with each parameter pair to
 1416 generate the simulated datasets. We randomly took a dataset under each parameter pair.
 1417 We plotted the normalized reads count of the non-selected and negatively selected
 1418 sgRNAs in the dataset and calculated the distribution of the log2 fold change of these

1419 sgRNAs. The scatter plot and distribution plot showed that the simulation is successful
1420 because they tend to fit the two plots in the real dataset (Figure S1a).

1421 To this step, we had generated a series of simulated datasets. Then, we tested the
1422 performance of NSgRNAShot on detecting the negatively-selected sgRNAs at the
1423 simulated datasets. We used two indicators-precision (*prec*) and recall (*rec*), which are
1424 calculated with the following formula, to benchmark NSgRNAShot,

$$1425 \quad \textit{prec} = \frac{TP}{TP+FP}, \text{ and } \textit{rec} = \frac{TP}{TP+FN},$$

1426 where *TP* is the number of negatively-selected sgRNA identified as ssgRNA; *FP* is
1427 the number of non-selected sgRNA identified as ssgRNA; *FN* is the number of
1428 negatively-selected sgRNA failed to be identified as ssgRNA.

1429 We run NSgRNAShot on the simulated datasets to detect negatively-selected sgRNAs.
1430 We can see that in the datasets whose parameter pairs are (0.05,8), (0.1,6), (0.1,8),
1431 (0.2,6), and (0.2,8), the precision and recall indicators are both higher than 80% (Table
1432 S1). We believed NSgRNAShot could successfully identify the negatively-selected
1433 sgRNAs in the circumstances.

1434

1435

1436

1437

1438 **References**

- 1439 1. Shalem, O. et al. Genome-scale CRISPR-Cas9 knockout screening in human cells.
1440 *Science* **343**, 84-87 (2014).
- 1441 2. Chen, S. et al. Genome-wide CRISPR screen in a mouse model of tumor growth
1442 and metastasis. *Cell* **160**, 1246-1260 (2015).
- 1443 3. Kim, J., Cantor, A.B., Orkin, S.H. & Wang, J. Use of in vivo biotinylation to study
1444 protein-protein and protein-DNA interactions in mouse embryonic stem cells. *Nature*
1445 *protocols* **4**, 506-517 (2009).

- 1446 4. Tzelepis, K. et al. A CRISPR Dropout Screen Identifies Genetic Vulnerabilities and
1447 Therapeutic Targets in Acute Myeloid Leukemia. *Cell Rep* **17**, 1193-1205 (2016).
1448 5. Li, W. et al. MAGeCK enables robust identification of essential genes from
1449 genome-scale CRISPR/Cas9 knockout screens. *Genome biology* **15**, 554 (2014).
1450 6. Bodapati, S., Daley, T.P., Lin, X., Zou, J. & Qi, L.S. A benchmark of algorithms for the
1451 analysis of pooled CRISPR screens. *Genome biology* **21**, 62 (2020).
1452 7. Shohat, S. & Shifman, S. Genes essential for embryonic stem cells are associated
1453 with neurodevelopmental disorders. *Genome research* **29**, 1910-1918 (2019).
1454
1455

Figures

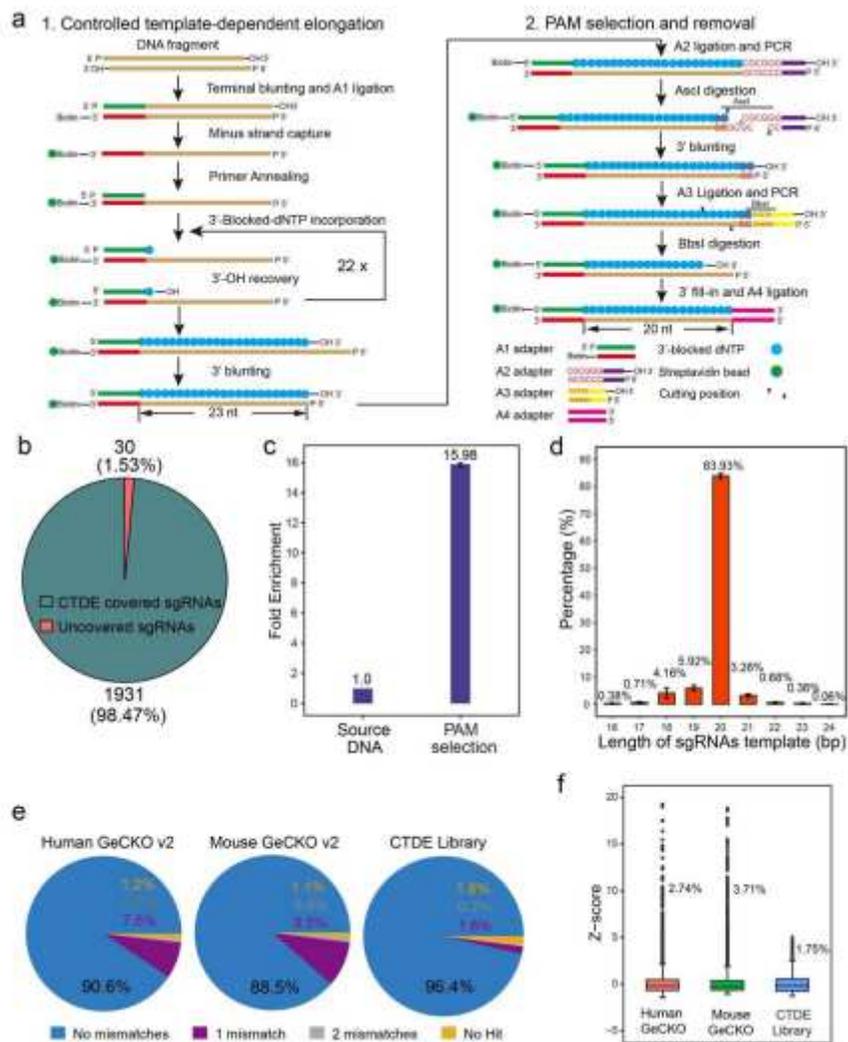


Figure 1

[Please see the manuscript file to view the figure caption.]

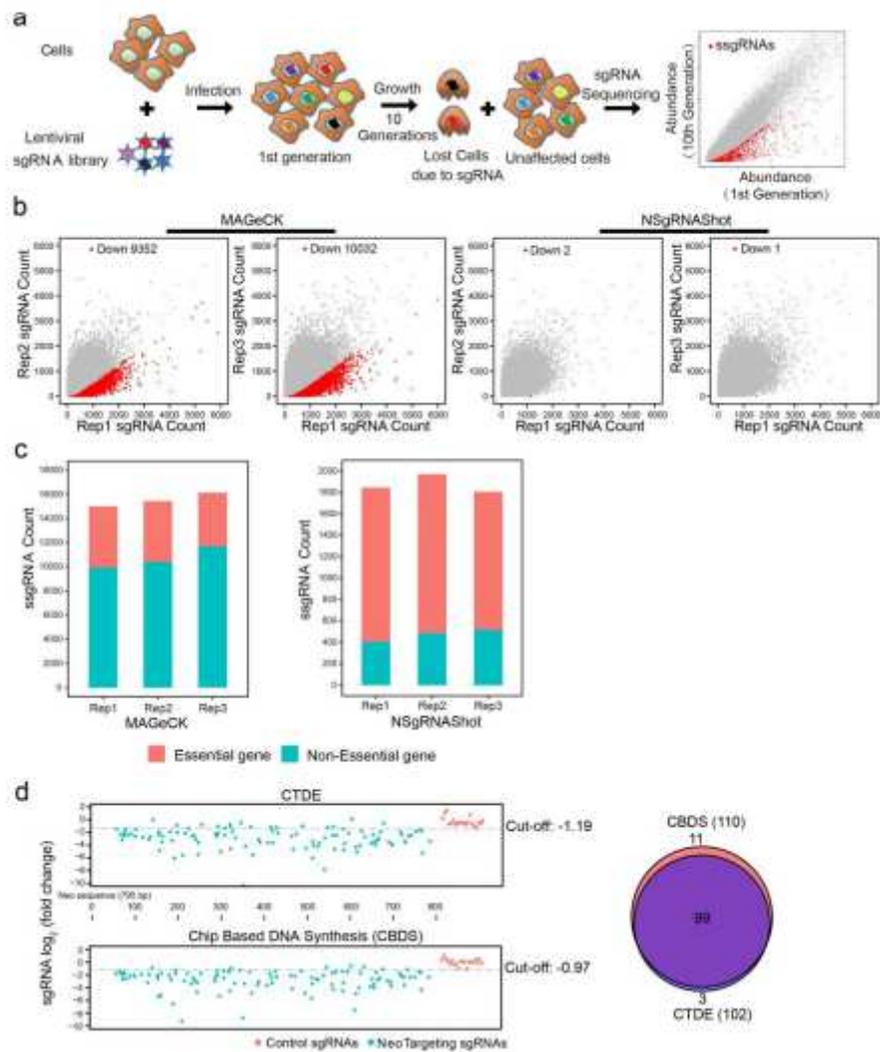


Figure 2

[Please see the manuscript file to view the figure caption.]

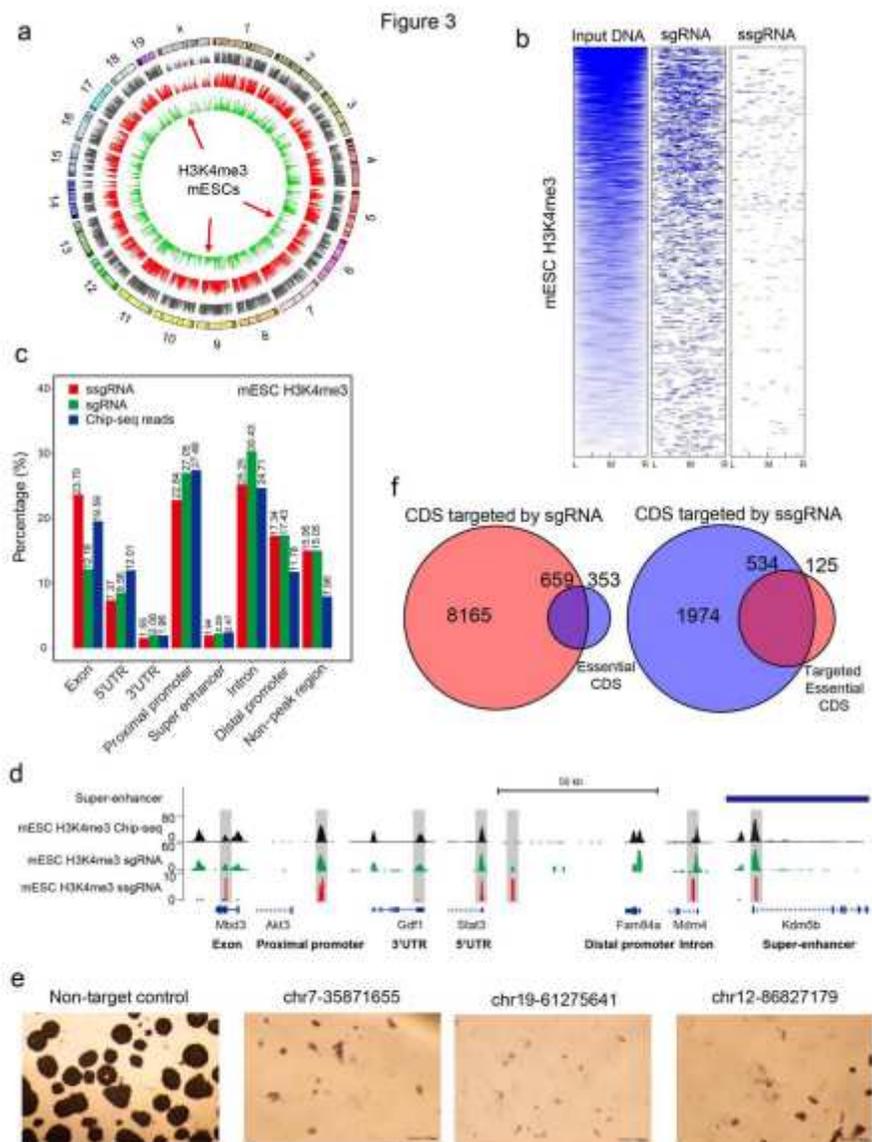


Figure 3

[Please see the manuscript file to view the figure caption.]

Figure 4

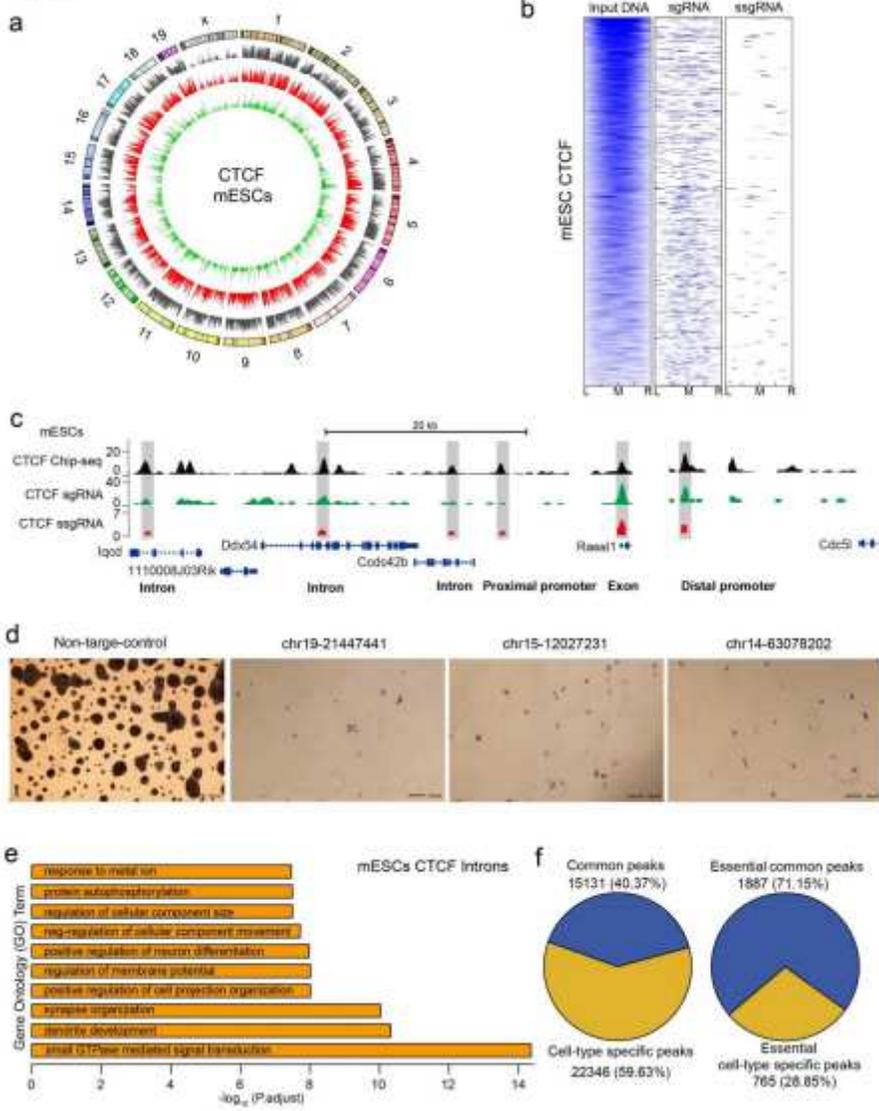


Figure 4

[Please see the manuscript file to view the figure caption.]

Figure 5

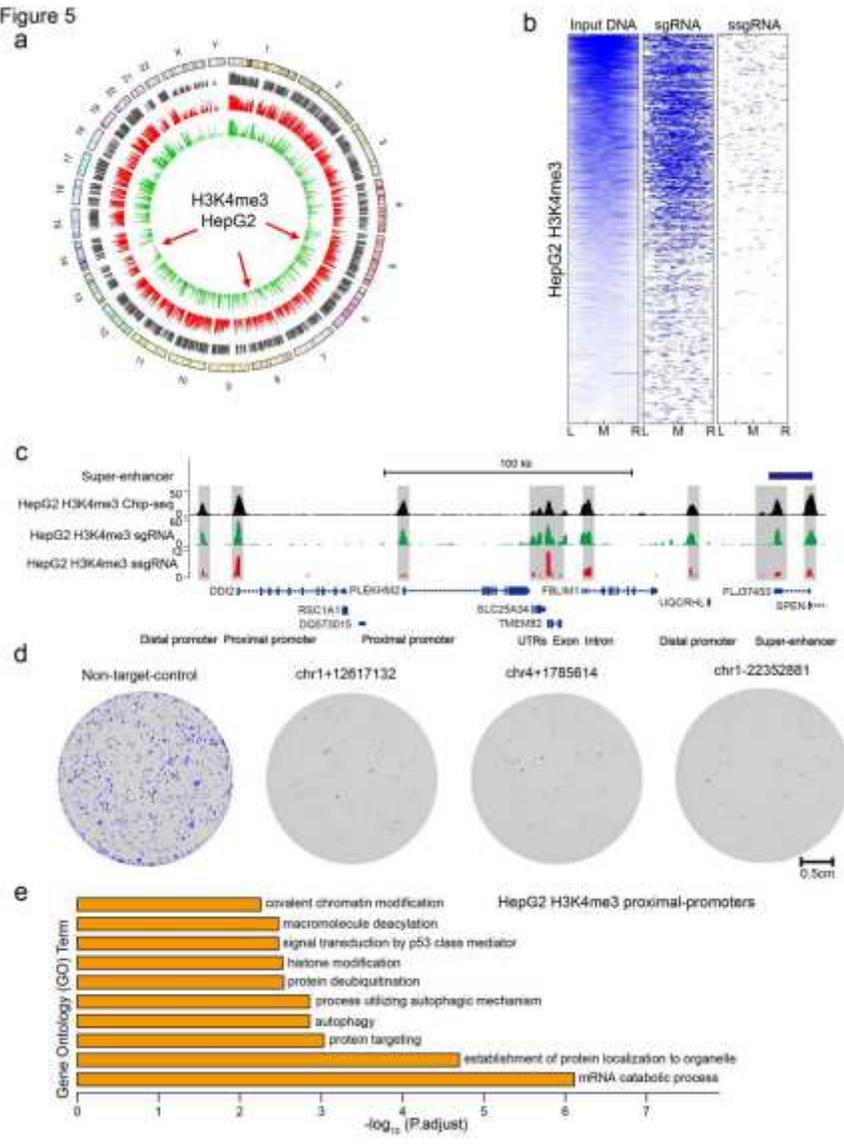


Figure 5

[Please see the manuscript file to view the figure caption.]

Figure 6

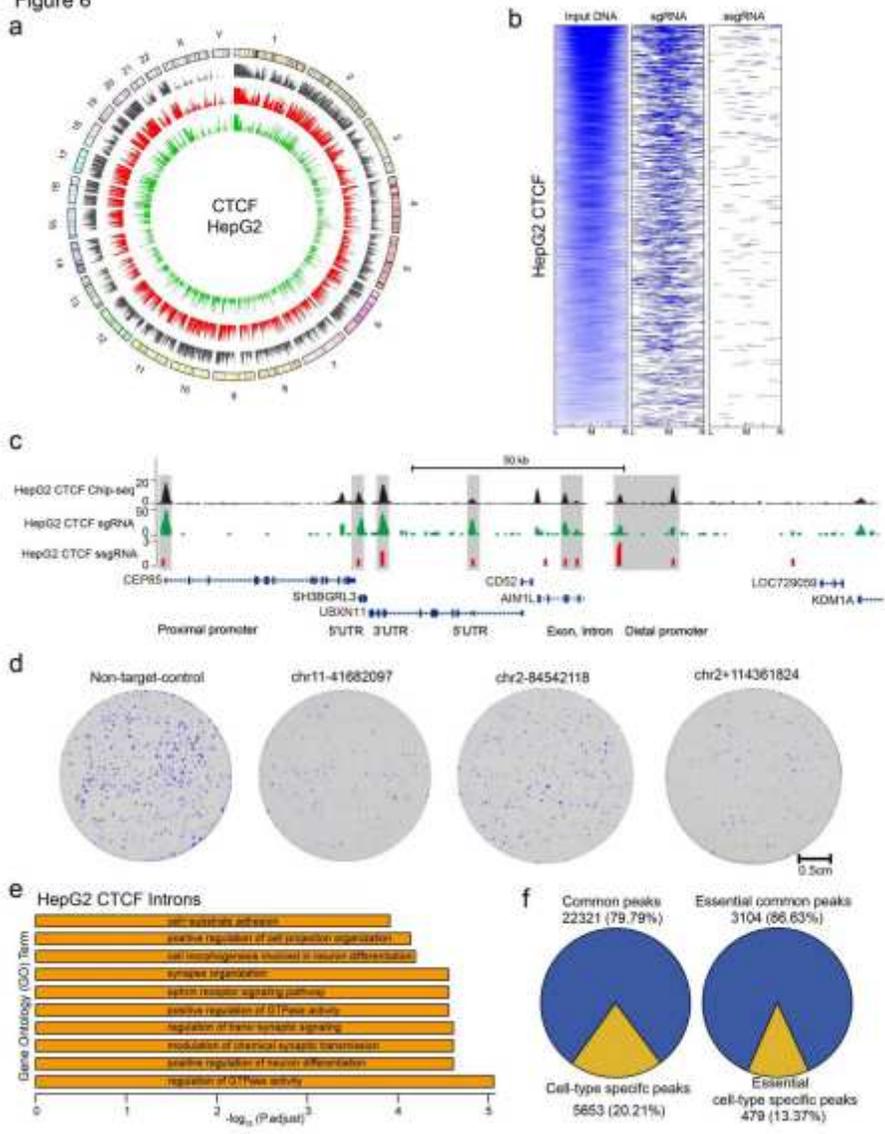


Figure 6

[Please see the manuscript file to view the figure caption.]

Figure 7

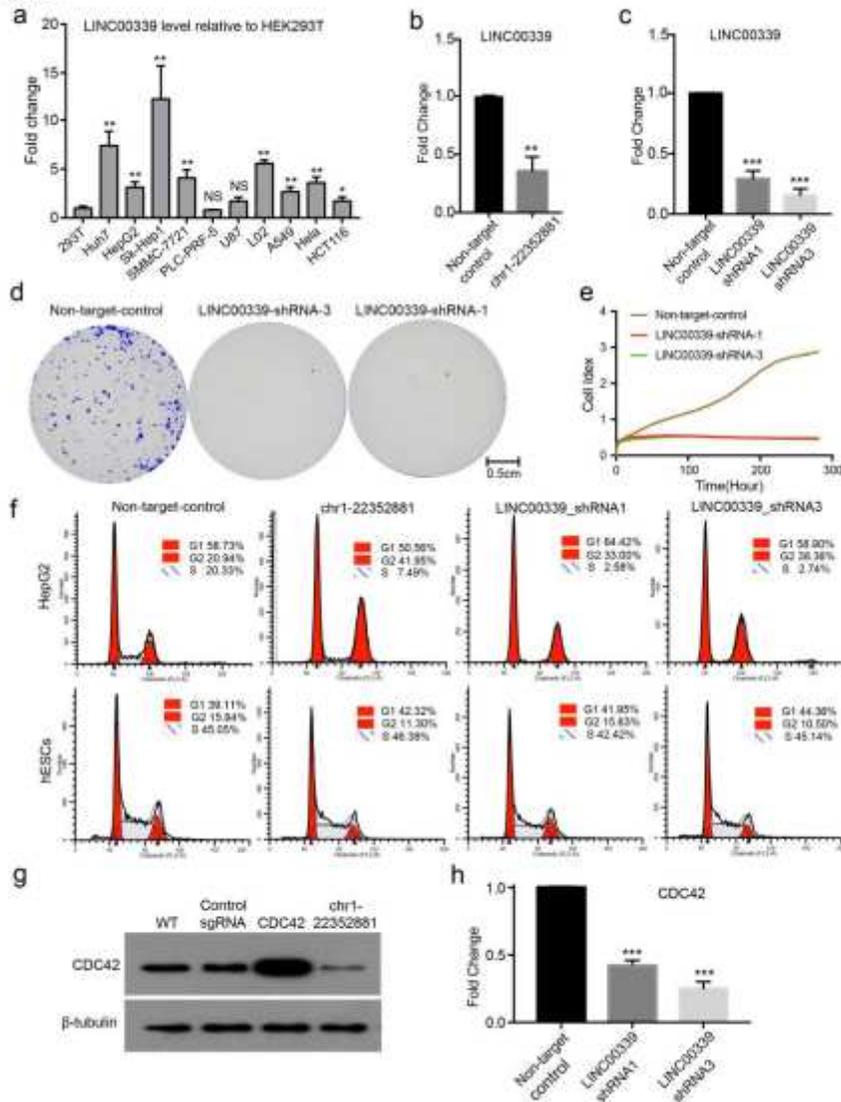


Figure 7

[Please see the manuscript file to view the figure caption.]

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementalTable.zip](#)