

1 **HYENA detects oncogenes activated by distal enhancers in**
2 **cancer**

3
4 Anqi Yu ^{1,*}, Ali E. Yesilkanal ^{1,*}, Ashish Thakur ², Fan Wang ¹, Yang Yang ¹, William Phillips
5 ¹, Xiaoyang Wu ^{1,3}, Alexander Muir ^{1,3}, Xin He ², Francois Spitz ², Lixing Yang ^{1,2,3,**}

- 6 1. Ben May Department for Cancer Research, University of Chicago, Chicago IL, USA
7 2. Department of Human Genetics, University of Chicago, Chicago IL, USA
8 3. University of Chicago Comprehensive Cancer Center, Chicago, IL, USA

9 * these authors contributed equally

10 ** correspondence: lixingyang@uchicago.edu

11

12

13 **Abstract**

14 Somatic structural variations (SVs) in cancer can shuffle DNA content in the genome, relocate
15 regulatory elements, and alter genome organization. Enhancer hijacking occurs when SVs
16 relocate distal enhancers to activate proto-oncogenes. However, most enhancer hijacking studies
17 have only focused on protein-coding genes. Here, we develop a computational algorithm
18 “HYENA” to identify candidate oncogenes (both protein-coding and non-coding) activated by
19 enhancer hijacking based on tumor whole-genome and transcriptome sequencing data. HYENA
20 detects genes whose elevated expression is associated with somatic SVs by using a rank-based
21 regression model. We systematically analyze 1,146 tumors across 25 types of adult tumors and
22 identify a total of 108 candidate oncogenes including many non-coding genes. A long non-
23 coding RNA *TOBI-ASI* is activated by various types of SVs in 10% of pancreatic cancers
24 through altered 3-dimensional genome structure. We find that high expression of *TOBI-ASI* can
25 promote cell invasion and metastasis. Our study highlights the contribution of genetic alterations
26 in non-coding regions to tumorigenesis and tumor progression.

27

28 Introduction

29 At mega-base-pair scale, linear DNA is organized into topologically associating domains (TADs)
30 ¹, and gene expression is regulated by DNA and protein interactions governed by 3D genome
31 organization. Enhancer-promoter interactions are mostly confined within TADs ²⁻⁴. Non-coding
32 somatic single nucleotide variants (SNVs) in promoters and enhancers have been linked to
33 transcriptional changes in nearby genes and tumorigenesis ⁵. Structural variations (SVs),
34 including deletions, duplications, inversions, and translocations, can dramatically change TAD
35 organization and gene regulation ⁶ and subsequently contribute to tumorigenesis. Previously, we
36 discovered that *TERT* is frequently activated in chromophobe renal cell carcinoma by relocation
37 of distal enhancers ⁷, a mechanism referred to as enhancer hijacking (**Fig. 1A**). In fact, many
38 oncogenes, such as *BCL2* ⁸, *MYC* ⁹, *TALI* ¹⁰, *MECOM/EVII* ¹¹, *GFII* ¹², *IGF2* ¹³, *PRDM6* ¹⁴, and
39 *CHD4* ¹⁵, can be activated through this mechanism. These examples demonstrate that genomic
40 architecture plays an important role in cancer pathogenesis. However, the vast majority of the
41 known enhancer hijacking target oncogenes are protein-coding genes, and few non-coding genes
42 have been reported to promote diseases through enhancer hijacking. Here, we refer to non-
43 coding genes as all genes that are not protein-coding. They include long non-coding RNAs
44 (lncRNAs), pseudogenes, and other small RNAs such as microRNAs, small nuclear RNAs
45 (snRNAs), small nucleolar RNAs (snoRNAs), etc. They are known to play important roles in
46 many biological processes ¹⁶, and some are known to drive tumorigenesis ^{17,18}. In this study, we
47 will focus on identifying oncogenes, including oncogenic non-coding genes, activated by
48 enhancer hijacking.

49 Several existing algorithms can detect enhancer hijacking target genes based on patient cohorts,
50 such as CESAM ¹³ and PANGEA ¹⁵. These two algorithms implemented linear regression and
51 elastic net model (also based on linear regression) to associate elevated gene expression with
52 nearby SVs, respectively. PANGEA also considers the effects of somatic SNVs on gene
53 expression. However, a major drawback of these algorithms is that linear regression is quite
54 sensitive to outliers. Outliers are very common in gene expression data from cancer samples and
55 can seriously impair the performances of these algorithms. In addition, CESAM is optimized for
56 microarray data, while PANGEA depends on the annotation of tissue-specific promoter-enhancer
57 pairs, which are not readily available for many tumor types. Cis-X ¹⁹ and NeoLoopFinder ²⁰ can
58 detect enhancer hijacking target genes based on individual samples. However, these tools have
59 limitations in detectable genes and input data. Cis-X detects *cis*-activated genes based on allele-
60 specific expression, which requires the genes to carry heterozygous SNVs. NeoLoopFinder takes
61 Hi-C, Chromatin Interaction Analysis with Paired-End Tag (ChIA-PET), or similar data
62 measuring chromatin interactions as input, which remain very limited. Furthermore, the
63 identification of recurrent mutational events that result in oncogenic activation requires large
64 patient cohorts. Therefore, tools that use whole-genome and transcriptome sequencing data,
65 which are available at much larger sample sizes, would be more useful in identifying SV-driven
66 oncogene activation. Finally, no non-coding oncogenes have been reported as enhancer hijacking
67 targets by the above algorithms. A recent study on SVs altering gene expression in Pan-Cancer
68 Analysis of Whole Genomes (PCAWG) samples ²¹ only considered protein-coding genes but not
69 non-coding genes.

70 Here, we developed Hijacking of Enhancer Activity (HYENA) using normal-score regression
71 and permutation test to detect candidate enhancer hijacking genes (both protein-coding and non-
72 coding genes) based on tumor whole-genome and transcriptome sequencing data from patient
73 cohorts. Among the 108 putative oncogenes detected by HYENA, we studied the oncogenic
74 functions of a lncRNA, *TOBI-ASI*, and demonstrated that it is a regulator of cancer cell invasion
75 in vitro and tumor metastasis in vivo.
76

77 **Methods**

78 **Datasets**

79 This study used data generated by the Pan-Cancer Analysis of Whole Genomes (PCAWG). We
80 limited our study to a total of 1,146 tumor samples for which both whole-genome sequencing
81 (WGS) and RNA-Seq data were available. The data set was composed of cancers from 25 tumor
82 types including 23 bladder urothelial cancers (BLCA), 88 breast cancers (BRCA), 20 cervical
83 squamous cell carcinomas (CESC), 68 chronic lymphocytic leukemias (CLLE), 51 colorectal
84 cancers (COAD/READ), 20 glioblastoma multiforme (GBM), 42 head and neck squamous cell
85 carcinomas (HNSC), 43 chromophobe renal cell carcinomas (KICH), 37 renal clear cell
86 carcinomas from the United States (KIRC), 31 renal papillary cell carcinomas (KIRP), 18 low-
87 grade gliomas (LGG), 51 liver cancers from United States (LIHC), 67 liver cancers from Japan
88 (LIRI), 37 lung adenocarcinomas (LUAD), 47 lung squamous cell carcinomas (LUSC), 95
89 malignant lymphomas (MALY), 80 ovarian cancers (OV), 74 pancreatic cancers (PACA), 19
90 prostate adenocarcinomas (PRAD), 49 renal clear cell carcinomas from European Union/France
91 (RECA), 34 sarcomas (SARC), 34 skin cutaneous melanomas (SKCM), 29 stomach
92 adenocarcinomas (STAD), 47 thyroid cancers (THCA), and 42 uterine corpus endometrial
93 carcinomas (UCEC). More detailed information on the sample distribution and annotation can be
94 found in **Supplementary Table S1**.

95 WGS and RNA-Seq data analysis of tumor and normal samples were performed by the PCAWG
96 consortium as previously described²¹. Somatic and germline SNVs, somatic copy number
97 variations (CNVs), SVs, and tumor purity were detected by multiple algorithms and consensus
98 calls were made. Genome coordinates were based on the hg19 reference genome and GENCODE
99 v19 was used for gene annotation. Gene expression was quantified by HT-Seq (version 0.6.1p1)
100 as fragments per kilobase of million mapped (FPKM). Clinical data such as donor age and sex
101 were downloaded from the PCAWG data portal (<https://dcc.icgc.org/pcawg>). *TOBI* and *TOBI-AS1*
102 expression data in CCLE pancreatic cancer cell lines were downloaded from DepMap
103 Public 22Q2 version (<https://depmap.org/portal/download/all/>). Gene expression data of the
104 Cancer Genome Atlas (TCGA) PAAD cohort (TCGA.PAAD.sampleMap/HiSeqV2_PANCAN)
105 and International Cancer Genome Consortium (ICGC) PACA-CA cohort for 45 samples of
106 which “analysis-id” were labeled as “RNA” were downloaded from Xena Data Hubs
107 (<https://xenabrowser.net/datapages/>) and ICGC data portal ([https://dcc.icgc.org/projects/PACA-](https://dcc.icgc.org/projects/PACA-CA)
108 [CA](https://dcc.icgc.org/projects/PACA-CA)) respectively.

109 Significant expression quantitative trait loci (eQTL)-gene pairs (v8) were downloaded from the
110 Genotype-Tissue Expression (GTEx) data portal (<https://gtexportal.org/home/datasets>). Only
111 those eQTLs that had a hg19 liftover variant ID were included in the analysis and hg38 variants
112 without corresponding hg19 annotation were discarded.

113 The raw sequencing data for Hi-C and ATAC-Seq were available through NCBI Sequence Read
114 Archive (SRA) with accession number PRJNA1036282. The raw sequencing data for mouse
115 xenograft tumor RNA-Seq were available through NCBI SRA with accession number
116 PRJNA1011356.

117

118 **HYENA algorithm**

119 First, small tandem duplications (<10 kb) were discarded since they are unlikely to produce new
120 promoter-enhancer interactions. The remaining SVs were mapped to the flanking regions (500
121 kb upstream and downstream of transcription start sites [TSSs]) of annotated genes. SVs that fall
122 entirely within a gene body were also discarded. The SV status of each gene was defined by the
123 presence or absence of SV breakpoints within the gene or its flanking regions for each tumor.
124 The binary variable SV status was used in the normal-score regression model below. Only genes
125 carrying SVs in at least 5% of samples carrying SVs were tested. For each gene, samples with
126 that gene highly amplified (>10 copies) were removed from the regression model.

127 ***Gene expression normal scores***

128 Gene expression quantifications (fragments per kilobase per million [FPKM]) were quantile
129 normalized (FPKM-QN) using the *quantile.normalize()* function from the *preprocessCore* R
130 package to enhance cross-sample comparison. For each gene, samples were ranked based on
131 their expression values, the ranks were mapped to a standard normal distribution and the
132 corresponding z scores were gene expression normal scores. Normal-score conversion forced the
133 expression data into a Gaussian distribution, allowing for parametric comparisons between
134 samples.

135 ***Normal-score regression***

136 A generalized linear model was used to test associations between gene expression normal scores
137 and SV status and control for confounding variables such as gene copy number, tumor sample
138 purity, donor age, and sex. To capture unobserved variations in gene expression, the first n
139 principal components (PCs) of the expression data were also included in the regression model,
140 where n was determined as 10% of the sample size of the cohort and up to 20 if the sample size
141 was more than 200. The regression model was as shown below:

142 $\text{Expression_normal_score} \sim \text{sv_status} + \text{copy_number} + \text{purity} + \text{age} + \text{sex} + \text{PC}_1 + \text{PC}_2 \dots + \text{PC}_n$

143 For each gene, all PCs were tested for associations with the SV status of that gene, and those PCs
144 that significantly correlate (Mann-Whitney test, $P < 0.05$) with SV status were not used in
145 regression. A similar strategy was used to detect eQTLs in normal tissues²².

146 ***Calculating empirical P values and model selection***

147 Gene expression data were permuted 1000 times by randomly shuffling expression values within
148 the cohort. For tumor types with more than 10,000 genes to test (**Supplementary Table S1**),
149 only 100 permutations were performed to reduce run time. The normal-score regression was
150 performed in the same way on observed gene expression and permuted expression. *P* values for
151 SV status from permuted expression were pooled as a null distribution. Then the *P* values for SV
152 status from observed expression and the *P*-value null distribution were used to calculate
153 empirical *P* values. One-sided *P* values were used since we were only interested in elevated gene
154 expression. False discovery rates (FDRs) were calculated using the Benjamini-Hochberg

155 procedure. Genes with FDR less than 0.1 were considered candidate genes. For example, in
156 MALY, there were 1,863 genes reaching 5% SV frequency and 1,863 P values were obtained in
157 each permutation. After 1000 permutations, 1,863,000 P values were generated and should
158 represent the null distribution very well. Empirical P values were calculated using these
159 1,863,000 permuted P values.

160 The above empirical P value calculation and candidate gene detection were performed iteratively
161 with no PCs and up to n PCs in the regression model. When different numbers of PCs were
162 included in the model, the numbers of candidate genes varied. The regression model with the
163 lowest number of PCs reaching 80% of the maximum number of candidate genes in all
164 regression models tested was selected as the final model to avoid over fitting. For example, the
165 sample size for PCAWG UCEC was 42; therefore, we tested from 0 to 4 PCs. Among these, the
166 model including 4 PCs gave the highest number (4) of candidate genes. Therefore, the model
167 including 4 PCs with 4 candidate genes was selected as the final model (**Supplementary Table**
168 **S2**).

169 In our normal-score regression, we essentially attempt to model variations in gene expression.
170 Including confounding factors will improve performance. Tumor purity, gene copy number,
171 patient age, and sex are factors known to affect gene expression. Therefore, they were included
172 in the regression model. Unobserved variations may include tumor subtype, tumor stage, patient
173 ethnicity, smoking status, alcohol consumption, and other unknown factors that may alter gene
174 expression. Since HYENA was designed for wide applications, we did not require users to
175 provide information on tumor subtype, tumor stage, patient ethnicity, smoking status, alcohol
176 consumption, etc. Principle component analysis is a linear decomposition of gene expression
177 variations. Therefore, including PCs in a regression model was suitable for removing systematic
178 variations and could better model the effects of SV status. However, some enhancer hijacking
179 target genes are master transcription factors, such as *MYC*, and have a profound impact on the
180 gene expression of multiple pathways. Hence, it is possible that some PCs capture the activities
181 of transcription factors. If these transcription factors were activated by somatic SVs, the PCs
182 would be correlated with SV status. Including these PCs would diminish our ability to detect the
183 effects of SV status. Therefore, we excluded these PCs from the regression model.

184 ***Testing eQTL-SV associations***

185 Known germline eQTLs from the matching tissues were obtained from GTEx (**Supplementary**
186 **Table S3**). The associations between germline genotypes of eQTLs and SV status of the
187 candidate genes in the PCAWG cohort were tested using a Chi-squared test. Genes with
188 significant correlations ($P < 0.05$) between their SV status and at least one eQTL were removed.
189 The remaining genes were our final candidate enhancer-hijacking target genes.

190

191 **Benchmarking**

192 Known enhancer hijacking target genes in PCAWG tumor types were selected to test the
193 sensitivity of HYENA, CESAM and PANGEA. The genes included *MYC* in malignant

194 lymphoma, *BCL2* in malignant lymphoma, *CCNE1* in stomach/gastric adenocarcinoma, *TERT* in
195 chromophobe renal carcinoma, *IGF2* in colorectal cancer, *IGF2* in stomach/gastric
196 adenocarcinoma, *IGF2BP3* in thyroid cancer, and *IRS4* in lung squamous cell carcinoma. The
197 same SVs, CNVs, and SNVs were used as input for all three algorithms. For CESAM and
198 PANGEA, upper-quantile normalized fragments per kilobase per million (FPKM-UQ) were
199 normalized by tumor purity and gene copy number, and then used as gene expression inputs.
200 CESAM was run using default parameters, and FDR of 0.1 was used to select significant genes.
201 PANGEA requires predicted enhancer-promoter (EP) interactions based on ChIP-Seq and RNA-
202 Seq data. The EP interactions were downloaded from EnhancerAtlas 2.0
203 (<http://www.enhanceratlas.org/>) (**Supplementary Table S4**). EP interactions from multiple cell
204 lines of the same type were merged. PANGEA was run with default parameters as well and
205 significant genes were provided by PANGEA (multiple testing adjusted *P* value <0.05). To test
206 HYENA, CESAM, and PANGEA for false positives, 20 random gene expression datasets for
207 malignant lymphoma and breast cancer were generated by randomly shuffling sample IDs in
208 gene expression data. HYENA, CESAM, and PANGEA were run with random expressions in
209 the same way as above.

210

211 **Predicting 3D genome organization**

212 A 1 Mb sequence was extracted from the reference genome centered at each somatic SV
213 breakpoint and was used as input for Akita²³ to predict the 3D genome organization. Two 500
214 kb sequences were merged according to the SV orientation to construct the sequence of the
215 rearranged genome fragments. Akita was used to predict the genome organization for the
216 rearranged sequence. High-resolution Micro-C data obtained from human H1-ESCs and HFF
217 cells²⁴ were used to facilitate TAD annotation together with predicted genome organization.
218 H3K27Ac and CCCTC-binding factor (CTCF) ChIP-Seq data from the PANC-1 cell line were
219 downloaded from the ENCODE data portal (<https://www.encodeproject.org/>). SV breakpoints
220 were provided to Orca²⁵ to predict 3D genome structures through its web interface
221 (<https://orca.zhoulab.io/>).

222

223 **In situ Hi-C and ATAC-Seq**

224 Ten million cells of Panc 10.05, PANC-1, PATU-8988S, and PATU-8988T cell lines were
225 collected to construct Hi-C libraries²⁶. The Hi-C libraries were sequenced on Illumina NovaSeq
226 X Plus platform with 1% phix. About 2 billion reads were obtained from Panc 10.05, PATU-
227 8988S, and PATU-8988T, and 1 billion reads were obtained from PANC-1. The paired-end reads
228 were aligned to chromosomes 1-22, X, Y and M by bwa-mem. SVs were identified by EagleC²⁷
229 at 5 kb, 10 kb and 50 kb resolutions. The non-redundant SVs in **Supplementary Table S5** were
230 combined for the three resolutions. Chromatin loops were identified by NeoLoopFinder²⁰. A
231 probability threshold of 0.95 was used, and default values were used for all other parameters.
232 Fifty thousand cells of Panc 10.05, PATU-8988S, and PATU-8988T cell lines were harvested to
233 construct ATAC-Seq libraries²⁸. The libraries were sequenced using Illumina NovaSeq. About

234 60 million reads were generated from each library. The paired-end reads were aligned to the
235 reference genome by hisat2. Hi-C and ATAC-Seq read coverages were generated by deepTools
236 with 10 bp bin-size, RPGC normalization, and an effective genome size of 2,864,785,220.

237

238 **Cell lines**

239 HEK293T, PANC-1, and PATU-8988T cells were obtained from Dr. Alexander Muir
240 (University of Chicago). Panc 10.05 was purchased from ATCC (American Type Culture
241 Collection, USA) (<https://www.atcc.org/products/crl-2547>) and PATU-8988S was purchased
242 from DSMZ (<https://www.dsmz.de/collection/catalogue/details/culture/ACC-204>). All cell lines
243 were cultured at 37°C/5% CO₂. HEK293T cells and PANC-1 cells were cultured in Dulbecco's
244 Modified Eagle Medium (DMEM) (Gibco, 21041025) containing 10% fetal bovine serum (FBS)
245 (Gibco, A4766), and Panc 10.05 cells were cultured in RPMI-1640 medium (Gibco, 11875093)
246 containing 10% FBS, as per ATCC instructions (<https://www.atcc.org/products/crl-3216>,
247 <https://www.atcc.org/products/crl-1469>, <https://www.atcc.org/products/crl-2547>). PATU-8988T
248 and PATU-8988S cells were cultured with DMEM containing 5% FBS, 5% horse serum (Gibco,
249 26050088), and 2 mM L-glutamine as recommended by DSMZ (Deutsche Sammlung von
250 Mikroorganismen und Zellkulturen, Germany)
251 (<https://www.dsmz.de/collection/catalogue/details/culture/ACC-162>). All cell lines have been
252 regularly monitored and tested negative for mycoplasma using a mycoplasma detection kit
253 (Lonza, LT07-218).

254

255 ***TOBI-AS1* and luciferase overexpression**

256 A 1,351 bp *TOBI-AS1* complementary DNA (cDNA) (ENST00000416263.3) was synthesized
257 by GenScript (New Jersey, USA) and subcloned into the lentiviral pCDH-CMV-MCS-EF1-Puro
258 plasmid (SBI, CD510B-1). The cDNA sequence in the plasmid was verified by Sanger
259 sequencing at University of Chicago Medicine Comprehensive Cancer Center core facility. The
260 *TOBI-AS1* overexpression plasmid was amplified by transforming Stellar™ Competent Cells
261 (Takara, 636763) with the plasmid as per instructions and isolated by QIAGEN HiSpeed Plasmid
262 Midi Kit (QIAGEN, 12643). LucOS-Blast vector was obtained from Dr. Yuxuan Phoenix Miao
263 (University of Chicago), cloned, and amplified as described above.

264 HEK293T cells were plated in T-25 flasks and grown to 75% confluence prior to transfection.
265 For each T-25 flask, 240µl Opti-MEM (Gibco, 31985070), 1.6µg pCMV-VSV-G, 2.56µg
266 pMDLg/pRRE, 2.56µg pRSV-Rev, 3.4µg *TOBI-AS1* overexpression vector and 22.8µl TransIT-
267 LT1 Transfection Reagent (Mirus, MIR 2306) were mixed and incubated at room temperature
268 for 30 minutes, then added to the plated HEK293T cells with fresh medium. The luciferase
269 vector was packaged into lentivirus with the same method. Upon 48 hours of incubation,
270 lentiviral supernatant was collected, filtered through 0.45-µmpolyvinylidene difluoride filter
271 (Millipore), and mixed with 8µg/ml polybrene. PANC-1 or PATU-8988T cells at 60%

272 confluence were transduced with the lentiviral supernatant for 48 hours followed by three rounds
273 of antibiotic selection with 4µg/ml puromycin for *TOBI-ASI* overexpression and 10µg/ml
274 blasticidin for the luciferase expression. *TOBI-ASI* expression was validated by quantitative
275 reverse transcription polymerase chain reaction (qRT-PCR), and luciferase expression was
276 validated by in vitro bioluminescence imaging in black wall 96-well plates (Corning, 3603). D-
277 luciferin potassium salt (Goldbio, LUCK-100) solution with 0, 1.25, 2.5, 5 and 10µl 15mg/ml
278 was added into the wells as serial dilutions, and imaging was obtained after 5 minutes. Finally,
279 *TOBI-ASI* overexpression or empty pCDH transduced cell lines with luciferase co-expression
280 were built for both PATU-8988T and PANC-1 cells.

281

282 ***TOBI-ASI* transient knock-down using antisense oligonucleotides (ASOs)**

283 Three Affinity Plus® ASOs were synthesized by Integrated DNA Technologies (IDT), with two
284 targeting *TOBI-ASI* and one non-targeting negative control. The ASO sequences were:

285 Non-targeting ASO (NC): 5' -GGCTACTACGCCGTCA- 3'

286 *TOBI-ASI* ASO1: 5' -GCCGATTTGGTAGCTA- 3'

287 *TOBI-ASI* ASO2: 5' -CTGCGGTTTAACTTCC- 3'

288 The ASOs were transfected into PATU-8988S and Panc 10.05 cells with Lipofecatmine™ 2000
289 (Invitrogen, 11668019) using reverse-transfection method according to IDT protocol
290 (<https://www.idtdna.com/pages/products/functional-genomics/antisense-oligos>) with a final ASO
291 concentration of 9 nM. Cells were transfected in 6-well plates and incubated for 48 hours to
292 reach 60% confluence before RNA extraction or Transwell assay.

293

294 **RNA isolation and qRT-PCR**

295 Cells were plated in 6-well plates and allowed to reach 80% confluence, or transfected by ASOs
296 as described above, prior to RNA extraction. After cells lysis in 300µl/well TRYzol™
297 (Invitrogen, 15596026), RNA samples were prepared following the Direct-zol RNA Miniprep kit
298 manual (RPI, ZR2052). Reverse transcription was performed using Applied Biosystems High-
299 Capacity cDNA Reverse Transcription Kit (43-688-14) following manufacturer's instructions.
300 Quantitative PCR (qPCR) was conducted on StepOnePlus Real-Time PCR System (Applied
301 Biosystems, 4376600), using PowerUp SYBR Green Master Mix (A25742) following the
302 manufacturer's instructions with a primer concentration of 300nM in 10µl reaction systems.
303 Primers were ordered from Integrated DNA Technologies. Primer sequences used in this study
304 are as follows:

305 *TOBI* forward: 5' -GGCACTGGTATCCTG AAA AGCC- 3'

306 *TOBI* reverse: 5' – GTGGCAGATTGCCACGAACATC- 3'

307 *TOBI-ASI* forward: 5' -GGAGTGGTCAGGTGACTGATT- 3'

308 *TOBI-ASI* reverse: 5' -ATTCCACTCCTGTTTGCAACT- 3'

309 *GAPDH* forward: 5' – ACCACAGTCCATGCCATCAC- 3'

310 *GAPDH* reverse: 5' -TCCACCACCCTGTTGCTGTA- 3'

311 Relative expression levels for *TOBI-ASI* and *TOBI* were calculated by the $2^{(-\Delta\Delta C_T)}$ method
312 based on *GAPDH* expression as an endogenous control.

313

314 **Transwell assay for cell invasion in vitro**

315 Transparent PET membrane culture inserts of 24-well plate (Falcon, 353097) were coated with
316 Cultrex Reduced Growth Factor Basement Membrane Extract (BME) (R&D Systems, 3533-010-
317 02) at 50 μ g per membrane (200 μ l of 0.25mg/ml BME stock per membrane) at 37°C for an hour.
318 A total of 100,000 PANC-1 cells/well, 50,000 PATU-8988T cells/well, 50,000 Panc 10.05
319 cells/well, or 50,000 PATU-8988S cells were resuspended in serum-free, phenol-red free DMEM
320 medium and seeded into the coated inserts. Phenol-red free DMEM of 500 μ l (Gibco, A1443001)
321 with 10% FBS was added to the bottom of the wells and the cells were allowed to invade for 16
322 hours. Additional wells with 500 μ l serum-free, phenol-red free DMEM medium without FBS in
323 the bottom chamber were seeded with the same number of cells as indicated above as a negative
324 control. At the end of the assay, the membranes were stained with 500 μ l 4 μ g/ml Calcein AM
325 (CaAM) (Corning, 354216) for one hour at 37°C. The cells that failed to invade were removed
326 from the top chamber with a cotton swab and all inserts were transferred into 1x Cell
327 Dissociation Solution (Bio-Techne, 3455-05-03) and shaken at 150rpm for an hour at 37°C.
328 Finally, CaAM signal from the invaded cells was measured by a plate reader (Perkin Elmer
329 Victor X3) at 465/535nm.

330

331 **Tumor metastasis in vivo**

332 All animal experiments for this study were approved by the University of Chicago Institutional
333 Animal Care and Use Committee (IACUC) prior to execution. Male NOD *scid* gamma (NSG)
334 mice were ordered from the Jackson Laboratory (strain#005557). For tail vein inoculation, mice
335 were injected intravenously through the tail vein with luciferase-expressing at 400,000
336 cells/mouse for PANC-1 cells in cold phosphate buffered saline (PBS) (Gibco, 10010-023). For
337 orthotopic inoculation, mice were injected with 200,000 PANC-1 cells/mouse into the pancreas
338 under general anesthesia. Cells were resuspended in cold PBS containing 5.6mg/mL Cultrex
339 Reduced Growth Factor BME (R&D Systems, 3533-010-02). Primary tumor and metastatic
340 tumor burdens were measured weekly for 4 and 6 weeks for tail vein injection models and
341 orthotopic models, respectively, via bioluminescence imaging using Xenogen IVIS 200 Imaging
342 System (PerkinElmer) at the University of Chicago Integrated Small Animal Imaging Research
343 Resource (iSAIRR) Facility. Each mouse was weighed and injected intra-peritoneally with D-
344 luciferin solution at a concentration of 150 μ g/g of body weight 14 minutes prior to image
345 scanning ventral side up.

346

347 **Ex vivo IVIS imaging**

348 Ex vivo imaging was done for the PANC-1 orthotopic injection mice after 8 weeks of orthotopic
349 inoculation. Mice were injected intra-peritoneally with D-luciferin solution at a concentration of
350 150µg/g of body weight immediately before euthanasia. Immediately after necropsy, mice were
351 dissected, and tissues of interest (primary tumors, livers and spleens) were placed into individual
352 wells of 6-well plates covered with 300 µg/mL D-luciferin. Tissues were imaged using Xenogen
353 IVIS 200 Imaging System (PerkinElmer) and analysis was performed (Living Image Software,
354 PerkinElmer) maintaining the regions of interest (ROIs) over the tissues as a constant size.

355

356 **Tumor RNA sequencing and gene expression analysis**

357 RNA was isolated from mouse subcutaneous tumors (six *TOBI-AS1* overexpression and six
358 control mice) after 6 weeks of PANC-1 cell subcutaneous injection using Direct-zol RNA
359 Miniprep kit (RPI, ZR2052). The quality and quantity of the RNA were assessed using Qubit.
360 Sequencing was performed using the Illumina NovaSeq 6000. About 40 million reads were
361 sequenced per sample. The pair-end reads were aligned to mouse genome (mm10) and human
362 genome (hg19) with hisat2, and the reads mapped to mouse or human genomes were
363 disambiguated using AstraZeneca-NGS disambiguate package. Gene counts were generated with
364 htseq-count. Differential gene expression was analyzed using DESeq2. Differentially expressed
365 genes were defined as genes with an FDR smaller than 0.1 and a fold change greater than 1.5.

366

367 **Code availability**

368 The HYENA package is available at [https://github.com/yanglab-](https://github.com/yanglab-computationalgenomics/HYENA)
369 [computationalgenomics/HYENA](https://github.com/yanglab-computationalgenomics/HYENA).

370

371 **Results**

372 **HYENA workflow**

373 Conceptually, the SVs leading to elevated gene expression are eQTLs. The variants are SVs
374 instead of commonly used germline single nucleotide polymorphisms (SNPs) in eQTL analysis.
375 With somatic SVs and gene expression measured from the same tumors through WGS and RNA-
376 Seq, we can identify enhancer hijacking target genes by eQTL analysis. However, the
377 complexities of cancer and SVs pose many challenges. For instance, there is tremendous inter-
378 tumor heterogeneity—no two tumors are identical at the molecular level. In addition, there is
379 substantial intra-tumor heterogeneity as tumor tissues are always mixtures of tumor, stromal, and
380 immune cells. Moreover, genome instability is a hallmark of cancer, and gene dosages are
381 frequently altered²⁹. Furthermore, gene expression networks in cancer are widely rewired³⁰, and
382 outliers of gene expression are common.

383 Here, we developed an algorithm HYENA to overcome the challenges described above (see
384 more details in Methods Section). We used a gene-centric approach to search for elevated
385 expression of genes correlated with the presence of SVs within 500 kb of transcription start sites
386 (**Fig. 1B**). Although promoter-enhancer interaction may occur as far as several mega-bases,
387 mega-base-level long-range interactions are extremely rare. In addition, although duplicated
388 enhancers can upregulate genes^{31,32}, we do not consider these as enhancer hijacking events since
389 no neo-promoter-enhancer interactions are established. However, small deletions can remove
390 TAD boundaries or repressive elements and lead to neo-promoter-enhancer interactions (**Fig.**
391 **1A**). Therefore, small tandem duplications were discarded, and small deletions were retained.
392 For each gene, we annotated SV status (presence or absence of nearby SVs) for all samples.
393 Samples in which the testing genes were highly amplified were discarded since many of these
394 genes are amplified by circular extrachromosomal DNA (ecDNA)³³, and ecDNA can promote
395 accessible chromatin³⁴ with enhancer rewiring³⁵. Only genes with nearby SVs in at least 5% of
396 tumors were further considered. In contrast to CESAM and PANGEA, we did not use linear
397 regression to model the relationships between SV status and gene expression because linear
398 regression is sensitive to outliers and many false positive associations would be detected³⁶.
399 Instead, we used a rank-based normal-score regression approach. After quantile normalization of
400 gene expression for both protein-coding and non-coding genes, we ranked the genes based on
401 quantile-normalized expression and transformed the ranks to the quantiles of the standard normal
402 distribution. We used the z scores (normal scores) of the quantiles as dependent variables in
403 regression. In the normal-score regression model, tumor purity, copy number of the tested gene,
404 patient age, and sex were included as covariates since these factors confound gene expression.
405 We also included gene expression principal components (PCs) that were not correlated with SV
406 status to model unexplained variations in gene expression. To deduce a better null distribution,
407 we permuted the gene expression 100 to 1000 times (**Supplementary Table S1**) and ran the
408 same regression models. All *P* values from the permutations were pooled together and used as
409 the null distribution to calculate empirical *P* values. Then, multiple testing corrections were
410 performed on one-sided *P* values since we are only interested in elevated gene expression under
411 the influence of nearby SVs. Finally, genes were discarded if their elevated expression could be

412 explained by germline eQTLs. The remaining genes were candidate enhancer hijacking target
413 genes.

414 **Benchmarking performances**

415 There is no gold standard available to comprehensively evaluate the performance of HYENA.
416 We compared HYENA's performance to two other algorithms—CESAM and PANGEA. All
417 three algorithms were run on the same somatic SVs and gene expression data from six types of
418 adult tumors profiled by the PCAWG (**Supplementary Table S1**): malignant lymphoma
419 (MALY), stomach/gastric adenocarcinoma (STAD), chromophobe renal cell carcinoma (KICH),
420 colorectal cancer (COAD/READ), thyroid cancer (THCA), and lung squamous cell carcinoma
421 (LUSC)²¹. Note that PANGEA depends on promoter-enhancer interactions predicted from cell
422 lines, and such data were not available for thyroid tissue. Therefore, thyroid cancer data were not
423 analyzed by PANGEA. To compare the performance of HYENA to the other algorithms, we
424 used the following three strategies.

425 First, we used eight known enhancer hijacking target genes including *MYC*⁹, *BCL2*⁸, *CCNE1*³⁷,
426 *TERT*⁷, *IGF2*^{13,37} (in two tumor types), *IGF2BP3*³⁸ and *IRS4*¹³ to test the sensitivities. Out of
427 the eight known enhancer hijacking genes, HYENA detected four (*MYC*, *BCL2*, *TERT*, and
428 *IGF2BP3*) (**Fig. 2A** and **Supplementary Fig. S1A**), CESAM detected three (*MYC*, *BCL2*, and
429 *TERT*), and PANGEA did not detect any (**Fig. 2A**). In the five tumor types analyzed by all three
430 algorithms, HYENA identified a total of 25 candidate genes, CESAM identified 19, whereas
431 PANGEA identified 255 genes (**Fig. 2B**, **Supplementary Tables S6**, **S7**, and **S8**). Six genes
432 were detected by both HYENA and CESAM, while PANGEA had little overlap with the other
433 algorithms (**Fig. 2B**). The ability of the algorithms to detect known target genes seems to be
434 sensitive to sample size. Both *IGF2* and *IRS4* were initially discovered by CESAM as enhancer
435 hijacking target genes using CNV breakpoints profiled by microarray with much larger sample
436 sizes (378 colorectal cancers and 497 lung squamous cell carcinomas)¹³. In the PCAWG, there
437 were many fewer samples with both WGS and RNA-Seq data available (51 colorectal cancers
438 and 47 lung squamous cell carcinomas). Neither *IGFR* nor *IRS4* was detected by any algorithms.
439 *IGF2* reached 5% SV frequency cutoff required by HYENA, however its FDR did not reach the
440 significance cutoff (**Supplementary Fig. S1B**). In stomach/gastric adenocarcinoma, *IGF2* and
441 *CCNE1* were identified as enhancer hijacking target genes in a cohort of 208 samples³⁷. Neither
442 of these genes was detected by any algorithms because there were only 29 stomach tumors in the
443 PCAWG. Therefore, known target genes missed by HYENA were likely due to small sample
444 size. In summary, HYENA had the best sensitivity of the three algorithms.

445 Second, we also expect immunoglobulin genes to be detected as enhancer hijacking candidates in
446 B-cell lymphoma due to V(D)J recombination. In B cells, V(D)J recombination occurs to join
447 different variable (V), joining (J), and constant (C) segments to produce antibodies with a wide
448 range of antigen recognition ability. Therefore, certain segments have elevated expression and
449 the recombination events can be detected as somatic SVs. Of the 16 genes detected by HYENA
450 in malignant lymphoma (B-cell derived Burkitt lymphomas³⁹), there were two immunoglobulin
451 light chain genes from the lambda cluster (*IGLC7* and *IGLJ7*) and an immunoglobulin-like gene
452 *IGSF3* (**Supplementary Table S6**). CESAM detected 11 genes, one of which was an

453 immunoglobulin gene (*IGLC7*) (**Supplementary Table S7**). In contrast, PANGEA detected 30
454 candidate genes, but none were immunoglobulin genes (**Supplementary Table S8**). These data
455 further support HYENA as the algorithm with the best sensitivity among the three algorithms.

456 Third, to evaluate the specificity of the algorithms, we ran each algorithm on 20 datasets
457 generated by randomly shuffling gene expression data in both MALY and breast cancer (BRCA).
458 Since these gene expression data were random, there should be no associations between SVs and
459 gene expression, and all genes detected should be false positives. In malignant lymphoma with
460 observed gene expression, HYENA, CESAM, and PANGEA detected 16, 11, and 30 candidate
461 genes respectively (**Supplementary Tables S6, S7 and S8**). In the 20 random gene expression
462 datasets for malignant lymphoma, HYENA detected an average of 0.55 genes per dataset (**Fig.**
463 **2C**), and CESAM detected an average of 0.5 genes per dataset, whereas PANGEA detected an
464 average of 40 genes per dataset (**Supplementary Fig. S2**). In breast cancer with observed gene
465 expression, HYENA, CESAM, and PANGEA detected 7, 9, and 2,309 candidate genes,
466 respectively (**Supplementary Tables S6, S7 and S8**). In 20 random gene expression datasets for
467 breast cancer, HYENA, CESAM, and PANGEA detected 0.45, 0.9 and 2,296 genes on average
468 (**Fig. 2C** and **Supplementary Fig. S2**). In both tumor types, the numbers of false positives called
469 by PANGEA in random datasets were comparable to the numbers of genes detected with
470 observed gene expression (**Supplementary Fig. S2**). In summary, HYENA predicted the least
471 number of false positives among the three algorithms.

472 Overall, HYENA has superior sensitivity and specificity in the detection of enhancer hijacking
473 genes. Although the performances of CESAM were similar to HYENA, the genes detected by
474 HYENA and CESAM in the six benchmarking tumor types had little overlap (**Fig. 2B**). We
475 performed extensive validation on one gene detected only by HYENA.

476 **Enhancer hijacking candidate genes in the PCAWG**

477 We used HYENA to analyze a total of 1,146 tumors across 25 tumor types in the PCAWG with
478 both WGS and RNA-Seq data. When each tumor type was analyzed individually, we identified
479 108 candidate enhancer hijacking target genes in total (**Supplementary Tables S1 and S6**), four
480 of which were known enhancer hijacking targets (**Fig. 3A**). *TERT* was detected in kidney cancers
481 both from the US cohort (KICH) and the European cohort (RECA) which further demonstrated
482 the reproducibility of HYENA. All other candidate genes were only detected in one tumor type,
483 highlighting high tumor type specificity of the findings. The number of genes detected in each
484 tumor type also differed dramatically (**Fig. 3B**). No genes were detected in bladder cancer
485 (BLCA), cervical cancer (CESC), glioblastoma multiforme (GBM), or low-grade glioma (LGG),
486 probably due to their small sample sizes. Pancreatic cancer (PACA) had the greatest number of
487 candidate genes. There were two liver cancer cohorts with comparable sample sizes—LIHC
488 from the US and LIRI from Japan. Interestingly, a total of 14 genes were identified in the US
489 cohort whereas no genes were found in the Japanese cohort. One possible reason for such a
490 drastic difference could be that hepatitis B virus (HBV) infection is more common in liver cancer
491 in Japan⁴⁰, and virus integration into the tumor genome can result in oncogene activation⁴¹. In
492 Chronic Lymphocytic Leukemia (CLLE), a total of six genes were detected, and three were
493 immunoglobulin genes from both the lambda and kappa clusters (**Supplementary Table S6**).

494 Given that sample size and genome instability can only explain a small fraction of the variations
495 of enhancer hijacking target genes detected in different tumor types, the landscape of enhancer
496 hijacking in cancer seems to be mainly driven by the underlying disease biology. Intriguingly,
497 out of the 108 candidate genes, 54 (50%) were non-coding genes including lncRNAs and
498 microRNAs (**Fig. 3B**).

499 **Neo-TADs formed through somatic SVs**

500 Next, we focused on the most frequently altered candidate non-coding enhancer-hijacking target
501 gene in pancreatic cancer: *TOBI-ASI* (**Fig. 4A**), a lncRNA. *TOBI-ASI* was not detected as a
502 candidate gene by either CESAM (**Supplementary Table S7**) or PANGEA (**Supplementary**
503 **Table S8**) using the same input data. Seven (9.6%) out of 74 tumors had some form of somatic
504 SVs near *TOBI-ASI* including translocations, deletions, inversions, and tandem duplications
505 (**Fig. 4B** and **Supplementary Table S9**). For example, tumor 9ebac79d-8b38-4469-837e-
506 b834725fe6d5 had a translocation between chromosomes 17 and 19 (**Fig. 4C**). The breakpoints
507 were upstream of *TOBI-ASI* and upstream of *UQCRFS1* (**Fig. 4D**). In tumor 748d3ff3-8699-
508 4519-8e0f-26b6a0581bff, there was a 19.3 Mb deletion which brought *TOBI-ASI* next to a
509 region downstream of *KCNJ2* (**Fig. 4C** and **4E**).

510 We used Akita²³, a convolutional neural network that predicts 3D genome organization, to
511 assess the 3D architecture of the loci impacted by SVs. While 3D structures are dynamic and
512 may change with cell-type and gene activity, TAD boundaries are often more stable and remain
513 similar across different cell-types¹. TAD boundaries are defined locally by the presence of
514 binding sites for CTCF, a ubiquitously expressed DNA-binding protein^{1,26}, and TAD formation
515 arises from the stalling of the cohesin-extruded chromatin loop by DNA-bound CTCF at these
516 positions⁴². For this reason, it is expected that upon chromosomal rearrangements, normal TADs
517 can be disrupted, and new TADs can form by relocation of TAD boundaries. This assumption
518 has been validated with direct experimental evidence from examining the “neo-TADs”
519 associated with SVs at different loci⁴³⁻⁴⁵. The wildtype *TOBI-ASI* locus had a TAD between a
520 CTCF binding site in *RSAD1* and another one upstream of *SPAG9* (**Fig. 4D** and **Supplementary**
521 **Fig. S3**). There were TADs spanning *UQCRFS1* and downstream of *KCNJ2* in the two partner
522 regions (**Fig. 4D, 4E** and **Supplementary Fig. S3**). In tumor 9ebac79d-8b38-4469-837e-
523 b834725fe6d5, the translocation was predicted to lead to a neo-TAD resulting from merging the
524 TADs of *TOBI-ASI* and *UQCRFS1* (**Fig. 4D**). In tumor 748d3ff3-8699-4519-8e0f-
525 26b6a0581bff, another neo-TAD was predicted to form as a result of the deletion that merged the
526 TADs of *TOBI-ASI* and the downstream portion of *KCNJ2* (**Fig. 4E**). In both cases, within these
527 predicted neo-TADs, Akita predicted strong chromatin interactions involving several CTCF
528 binding sites and H3K27Ac peaks between *TOBI-ASI* and its two SV partners (**Fig. 4D** and **4E**
529 black arrows in the right panels), indicating newly formed promoter-enhancer interactions. In the
530 vicinity of the *TOBI-ASI* locus, *TOBI-ASI* was the only gene with significant changes in gene
531 expression. Similar neo-TADs could be observed in two additional tumors (**Supplementary Fig.**
532 **S4**). In two tumors harboring tandem duplications of *TOBI-ASI* of 317 kb and 226 kb, the
533 *TOBI-ASI* TADs were expanded (**Supplementary Fig. S5A**). However, not all SVs near *TOBI-*
534 *ASI* led to alterations in TAD architecture; for example, in tumor a3edc9cc-f54a-4459-a5d0-

535 097879c811e5, *TOBI-ASI* was predicted to remain in its original TAD after a 4 Mb tandem
536 duplication (**Supplementary Fig. S5B**). In summary, at least four out of the seven tumors
537 harboring somatic SVs near *TOBI-ASI* were predicted to result in neo-TADs including *TOBI-*
538 *ASI*. We then used another deep-learning algorithm called Orca²⁵ to predict 3D genome
539 structure based on DNA sequences. Orca-predicted 3D genome architectures were very similar
540 to Akita predictions (**Supplementary Fig. S6**) in neo-TAD formation due to SVs in the *TOBI-*
541 *ASI* locus.

542 To further study the 3D genome structure of the *TOBI-ASI* locus, we performed high-resolution
543 in situ Hi-C sequencing for four pancreatic cancer cell lines. Among these, two cell lines (Panc
544 10.05 and PATU-8988S) had high expression of *TOBI-ASI*, whereas the other two (PANC-1
545 and PATU-8988T) had low expression (**Fig. 5A**). At mega-base-pair scale, three cell lines (Panc
546 10.05, PATU-8988S, and PATU-8988T) carried several SVs (black arrows in **Fig. 5B**). In Panc
547 10.05, a tandem duplication (chr17:43,145,000-45,950,000) was observed upstream of *TOBI-*
548 *ASI* (**Fig. 5B** black arrow in the left most panel and **Supplementary Table S10**). However, the
549 breakpoint was too far away (2 Mb) from *TOBI-ASI* (chr17:48,944,040-48,945,732) and
550 unlikely to regulate its expression. A neo chromatin loop was detected by NeoLoopFinder²⁰ near
551 *TOBI-ASI* (chr17:34,010,000-48,980,000) driven by a deletion (chr17:34,460,000-47,450,000)
552 detected by EagleC²⁷ (**Supplementary Fig. S7A, Supplementary Tables S5 and S10**). The
553 deletion breakpoint was also too far away (1.5 Mb) from *TOBI-ASI* and unlikely to regulate its
554 expression. No other SVs or neo chromatin loops were detected near *TOBI-ASI*
555 (**Supplementary Tables S5 and S10**). Interestingly, there was a CNV breakpoint
556 (chr17:48,980,000) 36 kb downstream of *TOBI-ASI* in Panc 10.05 (**Fig. 5C** left most panel)
557 which was also the boundary of the neo chromatin loop. In the high copy region (upstream of the
558 CNV breakpoint), heterozygous SNPs were present with allele ratios of approximately 4:1
559 (**Supplementary Fig. S8A**), whereas in the low copy region (downstream of the CNV
560 breakpoint), all SNPs were homozygous (**Supplementary Fig. S8B**). These data suggested that
561 the DNA copy number changed from five copies to one copy at the CNV breakpoint. The gained
562 copies must connect to some DNA sequences since there should not be any free DNA ends other
563 than telomeres. Given that no off-diagonal 3D genome interactions were observed at
564 chr17:48,980,000, we considered the possibilities that the high copy region was connected to
565 repetitive sequences or to sequences that were not present in the reference genome. If so, reads
566 mapped to the high copy region should have an excessive amount of non-uniquely mapped mates
567 or unmapped mates. However, this was not the case (**Supplementary Fig. S9**). The only possible
568 configuration was a foldback inversion in which two identical DNA fragments from the copy
569 gain region were connected head to tail (**Fig. 5D** bottom left panel). As a result, in Panc 10.05,
570 there was a wildtype chromosome 17, two foldback-inversion-derived chromosomes, and a
571 translocation-derived chromosome (**Fig. 5D** bottom left panel and **Supplementary Fig. S7B**).
572 Foldback inversions are very common in cancer. If DNA double strand breaks are not
573 immediately repaired, following replication, the two broken ends of sister chromatids can self-
574 ligate head to tail and sometimes result in dicentric chromosomes^{46,47}. Algorithms, such as hic-
575 breakfinder⁴⁸ and EagleC²⁷, rely on off-diagonal 3D genomic interactions in Hi-C contact
576 matrix to detect SVs. However, foldback inversions do not form any off-diagonal interactions
577 since the two connected DNA fragments have the same coordinates, so they are not detectable by

578 existing algorithms. The 3D genome structure of the *TOBI-ASI* locus in Panc 10.05 was quite
579 distinct from the other three cell lines (**Fig. 5C**). The region immediately involved in the
580 foldback inversion had homogeneous 3D interactions (**Fig. 5C** dashed blue triangle in the left
581 most panel) suggesting that a neo-subdomain was formed (**Fig. 5D** right panel). The high
582 expression of *TOBI-ASI* in Panc 10.05 was likely a combined effect of the copy gain and the
583 neo-subdomain. In PATU-8988S and PATU-8988T, a shared SV (chr17:48,880,000-52,520,000)
584 near *TOBI-ASI* was detected (**Fig. 5B** two right panels) since the two cell lines were derived
585 from the same pancreatic cancer patient⁴⁹. This shared SV could not regulate *TOBI-ASI* because
586 it pointed away from *TOBI-ASI* (**Supplementary Fig. S10**). No other SVs were found near
587 *TOBI-ASI* in these two cell lines. The high expression of *TOBI-ASI* in PATU-8988S was likely
588 due to transcriptional regulation since the promoter of *TOBI-ASI* in PATU-8988S was more
589 accessible than that in PATU-8988T (**Fig. 5E**). This result was consistent with a handful of
590 patient tumors that had high expression of *TOBI-ASI* without any SVs (**Fig. 4A**).

591 Taken together, our results demonstrated that HYENA can detect genes activated by
592 reorganization of 3D genome architecture.

593 **Oncogenic functions of *TOBI-ASI***

594 *TOBI-ASI* has been reported as a tumor suppressor in several tumor types^{50,51}. However,
595 HYENA predicted it to be an oncogene in pancreatic cancers. To test the potential oncogenic
596 functions of *TOBI-ASI* in pancreatic cancer, we performed both in vitro and in vivo
597 experiments. We surveyed pancreatic cancer cell line RNA-Seq data from Cancer Cell Line
598 Encyclopedia (CCLE) and identified that the commonly transcribed isoform of *TOBI-ASI* in
599 pancreatic cancers was ENST00000416263.3 (**Supplementary Fig. S11**). The synthesized
600 *TOBI-ASI* cDNA was cloned and overexpressed in two pancreatic cancer cell lines, PANC-1
601 and PATU-8988T, both of which had low expression of *TOBI-ASI* (**Fig. 5A** and
602 **Supplementary Fig. S12A**). In both cell lines, overexpression of *TOBI-ASI* (**Fig. 6A**) promoted
603 in vitro cell invasion (**Fig. 6B**). In addition, three weeks after tail vein injection, PANC-1 cells
604 with *TOBI-ASI* overexpression caused higher metastatic burden in immunodeficient mice than
605 the control cells (**Fig. 6C**). Six weeks after orthotopic injection, mice carrying *TOBI-ASI*
606 overexpressing PANC-1 cells showed exacerbated overall tumor burden (**Fig. 6D**), elevated
607 primary tumor burden, and elevated metastatic burden in the spleen (**Fig. 6E** and
608 **Supplementary Fig. S12B**). Liver metastasis was not affected (**Supplementary Fig. S12C**). In
609 addition, we knocked down *TOBI-ASI* in two other pancreatic cancer cell lines Panc 10.05 and
610 PATU-8988S, both of which had high expression of *TOBI-ASI* (**Fig. 5A** and **Supplementary**
611 **Fig. S12A**), using two antisense oligonucleotides (ASOs) (**Fig. 6F**). *TOBI-ASI* expression was
612 reduced by approximately 50% by both ASOs (**Fig. 6G**). Knockdown of *TOBI-ASI* substantially
613 suppressed cell invasion in vitro (**Fig. 6H**). Note that PATU-8988T and PATU-8988S were
614 derived from the same liver metastasis of a pancreatic cancer patient, and they had drastic
615 differences in *TOBI-ASI* expression (**Fig. 5A** and **Supplementary Fig. S12A**). It was reported
616 that PATU-8988S can form lung metastases in vivo with tail vein injection of nude mice,
617 whereas PATU-8988T cannot form any metastases in any organ⁴⁹. By altering the expression of
618 *TOBI-ASI*, we were able to reverse the cell invasion phenotypes in these two cell lines (**Fig. 6B**

619 and **6H**). These results suggested that *TOBI-ASI* has an important function in regulating cell
620 invasion.

621 It is possible that *TOBI-ASI*, as an anti-sense lncRNA, transcriptionally regulates the expression
622 of the sense protein-coding gene *TOBI*. However, we did not find consistent correlations
623 between *TOBI-ASI* and *TOBI* expression in different pancreatic cancer cohorts and pancreatic
624 cancer cell lines (**Supplementary Fig. S12D**). Hence, it is unlikely that *TOBI-ASI* functions
625 through transcriptional regulation of *TOBI*. Although knocking down *TOBI-ASI* resulted in
626 down regulation of *TOBI* expression, this is an expected result given that the ASOs also targeted
627 the introns of *TOBI* (**Fig. 6F**). The decrease in *TOBI* expression was relatively mild at 10-20%
628 (**Fig. 6G**). Overexpression of *TOBI-ASI* did not have a major impact on *TOBI* expression (**Fig.**
629 **6A**). Therefore, the oncogenic functions of *TOBI-ASI* that we observed in vitro and in vivo are
630 likely independent of *TOBI*. To gain further insights into the pathway that *TOBI-ASI* is involved
631 in and its downstream targets, we performed RNA-Seq on PANC-1-generated mouse tumors
632 with *TOBI-ASI* overexpression and found that the most significantly differentially expressed
633 gene was *CNNM1* (**Supplementary Fig. S12E**). *CNNM1* is a cyclin and CBS domain divalent
634 metal cation transport mediator and is predicted to be involved in ion transport⁵². How *TOBI-*
635 *ASI* promotes cell invasion and tumor metastasis and whether *CNNM1* plays a role require
636 further study.

637 Our results showed that the lncRNA *TOBI-ASI* is oncogenic and has a pro-metastatic function in
638 pancreatic cancer, and that HYENA is able to detect novel proto-oncogenes activated by distal
639 enhancers.

640

641 Discussion

642 Here, we report a computational algorithm HYENA to detect candidate oncogenes activated by
643 distal enhancers via somatic SVs. These SV breakpoints fell in the regulatory regions of the
644 genome and caused a shuffling of regulatory elements, altering gene expression. The candidate
645 genes we detected were not limited to protein-coding genes but also included non-coding genes.
646 Our in vitro and in vivo experiments showed that a lncRNA identified by HYENA, *TOBI-ASI*,
647 was a potent oncogene in pancreatic cancers.

648 HYENA detects candidate genes based on patient cohorts rather than individual samples. Genes
649 need to be recurrently rearranged in the cohort to be detectable, and HYENA aims to identify
650 oncogenes recurrently activated by somatic SVs since these events are under positive selection.
651 Therefore, sample size is a major limiting factor. Of the eight ground truth cases, HYENA only
652 detected four (**Fig. 2A**); undetected genes were likely due to the small sample size. However,
653 genes detected in individual tumors by tools such as cis-X and NeoLoopFinder may not be
654 oncogenes, and recurrent events would be required to identify candidate oncogenes.

655 The candidate genes identified by HYENA have statistically significant associations between
656 nearby somatic SVs and elevated expression. However, the relationship may not be causal. It is
657 possible that the presence of SVs and gene expression are unrelated, but both are associated with
658 another factor. We modeled other factors to the best of our ability including gene dosage, tumor
659 purity, patient sex, age, and principal components of gene expression. In addition, it is also
660 possible that the high gene expression caused somatic SVs. Open chromatin and double helix
661 regions unwound during transcription are prone to double-strand DNA breaks which may
662 produce somatic SVs. Therefore, it is possible that some of the candidate genes are not
663 oncogenes. Functional studies are required to determine the disease relevance of the candidate
664 genes.

665 Note that the predicted 3D genome organization is not cell-type-specific. Akita was trained on
666 five high quality Hi-C and Micro-C datasets (HFF, H1hESC, GM12878, IMR90 and HCT116)²³
667 and predicts limited cell-type-specific differences. Therefore, the predicted TADs reflect
668 conserved 3D genome structure in the five cell types (foreskin fibroblast, embryonic stem cell,
669 B-lymphocyte, lung fibroblast and colon cancer). There were minor differences between HFF
670 and H1hESC (**Supplementary Fig. S3**) in genome organization. For example, the left boundary
671 of the TAD at the *UQCRFS1* locus was different between HFF and H1hESC (**Supplementary**
672 **Fig. S3A**). Nonetheless, the translocation between chromosomes 17 and 19 removed the left
673 boundary and merged the right side of the *UQCRFS1* TAD with the *TOBI-ASI* TAD (**Fig. 4D**).
674 Therefore, the cell-type difference likely does not have a major impact on our results.

675

676 **Acknowledgements**

677 We thank the Center for Research Informatics at the University of Chicago for providing the
678 computing infrastructure, Matthew Stephens for helpful suggestions, Marsha Rosner for
679 assistance in lentiviral experiments, and Ani Solanki for assistance in animal experiments. The
680 work was supported by the Goldblatt Endowment (A.Y.), the National Institutes of Health grant
681 K22CA193848 (L.Y.), R01CA269977 (L.Y.) and University of Chicago and UChicago
682 Comprehensive Cancer Center (L.Y.).

683 **Disclosure**

684 The authors have no competing interests to declare.

685

686 **Figure Legends**

687 **Figure 1. Outline of enhancer hijacking and HYENA algorithm.** **A**, Mechanisms of gene
688 activation by SVs. SVs can activate genes by recruiting distal active enhancers (top panel) and
689 by removing TAD boundaries and forming de novo enhancer-promoter interactions (bottom
690 panel). **B**, HYENA workflow. Green and purple boxes denote input and output files,
691 respectively. Orange boxes denote intermediate steps. Numbers in parentheses represent default
692 values of HYENA.

693 **Figure 2. Benchmarking HYENA.** **A**, Comparison of HYENA, CESAM, and PANGEA in
694 detecting oncogenes known to be activated by enhancer hijacking in six tumor types from the
695 PCAWG cohort. **B**, UPSET plot demonstrating candidate genes identified and shared among the
696 three tools in five tumor types of PCAWG. The numbers of candidate genes predicted by three
697 algorithms are shown on the bottom left (19, 25, and 255). On the bottom right, individual dots
698 denote genes detected by one tool, and dots connected by lines denote genes detected by multiple
699 tools. The numbers of genes detected are shown above the dots and lines. For example, the dot
700 immediately on the right of “PANGEA” shows there are 254 candidate genes detected only by
701 PANGEA but not CESAM and HYENA. The left most line connecting two dots indicates that
702 there are six genes detected by both CESAM and HYENA but not by PANGEA. **C**, Number of
703 genes detected by HYENA in two PCAWG tumor types using observed gene expression and
704 randomized expression. Genes detected in random expression datasets are false positives.

705 **Figure 3. Enhancer hijacking candidate genes in PCAWG.** **A**, Candidate genes detected by
706 HYENA in individual tumor types of PCAWG. *TERT* is plotted twice since it is detected in two
707 cancer types. Genes labelled as red are known enhancer hijacking targets. **B**, Diverse types of
708 candidate genes identified by HYENA in PCAWG. Numbers after tumor type names denote
709 sample size in the corresponding tumor types.

710 **Figure 4. *TOBI-AS1* activated by various types of SVs in pancreatic cancer.** **A**, Normalized
711 expression of *TOBI-AS1* in samples with (n=7) and without (n=67) nearby SVs in pancreatic
712 cancers. The boxplot shows median values (thick black lines), upper and lower quartiles (boxes),
713 and 1.5× interquartile range (whiskers). Individual tumors are shown as black dots. **B**, Circos
714 plot summarizing intrachromosomal SVs (blue, n=5) and translocations (red, n=3) near *TOBI-*
715 *AS1*. **C**, Diagrams depicting putative enhancer hijacking mechanisms that activate *TOBI-AS1* in
716 one tumor with a 17:19 translocation (left panel) and another tumor with a large deletion (right
717 panel). **D**, Predicted 3D chromatin interaction maps of *TOBI-AS1* (left panel), *UQCRFS1*
718 (middle panel), and the translocated region in tumor 9ebac79d-8b38-4469-837e-b834725fe6d5
719 (right panel). The downstream fragment of the chromosome 19 SV breakpoint was flipped in
720 orientation and linked to chromosome 17. H3K27Ac and CTCF ChIP-Seq data of PANC-1 cell
721 line are shown at the bottom. The expected level of 3D contacts depends on linear distance
722 between two genomic locations. Longer distances correlate with fewer contacts. Akita predicts
723 3D contacts based on DNA sequences. The heatmaps are showing the ratio between predicted
724 and expected contacts. The darkest red represent regions having 100 times more contacts than
725 expected given the distance between the regions. **E**, Predicted 3D chromatin interaction maps of
726 *TOBI-AS1* (left panel) and *KCNJ2* (middle panel) loci without deletion as well as the same
727 region following deletion in tumor 748d3ff3-8699-4519-8e0f-26b6a0581bff (right panel).

728

729 **Figure 5. 3D genome structures in the *TOBI-ASI* locus in pancreatic cancer cell lines.** **A**,
730 *TOBI-ASI* expression in pancreatic cancer cell lines in CCLE. The cell lines in red are selected
731 for further studies. **B** and **C**, 3D genomic interactions in four pancreatic cancer cell lines. Black
732 arrows represent SVs with off-diagonal interactions. The locations of *TOBI-ASI* are marked by
733 blue lines. In Panc 10.05, the blue arrow points to the CNV breakpoint and the dashed blue
734 triangle represents the neo-subdomain formed due to the foldback inversion. **D**, The reference
735 chromosome 17 and derived chromosomes in Panc 10.05. The chromosomes are not to scale.
736 *TOBI-ASI* is shown as small blue boxes in the chromosomes. **E**, Open chromatin measured by
737 ATAC-Seq in PATU-8988S and PATU-8988T at the *TOBI-ASI* locus.

738 **Figure 6. *TOBI-ASI* promotes cell invasion and tumor metastasis.** **A**, *TOBI-ASI* and *TOBI*
739 relative expression levels in PATU-8988T and PANC-1 cells transduced with *TOBI-ASI*
740 overexpression vector (n=3) or control vector (n=3). **B**, *TOBI-ASI* overexpression in PATU-
741 8988T (4 biological replicates) and PANC-1 (3 biological replicates) promoted in vitro cell
742 invasion using Transwell assay. Each biological replicate was an independent experiment with 7
743 technical replicates per experimental group. The average fold change of cell invasion was
744 calculated after the background invasion measured in the absence of any chemotactic agent was
745 subtracted from each technical replicate. *P* values were calculated by two-sided student t test. **C**,
746 *TOBI-ASI* overexpression in PANC-1 cells promoted in vivo tumor metastasis in the tail vein
747 injection model. **D**, *TOBI-ASI* overexpression in PANC-1 cells exacerbated in vivo tumor
748 growth and spontaneous metastasis in the orthotopic tumor model. Images of radiance in
749 immunodeficient mice are shown on the left while the quantifications of radiance are shown on
750 the right. Eight mice were used in both the overexpression group and the empty vector control.
751 The images were analyzed by setting the regions of interest (ROIs) to mouse torsos and
752 measuring the average radiance level (in p/sec/cm²/sr). **E**, Primary tumor burden and spleen
753 metastatic burden were higher in the mice that were orthotopically injected with *TOBI-ASI*
754 overexpression PANC-1 cells. The bar plots show quantified total radiance with a set area (in
755 p/sec). **F**, Targeting *TOBI-ASI* by two ASOs. **G**, *TOBI-ASI* knockdown in Panc 10.05 and
756 PATU-8988S cells transduced with ASO1 (n=3), ASO2 (n=3) or non-targeting control ASO
757 (NC) (n=3). **H**, *TOBI-ASI* knockdown suppressed Panc 10.05 (3 biological replicates) and
758 PATU-8988S (3 biological replicates) cell invasion in vitro. Cell invasion fold change
759 calculation is the same as in **B**. Two-sided student t test was used. Error bars in all panels
760 indicate standard error of the mean.

761

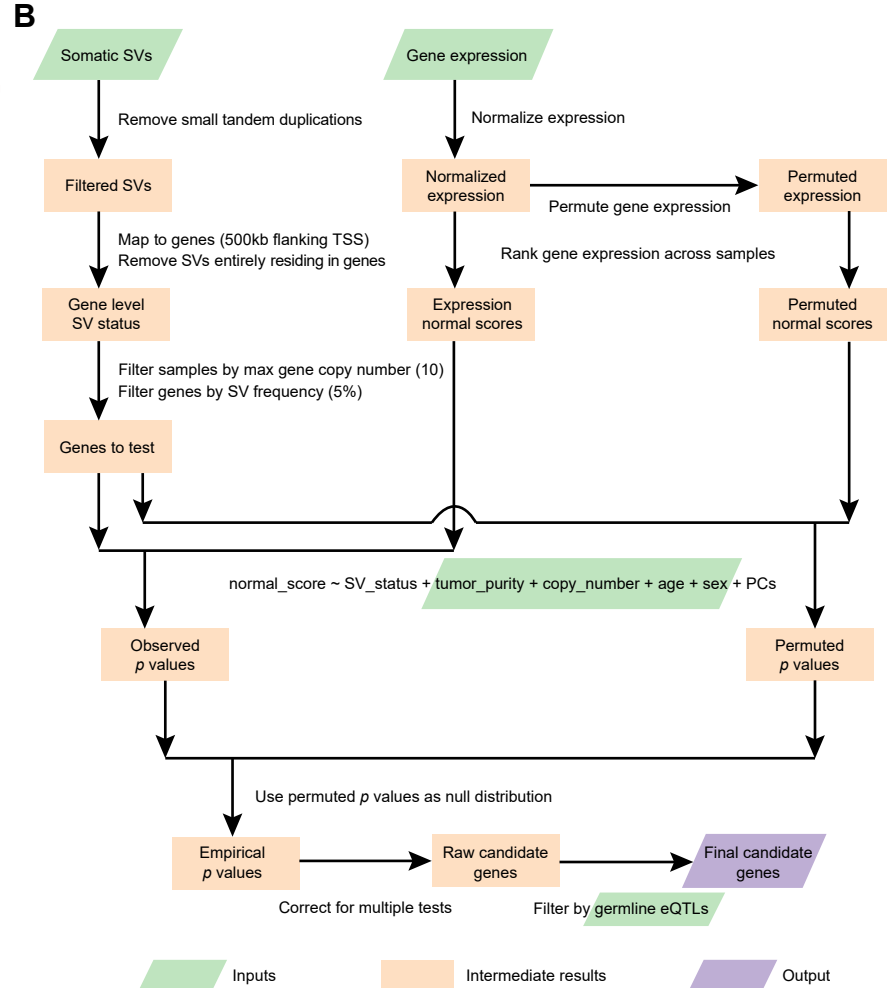
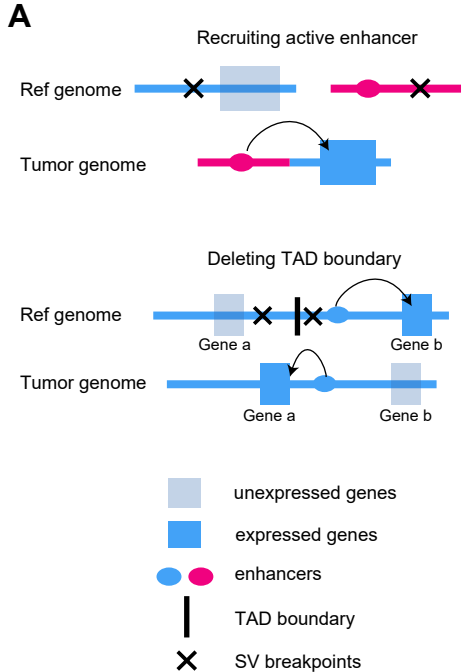
762 **References**

- 763 1. Dixon, J. R. *et al.* Topological domains in mammalian genomes identified by analysis of
764 chromatin interactions. *Nature* **485**, 376–380 (2012).
- 765 2. Krijger, P. H. L. & De Laat, W. Regulation of disease-associated gene expression in the
766 3D genome. *Nat. Rev. Mol. Cell Biol.* **17**, 771–782 (2016).
- 767 3. Symmons, O. *et al.* Functional and topological characteristics of mammalian regulatory
768 domains. *Genome Res.* **24**, 390–400 (2014).
- 769 4. Lupiáñez, D. G. *et al.* Disruptions of topological chromatin domains cause pathogenic
770 rewiring of gene-enhancer interactions. *Cell* **161**, 1012–1025 (2015).
- 771 5. Zhang, W. *et al.* A global transcriptional network connecting noncoding mutations to
772 changes in tumor gene expression. *Nat. Genet.* *2018* **50**, 613–620 (2018).
- 773 6. Weischenfeldt, J., Symmons, O., Spitz, F. & Korbel, J. O. Phenotypic impact of genomic
774 structural variation: insights from and for human disease. *Nat. Rev. Genet.* **14**, 125–138
775 (2013).
- 776 7. Davis, C. F. *et al.* The somatic genomic landscape of chromophobe renal cell carcinoma.
777 *Cancer Cell* **26**, 319–330 (2014).
- 778 8. Bakhshi, A. *et al.* Cloning the chromosomal breakpoint of t(14;18) human lymphomas:
779 clustering around Jh on chromosome 14 and near a transcriptional unit on 18. *Cell* **41**,
780 899–906 (1985).
- 781 9. Gostissa, M. *et al.* Long-range oncogenic activation of Igh-c-myc translocations by the
782 Igh 3' regulatory region. *Nature* **462**, 803–807 (2009).
- 783 10. Hnisz, D. *et al.* Activation of proto-oncogenes by disruption of chromosome
784 neighborhoods. *Science* **351**, 1454–1458 (2016).
- 785 11. Gröschel, S. *et al.* A single oncogenic enhancer rearrangement causes concomitant EVI1
786 and GATA2 deregulation in leukemia. *Cell* **157**, 369–381 (2014).
- 787 12. Northcott, P. A. *et al.* Enhancer hijacking activates GFII family oncogenes in
788 medulloblastoma. *Nature* **511**, 428–434 (2014).
- 789 13. Weischenfeldt, J. *et al.* Pan-cancer analysis of somatic copy-number alterations implicates
790 IRS4 and IGF2 in enhancer hijacking. *Nat. Genet.* **49**, 65–74 (2016).
- 791 14. Northcott, P. A. *et al.* The whole-genome landscape of medulloblastoma subtypes. *Nature*
792 **547**, 311–317 (2017).
- 793 15. He, B. *et al.* Diverse noncoding mutations contribute to deregulation of cis-regulatory
794 landscape in pediatric cancers. *Sci. Adv.* **6**, eaba3064 (2020).
- 795 16. Kopp, F. & Mendell, J. T. Functional Classification and Experimental Dissection of Long
796 Noncoding RNAs. *Cell* **172**, 393–407 (2018).
- 797 17. Lin, C.-P. & He, L. Noncoding RNAs in Cancer Development. *Annu. Rev. Cancer Biol.* **1**,
798 163–184 (2017).

- 799 18. Liu, S. J., Dang, H. X., Lim, D. A., Feng, F. Y. & Maher, C. A. Long noncoding RNAs in
800 cancer metastasis. *Nat. Rev. Cancer* 2021 217 **21**, 446–460 (2021).
- 801 19. Liu, Y. *et al.* Discovery of regulatory noncoding variants in individual cancer genomes by
802 using cis-X. *Nat. Genet.* **52**, 811–818 (2020).
- 803 20. Wang, X. *et al.* Genome-wide detection of enhancer-hijacking events from chromatin
804 interaction data in rearranged genomes. *Nat. Methods* **18**, 661–668 (2021).
- 805 21. Campbell, P. J. *et al.* Pan-cancer analysis of whole genomes. *Nature* **578**, 82–93 (2020).
- 806 22. Aguet, F. *et al.* Genetic effects on gene expression across human tissues. *Nat.* 2017
807 5507675 **550**, 204–213 (2017).
- 808 23. Fudenberg, G., Kelley, D. R. & Pollard, K. S. Predicting 3D genome folding from DNA
809 sequence with Akita. *Nat. Methods* **17**, 1111–1117 (2020).
- 810 24. Krietenstein, N. *et al.* Ultrastructural Details of Mammalian Chromosome Architecture.
811 *Mol. Cell* **78**, 554-565.e7 (2020).
- 812 25. Zhou, J. Sequence-based modeling of three-dimensional genome architecture from
813 kilobase to chromosome scale. *Nat. Genet.* **54**, 725–734 (2022).
- 814 26. Rao, S. S. P. *et al.* A 3D Map of the Human Genome at Kilobase Resolution Reveals
815 Principles of Chromatin Looping. *Cell* **159**, 1665–1680 (2014).
- 816 27. Wang, X., Luan, Y. & Yue, F. EagleC: A deep-learning framework for detecting a full
817 range of structural variations from bulk and single-cell contact maps. *Sci. Adv.* **8**, 9215
818 (2022).
- 819 28. Grandi, F. C., Modi, H., Kampman, L. & Corces, M. R. Chromatin accessibility profiling
820 by ATAC-seq. *Nat. Protoc.* 2022 176 **17**, 1518–1552 (2022).
- 821 29. Beroukhim, R. *et al.* The landscape of somatic copy-number alteration across human
822 cancers. *Nature* **463**, 899–905 (2010).
- 823 30. Uhlen, M. *et al.* A pathology atlas of the human cancer transcriptome. *Science* **357**,
824 eaan2507 (2017).
- 825 31. Zhang, X. *et al.* Identification of focally amplified lineage-specific super-enhancers in
826 human epithelial cancers. *Nat. Genet.* **48**, 176–182 (2015).
- 827 32. Takeda, D. Y. *et al.* A Somatically Acquired Enhancer of the Androgen Receptor Is a
828 Noncoding Driver in Advanced Prostate Cancer. *Cell* **174**, 422-432.e13 (2018).
- 829 33. Turner, K. M. *et al.* Extrachromosomal oncogene amplification drives tumour evolution
830 and genetic heterogeneity. *Nature* **543**, 122–125 (2017).
- 831 34. Wu, S. *et al.* Circular ecDNA promotes accessible chromatin and high oncogene
832 expression. *Nature* **575**, 699–703 (2019).
- 833 35. Morton, A. R. *et al.* Functional Enhancers Shape Extrachromosomal Oncogene
834 Amplifications. *Cell* **179**, 1330-1341.e13 (2019).

- 835 36. Li, Y., Ge, X., Peng, F., Li, W. & Li, J. J. Exaggerated false positives by popular
836 differential expression methods when analyzing human population samples. *Genome Biol.*
837 **23**, 79 (2022).
- 838 37. Ooi, W. F. *et al.* Integrated paired-end enhancer profiling and whole-genome sequencing
839 reveals recurrent CCNE1 and IGF2 enhancer hijacking in primary gastric
840 adenocarcinoma. *Gut* **69**, 1039–1052 (2020).
- 841 38. Yun, J. W. *et al.* Dysregulation of cancer genes by recurrent intergenic fusions. *Genome*
842 *Biol.* **21**, 166 (2020).
- 843 39. Richter, J. *et al.* Recurrent mutation of the ID3 gene in Burkitt lymphoma identified by
844 integrated genome, exome and transcriptome sequencing. *Nat. Genet.* **44**, 1316–1320
845 (2012).
- 846 40. Zapatka, M. *et al.* The landscape of viral associations in human cancers. *Nat. Genet.* 2020
847 523 **52**, 320–330 (2020).
- 848 41. Neuveut, C., Wei, Y. & Buendia, M. A. Mechanisms of HBV-related
849 hepatocarcinogenesis. *J. Hepatol.* **52**, 594–604 (2010).
- 850 42. Fudenberg, G. *et al.* Formation of Chromosomal Domains by Loop Extrusion. *Cell Rep.*
851 **15**, 2038–2049 (2016).
- 852 43. Franke, M. *et al.* Formation of new chromatin domains determines pathogenicity of
853 genomic duplications. *Nature* **538**, 265–269 (2016).
- 854 44. Melo, U. S. *et al.* Complete lung agenesis caused by complex genomic rearrangements
855 with neo-TAD formation at the SHH locus. *Hum. Genet.* **140**, 1459–1469 (2021).
- 856 45. de Bruijn, S. E. *et al.* Structural Variants Create New Topological-Associated Domains
857 and Ectopic Retinal Enhancer-Gene Contact in Dominant Retinitis Pigmentosa. *Am. J.*
858 *Hum. Genet.* **107**, 802–814 (2020).
- 859 46. Li, Y. *et al.* Constitutional and somatic rearrangement of chromosome 21 in acute
860 lymphoblastic leukaemia. *Nature* **508**, 98–102 (2014).
- 861 47. Maciejowski, J., Li, Y., Bosco, N., Campbell, P. J. & de Lange, T. Chromothripsis and
862 Kataegis Induced by Telomere Crisis. *Cell* **163**, 1641–1654 (2015).
- 863 48. Dixon, J. R. *et al.* Integrative detection and analysis of structural variation in cancer
864 genomes. *Nat. Genet.* **50**, 1388–1398 (2018).
- 865 49. Elsässer, H. P., Lehr, U., Agricola, B. & Kern, H. F. Establishment and characterisation of
866 two cell lines with different grade of differentiation derived from one primary human
867 pancreatic adenocarcinoma. *Virchows Arch. B. Cell Pathol. Incl. Mol. Pathol.* **61**, 295–
868 306 (1992).
- 869 50. Yao, J. *et al.* Long noncoding RNA TOB1-AS1, an epigenetically silenced gene,
870 functioned as a novel tumor suppressor by sponging miR-27b in cervical cancer. *Am. J.*
871 *Cancer Res.* **8**, 1483 (2018).
- 872 51. Shangguan, W. *et al.* TOB1-AS1 suppresses non-small cell lung cancer cell migration and

- 873 invasion through a ceRNA network. *Exp. Ther. Med.* **18**, (2019).
- 874 52. Wang, C. Y. *et al.* Molecular cloning and characterization of a novel gene family of four
875 ancient conserved domain proteins (ACDP). *Gene* **306**, 37–44 (2003).
- 876



A

Gene	Tumor type	HYENA	CESAM	PANGEA
<i>MYC</i>	malignant lymphoma	●	●	-
<i>BCL2</i>	malignant lymphoma	●	●	-
<i>CCNE1</i>	gastric/stomach cancer	-	-	-
<i>TERT</i>	chromophobe renal cell carcinoma	●	●	-
<i>IGF2</i>	colorectal cancer	-	-	-
<i>IGF2</i>	gastric/stomach cancer	-	-	-
<i>IGF2BP3</i>	thyroid cancer	●	-	NT
<i>IRS4</i>	lung squamous cell carcinoma	-	-	-

● detected - undetected NT not tested

