

1 **Title: Pseudogene-mediated DNA demethylation leads to oncogene activation**

2 **Authors:** Junsu Kwon¹, Yanjing V. Liu^{o1}, Chong Gao^{o2}, Mahmoud A. Bassal^{1,3}, Adrianna I.
3 Jones³, Junyu Yang², Zhiyuan Chen², Li Ying¹, Henry Yang¹, Leilei Chen¹,
4 Annalisa Di Ruscio^{4,5}, Yvonne Tay^{*1,6}, Li Chai^{*2}, Daniel G. Tenen^{*1,3,4}

5
6 **Affiliations**

7 ¹Cancer Science Institute of Singapore, National University of Singapore, 117599, Singapore

8 ²Department of Pathology, Brigham and Women's Hospital, Boston, MA 02115, USA

9 ³Harvard Stem Cell Institute, Harvard Medical School, Boston, MA 02115 USA

10 ⁴Harvard Medical School Initiative for RNA Medicine, Harvard Medical School, Boston, MA
11 02115, USA

12 ⁵Department of Translational Medicine, University of Eastern Piedmont, Novara, 28100, Italy

13 ⁶Department of Biochemistry, Yong Loo Lin School of Medicine, National University of
14 Singapore, 117599, Singapore

15 ^oThese authors contributed equally to this work as co-second authors

16 ^{*}Co-correspondence

17

18 **Running title:** Hepatitis B virus induces pseudogene to activate SALL4

19

20 **Keywords:** SALL4, hepatocellular carcinoma, hepatitis B infection, pseudogene-mediated
21 oncogene reactivation, DNA demethylation

22

23 **Co-corresponding authors:** Yvonne Tay, National University of Singapore, 14 Medical

24 Drive, Singapore, 117599. Phone: +65 6516-1160, E-mail: yvonnetay@nus.edu.sg; Li Chai,

25 Brigham and Women's Hospital, 77 Avenue Louis Pasteur, NRB630D. Phone: +1 617-840-
26 6711; Fax: 617-277-9013; E-mail: lchai@bwh.harvard.edu; Daniel G. Tenen, National
27 University of Singapore, 14 Medical Drive, Singapore, 117599. Phone: +65 6516-1160; E-
28 mail: daniel.tenen@nus.edu.sg.

29

30 **Declaration of interests:** The authors declare no competing interests.

31

32 This study was supported by Singapore Ministry of Health's National Medical Research
33 Council; National Institutes of Health; Xiu Research Fund, National Research Foundation
34 Fellowship; National University of Singapore President's Assistant Professorship; Singapore
35 Ministry of Education under its Research Centres of Excellence initiative; Singapore
36 Ministry of Education's AcRF Tier 3 grants; National Cancer Institute; Italian Association for
37 Cancer Research (AIRC).

38

39 **Word count: 3236 words for the abstract, statement of significance, introduction,**
40 **results, and discussion**

41

42 **Total number of figures: 4**

43

44 **Total number of supplementary figures/tables: 4/5**

45

46

47 **Abstract**

48 Despite being one of the leading causes of cancer-related deaths, there is an unmet clinical
49 need for hepatocellular carcinoma (HCC) patients. The lack of effective treatment is, at least
50 in part, due to our lack of understanding of the molecular pathogenesis of this disease.
51 Oncofetal protein SALL4 is re-activated in patients with aggressive HCC along with other
52 solid tumors and hematologic malignancies. This study identifies a previously unrecognized
53 mechanism of SALL4 reactivation which is mediated by pseudogene-induced demethylation.
54 Using a locus-specific demethylating technology, we identified the critical CpG region for
55 SALL4 expression. We showed that SALL4 pseudogene 5 hypomethylates this region
56 through interaction with DNMT1, resulting in SALL4 upregulation. Intriguingly,
57 pseudogene 5 is significantly upregulated in a hepatitis B virus (HBV) model prior to SALL4
58 induction, and both are increased in HBV-HCC patients. Our results suggest that pseudogene-
59 mediated demethylation represents a unique mechanism of oncogene activation in cancer.

60

61 **Significance**

62 Our study provides a mechanistic link between HBV infection, activation of the oncogene
63 SALL4, and HCC. We reveal a previously undescribed capability of a pseudogene to
64 epigenetically activate an oncogene by demethylation in a locus-specific manner.

65

66 **Introduction**

67 Hepatocellular carcinoma (HCC) is one of leading causes of cancer-related deaths
68 globally, with more than 700,000 new cases and 600,000 estimated HCC deaths each year.
69 Hepatitis B virus infection is one of the main causes of HCC, particularly in Asia. While
70 surgery, liver transplantation, or radiological intervention may be a viable option for early
71 stage disease, prognosis for advanced stage HCC remains bleak, with most patients
72 eventually dying within 20 months after diagnosis. Sorafenib, an oral multikinase inhibitor, is
73 the one of the few approved agents for patients with advanced HCC (1,2). However, the
74 effectiveness of Sorafenib for advanced HCC is debatable (2). There is an unmet clinical
75 need for the development of more effective therapies for the treatment of HCC. The lack of
76 effective treatment options for HCC is at least in part due to our lack of understanding the
77 pathogenesis of this disease. Identifying novel pathway(s) that are responsible for HCC could
78 be translated into targeted therapy and improve the outcomes of these patients.

79 SALL4 is a potent stem cell factor for self-renewal and pluripotency of embryonic stem
80 cells (ESCs) (3,4). During development, SALL4 expression diminishes gradually and is
81 eventually silenced in most normal tissues. Strikingly, high SALL4 expression levels have
82 been observed in many malignancies such as liver cancer, acute myeloid leukemia, breast
83 cancer and lung cancer (5-8). Re-expression of SALL4 in cancers is associated with a more
84 aggressive cancer phenotype, drug resistance and reduced patient survival (5,6,9-11). Of note,
85 HCC patients with detectable SALL4 expression have enriched hepatic progenitor-like gene
86 signatures and poorer prognoses (10). Furthermore, targeting SALL4 in HCC cell lines by
87 knocking down or using inhibitory peptides resulted in cellular death (12), suggesting that
88 SALL4 plays a crucial role in hepatocarcinogenesis and that targeting SALL4 may provide
89 an innovative therapeutic approach for this disease. However, mechanistically, how SALL4 is
90 re-activated in HCC is still unclear, although it has been reported that aberrant methylation

91 could be a contributing factor (13,14). By defining the mechanism of SALL4 reactivation in
92 HCC, we can better treat HCC.

93 DNA methylation is a frequently studied mechanism of epigenetic regulation in humans
94 that is mediated by DNA methyltransferases (DNMT); of which, DNMT1 has a structural
95 binding preference (15). Research by multiple groups including ours has demonstrated that
96 non-coding RNAs (ncRNAs) such as ecCEBP α , Dali, Dum, and Dacor1 can interact with
97 DNMT1 to inhibit its methylation activity. These ncRNAs thus indirectly alter local
98 methylation states in different cancers, acting as key tissue-specific epigenetic regulators of
99 gene expression (15-18). It was also reported that the exon 1-intron 1 region of the
100 *SALL4* gene locus is hypermethylated in non *SALL4*-expressing K562 leukemic cells.
101 Reprogramming of these cells resulted in demethylation of this region and a subsequent
102 increase in *SALL4* expression (14). Recently, a report described demethylation of specific
103 CpG sites downstream of the *SALL4* transcriptional start site (TSS) in hepatitis B (HBV)-
104 related HCC which could contribute to *SALL4* re-activation in HCC (13). However, it is still
105 unclear how HBV infection could initiate the demethylation and reactivation of specific
106 oncogenes.

107 Pseudogenes are a class of non-coding RNAs (ncRNAs) once regarded as insignificant
108 “junk” DNA relics due to their lack of coding potential. However, studies have demonstrated
109 that pseudogene transcripts can regulate gene expression of oncogenes and tumor-suppressors
110 by acting as antisense transcripts, processed small interfering RNAs (siRNAs) and competing
111 endogenous RNAs (ceRNAs) (19-21). Recent pseudogene expression analysis in over 2,800
112 patient samples showed strong concordance between pseudogene expression and tumor
113 subtypes, as well as patient prognoses, highlighting the clinical importance of pseudogenes
114 (22). Our group focused on characterizing regulatory functions of *SALL4* pseudogenes.
115 *SALL4*, a well-studied oncogene with high expression levels in several hematological

116 malignancies and solid tumors, has eight pseudogenes of different lengths varying from 500
117 nucleotides to 6,000 nucleotides, and yet there have been no studies investigating SALL4
118 pseudogenes (5-8).

119 As many pseudogenes are actively transcribed in cells, we postulated that they could
120 interact with RNA-binding proteins such as DNMT1 via highly homologous RNA motifs and
121 exert regulatory functions. We therefore tested the hypothesis as to whether pseudogenes are
122 involved in DNA methylation as DNMT1-interacting lncRNAs in an HBV-positive HCC
123 model.

124

125 **Results**

126 **SALL4 expression is negatively correlated with methylation of the**

127 **5' UTR - exon 1 - intron 1 region**

128 We hypothesized that the degree of methylation in the SALL4 locus is associated with
129 SALL4 expression, therefore, we examined HBV-positive HCC patients. Using a publicly
130 available dataset (23), the overall transcript levels and methylation status of SALL4 were
131 analysed using probes covered the entire SALL4 gene locus. A significant negative
132 correlation between SALL4 expression and methylation was observed in primary tumours at
133 the Probe 1 (Fig. 1A and 1B) which was only observed in the 5'UTR – exon 1 region. Sites
134 located either proximal or distal (Probe 2) to the 5'UTR-exon 1 locus showed poor to no
135 correlation with SALL4 expression (Fig. 1C).

136 As the methylome and transcriptome of cell lines could be different from those of primary
137 tumours, we further investigated the negative correlation between SALL4 methylation and
138 expression among HCC cell lines. The 5'UTR-exon 1-intron 1 region was first inspected and
139 found to have over 30 CpG dinucleotides (Supplementary S1). Bisulfite sequencing in the

140 HCC cell lines SNU398 and SNU387 revealed distinct and unique methylation profiles for
141 the two cell lines (Fig. 1D). Within the profiled region, SNU387 showed a near universal,
142 methylated profile in stark contrast to SNU398 which showed a completely demethylated
143 profile, with the exception of 4 CpG dinucleotides. The result was consistent with previous
144 reports observing that methylation of the SALL4 5'UTR-exon 1-intron 1 region is
145 differentially methylated in K562-induced pluripotency reprogrammed cells and HBV-related
146 HCC patients. (13,14). To investigate the relationship between the observed methylation
147 profiles and gene expression, we examined SALL4 expression in both SNU398 and SNU387
148 (Fig. 1E and 1F). The level of SALL4 transcription was substantive as more than 100 copies
149 of SALL4 mRNAs per cell were detected in SNU398, while SNU387 cells only expressed
150 about 10 copies per cell. It was also evident that SALL4 was expressed at much higher
151 magnitude in SNU398, in which the SALL4 loci was hypomethylated. Taken together, both
152 the cell line and patient data suggest that SALL4 methylation and expression are negatively
153 correlated. We therefore confirmed that DNA methylation could be a potential regulatory
154 mechanism for SALL4 expression in HCC.

155

156 **CRISPR-DiR demethylates and activates SALL4**

157 To further investigate the correlation between methylation of the SALL4
158 5' UTR- exon 1- intron 1 region and SALL4 expression, the CRISPR-DNMT1-interacting
159 RNA (CRISPR-DiR) technique was utilized to induce locus-specific demethylation by
160 blocking DNMT1 activity in SNU387 cells. Briefly, the single guide RNA (sgRNA) was
161 constructed to contain a SALL4 exon 1 targeting sequence, two RNA loops of ecCEBP α with
162 validated DNMT1 inhibitory function (15), and the dCAs9-interacting domain (Fig. 2A).
163 Numerous sgRNAs targeting the SALL4 5' UTR- exon 1- intron 1 locus were designed, with
164 the four most efficient candidates shortlisted (sgSALL4_1 - sgSALL4_4) via *in vitro* sgRNA

165 selection (Supplementary Fig. S2A). Transduced cells were also validated to express the
166 dCas9-mCherry through FACS (Supplementary Fig. S2B).

167 Methylation of the SALL4 5' UTR- exon 1- intron 1 CpG island was monitored in
168 SNU387 cells with four independent CRISPR-DiR inductions, one for each shortlisted
169 sgRNA. Of these inductions, sgSALL4_1 was the most potent sgRNA tested. Upon
170 transduction of SNU387 cells with sgSALL4_1, significant demethylation changes were
171 observed after 14 days, which continued for over 7 additional days (Fig. 2B). Conversely, no
172 change in methylation was observed in non-targeting, negative control transduced cells. To
173 examine potential off-target effects of CRISPR-DiR, we concurrently monitored methylation
174 of a region in SALL4 exon 4, and confirmed demethylation was localized to only the targeted
175 5' UTR - exon 1 - intron 1 CpG island (Supplementary Fig. S2C).

176 Following CRISPR-DiR targeted demethylation of the SALL4 5' UTR-exon 1-intron 1
177 CpG island, both SALL4 transcript and protein levels increased as predicted (Fig. 2C and
178 2D). The magnitude of SALL4 upregulation observed in these cells was comparable to that of
179 treatment with 5-aza-2'-deoxycytidine, a global demethylating agent. Furthermore, SALL4
180 expression was not significantly altered when the cells were transduced with other less
181 efficient sgRNAs (Supplementary Fig. S2D). As SALL4 overexpression promotes cancer cell
182 growth, we performed growth assays on sgSALL4-transduced SNU387 cells and observed
183 increased anchorage-independent and -dependent growth compared to negative control (Fig.
184 2E and F), suggesting that targeted demethylation of the SALL4 locus leads to upregulated
185 expression of SALL4, with concomitant enhanced cellular growth.

186

187 **SALL4P5 demethylates and activates SALL4, and associates with DNMT1**

188 There are 8 SALL4 pseudogenes, and since none are located on the same chromosome as
189 SALL4, it is unlikely that they will be transcribed as siRNAs or antisense transcripts to
190 deregulate SALL4. However, the identified pseudogenes do share high sequence homology
191 with the paralogous coding SALL4 gene. It is therefore possible that the SALL4 pseudogenes
192 could bind to other proteins with either matching DNA/RNA motifs or with comparable
193 secondary and tertiary structures owing to their high sequence homology. As previously
194 reported, ecCEBP α , a ncRNA that overlaps with and thus has regions of identity with its
195 paralog, CEBP α , can interact with DNMT1 and affect CEBP α gene expression. We therefore
196 postulated that SALL4 pseudogenes, which are highly homologous to SALL4, could
197 potentially mediate SALL4 demethylation.

198 Each SALL4 pseudogene was transiently overexpressed in SNU387 cells, and the
199 methylation profile of the 5' UTR-exon 1-intron 1 CpG island was assessed (Fig. 3A). Only
200 SALL4 pseudogene 5 (SALL4P5) overexpression resulted in a demethylation pattern
201 comparable to that seen using CRISPR-DiR. Consistently, SALL4P5 knock-down in
202 hypomethylated SNU398 cells led to increased methylation of the locus as predicted (Fig.
203 3B).

204 Additional evidence to suggest direct SALL4P5-DNMT1 interaction can be seen by their
205 matched cellular localization. Cellular localization of pseudogene transcripts is a critical
206 factor in determining their function, as they must be localized in the same cellular
207 compartment as their binding partners to exert specific biological functions. It is known that
208 DNMT1 facilitates methylation exclusively in the nucleus. Critically, SALL4P5 is also
209 primarily localized to the nucleus. This contrasts with SALL4P7, which has a predominant
210 cytoplasmic localization in SNU398 (Fig. 3C). Therefore, although SALL4P7 shares
211 sequence homology with SALL4P5 and SALL4, its primary localization in cytoplasm could,
212 in part, account for its inability to demethylate the SALL4 locus.

213 As DNMT1 is known as a maintenance DNA methylator, we therefore investigated
214 whether our observed SALL4P5 demethylation phenotype is due to a SALL4P5-DNMT1
215 interaction. As there are no known interacting SALL4P5-DNMT1 binding regions or pockets,
216 we performed an unbiased biotinylated pull-down assay using full-length SALL4P5. First, to
217 validate the efficacy of the pull-down, DNMT1 protein could successfully pull-down
218 ecCEPB α (Fig. 3D). Similarly, SALL4P5 was able to successfully pull-down DNMT1,
219 whereas SALL4P7 and the antisense negative control did not.

220 Having shown an association between SALL4 exon 1 - intron 1 demethylation and SALL4
221 expression up-regulation, we next investigated the effect of transiently overexpressing
222 SALL4P5 on SALL4 levels. SALL4P5 overexpression significantly upregulated SALL4
223 transcript (Fig. 3E) and protein levels, the latter of which was more striking and equivalent to
224 the overexpression of SALL4 itself (Fig. 3F). Interestingly, SALL4P7 overexpression also
225 increased SALL4 protein levels. Although SALL4P7 may not play a role in SALL4
226 demethylation, it could still contribute to gene regulation as homologous pseudogenes could
227 also function as ceRNAs (21) by sequestering bioavailable microRNAs that target and repress
228 SALL4. Consistent with the phenotype of elevated SALL4 levels, SALL4P5 overexpression
229 also significantly increased colony formation of SNU387 cells (Fig. 3G). The data suggest
230 that SALL4P5 could have oncogenic effects as it can directly upregulate SALL4 expression
231 and cell growth.

232

233 **SALL4P5 is upregulated in HCC patients and during hepatitis B induction**

234 The aforementioned results demonstrate that SALL4P5 upregulation can reactivate SALL4
235 expression in cell lines. We next sought to validate these findings in primary patient samples,
236 and measured SALL4P5 expression in HCC patients, who frequently have elevated levels of
237 SALL4 (10). Twenty HCC patients with paired non-disease samples were screened from a

238 Hong Kong cohort. Among these 20 patients, 19 of them were HBV positive and 7 had
239 increased SALL4 levels of over 1.5-fold (Fig. 4A). Interestingly, within these 7 patients, only
240 SALL4P5 expression was concomitantly upregulated, while SALL4P7 showed little to no
241 change. For patients, such as patient #1, with no SALL4 level change, SALL4P5 expression
242 was also unaltered.

243 HBV infection is the single most common risk factor of HCC, as more than 50% of
244 patients contract hepatitis B prior to HCC (24). We therefore sought to investigate whether
245 SALL4P5-mediated demethylation and subsequent reactivation of SALL4 during hepatitis B
246 infection could drive oncogenesis. The HepAD38B model was used, as it allows controlled
247 induction of hepatitis B virus production using the tet-off system (13). Using digital droplet
248 PCR (ddPCR), we validated that the HepAD38B cells produced the major hepatitis B viral
249 transcripts such as core, surface, hepatitis B antigen X (HBx) and polymerase transcripts
250 (Supplementary Fig. 3A). Upon hepatitis B induction, SALL4 and SALL4P5 transcript levels
251 also increased (Fig. 4B, Supplementary Fig. S3B and S3C). Interestingly, the expression
252 level for HBx increased first, SALL4P5 expression then followed at 54 hours, and lastly
253 SALL4 expression at 84 hours in step-wise manner. When performing bisulfite sequencing of
254 these critical time points, it was found that the average methylation across the CpG island in
255 the 5'UTR - exon 1 – intron 1 junction decreased upon hepatitis B induction, suggesting that
256 the infection-induced upregulation of SALL4P5 demethylates and reactivates SALL4. (Fig.
257 4C).

258

259 **Discussion**

260 In this study, we provided a mechanistic link between HBV infection, activation of the
261 oncogene SALL4, and HCC development. We, for the first time to our knowledge, showed
262 that hepatitis B viral infection could lead to pseudogene upregulation, which in turn could

263 epigenetically regulate oncogenes and drive tumorigenesis. We also demonstrated for the first
264 time that a pseudogene can associate with a DNA methyltransferase to inhibit its function and
265 subsequently influence expression of its coding paralog, oncogenic SALL4. We described a
266 strong negative correlation between SALL4 expression and exon 1-locus specific methylation
267 in HCC patients as well as in cell lines. By utilizing a novel CRISPR-DiR technology, we
268 validated the importance of exon 1 methylation to SALL4 expression and cell growth.
269 Interfering with DNMT1 activity at a specific CpG island in this region resulted in the
270 upregulation of SALL4. Furthermore, we identified and characterized SALL4 pseudogene 5
271 (SALL4P5), which shares high sequence homology with its paralogous coding gene SALL4,
272 as a critical regulator of SALL4 expression and function by interacting with DNMT1 to
273 demethylate and upregulate SALL4 expression. More importantly, we demonstrated that
274 SALL4P5 and SALL4 expression are sequentially upregulated in an HBV induction model as
275 well as positively correlated and upregulated in hepatitis B-infected HCC patients.

276 This work highlights the previously undescribed capability of a pseudogene to
277 epigenetically regulate an oncogene by interacting with DNMT1 and affecting its methylation
278 in a specific region. It is reported that there are at least 12,000 pseudogenes (25) in the human
279 genome. A recent pan-cancer analysis of pseudogene expression in different cancers
280 demonstrated that pseudogene expression alone can serve as a molecular and prognostic
281 factor for patients with different cancer subtypes, highlighting the functional and clinical
282 importance of pseudogenes (22). There may be other pseudogenes capable of mediating
283 similar DNMT1-interactions and exerting the demethylating function on oncogenes and
284 tumor-suppressors in different cancers. Monitoring expression levels of oncogene-
285 demethylating pseudogenes could enable predicting oncogene activation as well as disease
286 progression to improve patient outcomes.

287 Understanding the mechanism of pseudogene-mediated demethylation may provide
288 insights into SALL4 re-expression in HCC, which is of therapeutic value owing to SALL4's
289 prognostic significance for the disease (7). Here we elaborate how a pseudogene could play a
290 crucial role in epigenetic regulation of an important oncogene, thereby suggesting that there
291 could similarly significant and impactful pseudogenes potentially contributing to early-stage
292 gene regulation of tumor-suppressors and oncogenes. In addition, pseudogenes could exert
293 non-canonical functions by interacting with other RNA-binding proteins, with potential wide-
294 ranging implications in gene regulation and function as we have demonstrated here.

295 Previous studies from Di Ruscio et al., demonstrated that RNAs require distinct secondary
296 structures in order to associate with DNMT1 (15). This preferential interaction through
297 structure could potentially explain why DNMT1 interacts with SALL4P5, but not SALL4P7,
298 even though the two pseudogenes are highly homologous. Another striking aspect of
299 SALL4P5-SALL4 regulation is that unlike ecCEBP α , which resides in the CEBP α locus, the
300 SALL4P5 locus is on chromosome 3, while its paralogous coding gene SALL4 is on
301 chromosome 20. This *trans* regulation implies that it could be the sequence homology and
302 perhaps secondary structure that plays a more critical role than chromosomal location for a
303 ncRNA to work as a DNMT1-interacting RNA. Moreover, as DNMT1 plays a crucial role in
304 *de novo* methylation, the SALL4P5-DNMT1 interaction could contribute to SALL4
305 reactivation as well as other oncogene activation in cancers.

306 DNA hypermethylation of CpG islands in gene promoter regions has consistently
307 correlated with inactivation of tumor suppressors in cancers (26). Conversely, demethylation
308 of oncogene promoters leads to increased gene expression (27). Interestingly, it was the
309 methylation profile of the 5'UTR - Exon 1-Intron 1 region of SALL4 that was critical in
310 SALL4 upregulation leading to cell growth. Further investigation will be needed to determine

311 whether this non-canonical methylation site, downstream from the promoter, is only
312 significantly affected in the context of pseudogene-mediated demethylation.

313 These studies represent one of the first examples of gene locus specific demethylation
314 resulting in the activation of an oncogene. By an innovative strategy, we have identified the
315 specific CpG island in the locus that is required for SALL4 activation and expression. This
316 approach could be extended to other loci to identify the CpG “regulating” modules allowing
317 gene expression and controlled by an RNA-mediated mechanism.

318 Our studies suggest a model in which hepatitis B viral infection upregulates SALL4P5,
319 followed by SALL4 (Fig. 4D). This novel insight addresses unmet clinical need in HCC as
320 HCC is one of the leading causes of cancer-related deaths globally and chronic hepatitis B
321 virus infection accounts for more than 50% of HCC cases (1,28). Elucidating molecular
322 mechanisms of oncogene activation during hepatitis B virus infection could enhance our
323 understanding of the pathogenesis of HCC, and hence, aid in development of robust
324 therapies. Increased SALL4P5 expression is observed in HBV-related HCC patients, making
325 this pseudogene and its related function in oncogene SALL4 activation relevant to developing
326 novel therapeutics in HBV-related HCC.

327

328 **Methods**

329 **Cell Culture**

330 All HCC cell lines were obtained from ATCC and grown according to the manufacturer’s
331 instructions in the absence of antibiotics. Human hepatocellular carcinoma cell lines
332 (SNU398, SNU387 and HepAD38B) were maintained in Dulbecco's Modified Eagle's
333 medium (DMEM) and Roswell Park Memorial Institute 1640 medium (RPMI) (Life
334 Technologies, Carlsbad, CA) with 10% fetal bovine serum (FBS) (Invitrogen) and 2 mM L-

335 Glutamine (Invitrogen). These cell lines were cultured at 37°C in a humidified incubator with
336 5% CO₂. SNU398 was derived in 1990 from a 42-year-old, Asian male hepatocellular
337 carcinoma patient. SNU387 was derived in 1990 from a 41-year-old, Asian female
338 hepatocellular carcinoma patient. HepAD38B cell line was derived from a 15-year-old, male
339 hepatoblastoma patient.

340

341 **RNA extraction and gene expression analysis**

342 Total RNA was extracted from cells using the Trizol® reagent (Invitrogen) and purified
343 using the RNeasy Mini kit from Qiagen. 1 µg of purified RNA was used for cDNA synthesis
344 using the High-Capacity cDNA Reverse Transcription Kit (ThermoFisher Scientific)
345 according to the manufacturer's instructions. The QuantStudio 5 Real-Time PCR System
346 (Thermo Fisher Scientific) was used to assess the expression levels of the mRNAs, miRNAs,
347 and pseudogenes of interest. GoTaq® qPCR Master Mix (Promega) was used as a SYBR
348 master mix reagent for the qPCR procedures. The qRT-PCR data was analyzed using the
349 QuantStudio™ Design & Analysis Software Version 1.2 (ThermoFisher Scientific) and
350 represented as relative expression($\Delta\Delta C_t$), normalized against either GAPDH or β -actin. The
351 primer sequences used for the quantitative real-time PCR are provided in Supplementary
352 Table 1.

353

354 **Genomic DNA extraction**

355 Genomic DNA was extracted from HCC cell cultures using DNeasy Blood & Tissue
356 kit(Qiagen) for bisulfite-sequencing assay according to the manufacturer's protocols.

357

358 **Bisulfite treatment and sequencing**

359 SALL4 5'UTR Exon 1 3'UTR region methylation status was assessed using bisulfite
360 sequencing. In brief, 1 µg of genomic DNA extracted using the DNeasy Blood & Tissue
361 kit(Qiagen) was bisulfite-converted by using the EZ DNA Methylation kit (Zymo Research).
362 PCR products were gel-purified (Qiagen) from the 1.5% TAE gel and cloned into the pGEM-
363 T Easy Vector System (Promega) for transformation. The cloned vectors were transformed
364 into ECOS 101 DH5α cells and miniprep was performed to extract plasmids. Sequencing
365 results were analysed using BiQ analyser software. Samples with more than 90% conversion
366 rate and 70% sequences identity were analysed. The minimum number of clones for each
367 sequenced condition was 6. Primers used for the bisulfite sequencing of SALL4
368 5' UTR - exon 1 - intron 1 are listed in Supplementary Table 2.

369

370 **Protein extraction and immunoblotting**

371 Total cell lysates in protein lysis buffer (PLB) (100mM KCl (Ambion), 5mM MgCl₂
372 (Ambion), 25mM EDTA pH 8.0 (Life Technologies), 10mM HEPES (Life Technologies),
373 0.5% NP-40 (Roche), 20mM DTT(Fermentas), Proteinase inhibitor tablet (Roche)). PLB was
374 added to the cell pellets and incubated for 15 minutes on ice. The lysates were centrifuged for
375 10 minutes at 15,000 x g at 4 °C. Protein concentrations were measured using the Bradford
376 Protein Assay (Bio-Rad Laboratories) and absorbance was measured at 595 nm on the Tecan
377 Infinite® 2000 PRO plate reader (Tecan, Seestrasse, Switzerland). Equal amounts of protein
378 for each sample were diluted with 4X sample buffer (ThermoFisher Scientific) and heated at
379 95 °C for 5 minutes. The proteins were resolved by SDS-PAGE 12% self-cast gel and
380 transferred onto polyvinylidene difluoride (PVDF) membrane using the Mini Trans-Blot®
381 Electrophoretic Transfer Cell (Bio-Rad) in transfer buffer (25 mM Tris, 192 mM Glycine,
382 and 20% (v/v) methanol (Fischer Chemical)]. After blocking with 5% milk in Tris-NaCl
383 buffer (TBS) (ThermoFisher Scientific) with 0.1% Tween-20 (Sinopharm Chemical Reagent
384 Co., Ltd) (TBST), membranes were incubated with the appropriate primary and secondary

385 antibodies. The immune-reactive proteins were detected using the protein bands were
386 visualized using SuperSignal™ West Dura Extended Duration Substrate (ThermoFisher
387 Scientific) and visualized on Image Quant LAS 500 machine (GE Healthcare) according to
388 the manufacturer's instructions. SALL4 (Santa Cruz Biotechnology, EE30, #sc-101147),
389 GAPDH (Cell Signaling Technology, #5174), and DNMT1 (Abcam, #ab87656) antibodies
390 were used for immunoblotting as per manufacturer's instructions.

391

392 **Bacterial transformation**

393 ECOSTM 101 competent cells(DH5α) from Yeastern Biotech Co., Ltd. were used for
394 transformation following the manufacturer's instructions. 50μl of cells were thawed at room
395 temperature in a water bath. 2μl of pre-chilled DNA was immediately added. The tubes were
396 kept on ice for 5 minutes to increase the transformation efficiency. The cells went through
397 heat shock in a 42°C water bath for 30 seconds. The cells were kept on ice again for 5
398 minutes and plated on LB plates. The plates were incubated at 37°C for 16 hours.

399

400 **Plasmid extraction**

401 Plasmid was extracted from 1.5 ml of bacterial culture with the QIAamp DNA Mini
402 Kit(Qiagen) and purified for sequencing validation (1st BASE). For transfected cell line
403 DNA extraction, plasmid was extracted from the cell pellet that is suspended in 100 μl of
404 PBS.

405

406 **Plasmid transfection**

407 SNU398 cells were seeded at a density of 75,000 cells/well in 12-well plates 24 hours
408 before transfection. SNU387 and SNU182 cells were seeded at higher density of 100,000
409 cells/well in 12-well plates. 500 ng of plasmid was added to 3 μl of P3000 reagent (Life

410 Technologies) in 75 μ l of Opti-MEM prior to mixing with 2 μ l of Lipofectamine 3000 (Life
411 Technologies). The reagent mixture was incubated at room temperature for 10 minutes before
412 adding to each well.

413

414 **Cell viability assay**

415 24 hours post-transfection, cells were trypsinized and split into 5 individual wells of 5
416 separate 12-well plates. Upon adherence, cells were fixed using 10% neutral buffered
417 formalin solution (Sigma-Aldrich, HT501128-4L) and labelled as day 0. Subsequently, the
418 remaining plates were fixed daily from day 2 to day 5 (excluding day 1) prior to staining with
419 crystal violet (Sigma-Aldrich, C0775-100G) for 3-5 minutes at room temperature. Stained
420 wells were washed three times with Milli-Q water and left to dry. Crystal violet stain was
421 solubilized using 10% acetic acid (Sigma-Aldrich, A6283-2.5L). The plates were left on a
422 shaker at room temperature for at least 20 minutes. The absorbance reading for each well was
423 measured at 595 nm using the Tecan Infinite® 2000 PRO plate reader (Tecan, Seestrasse,
424 Switzerland).

425

426 **Soft agar assay**

427 A 0.6% agarose base was prepared by mixing 3.9 ml of 2% Ultrapure agarose (Invitrogen)
428 with 9.1 ml of cell culture medium. 2 ml of the mixture was added to individual wells of 6-
429 well plates. 24 hours post-transfection, cells were trypsinized, counted, and diluted to a
430 concentration of 15,000 cells/well. 450 μ l of the 2% agarose was added to 2.55 ml of cells for
431 a final agarose concentration of 0.3%, and 1 ml of the mixture was added to each well
432 containing the solidified 0.6% agarose base. 1 ml of culture media was added to the top agar
433 layer upon solidification and the cells were incubated at 37 °C in a humidified incubator.

434 Culture medium was changed every two days. Images of the colonies were taken under 4X
435 magnification every 4-5 days for a period of up to 14 days and quantified using ImageJ.

436

437 ***In vitro* transcription (IVT) and biotinylated RNA pulldown**

438 The DNA template was first amplified by PCR with primers containing a 5' T7-tag for in
439 vitro transcription. Antisense SALL4P5 control was also amplified by having the reverse
440 primers carrying the T7-tag. The IVT was performed as per manufacturer's guidelines. 1 µg
441 of purified PCR product was incubated with the transcription mix which was composed of
442 10X transcription buffer, 400 mM NTP mix, and 200 U T7 RNA polymerase for 5 hours at
443 37 °C. 140 µl of RNase-free water and 1000 µl of 100% ethanol were added to the
444 transcription product and incubated for at least 30 minutes at -20 °C. The reaction mixture
445 was centrifuged for 1 hour at 4 °C to precipitate the RNA. The RNA pellet was collected and
446 dissolved in 100 µl of ultra-pure water. The RNA was further purified using RNeasy Mini
447 250 columns (Qiagen) according to the manufacturer's instructions. The purified RNA
448 obtained from IVT was labelled with biotin at the 3' end using the Pierce™ RNA 3'End
449 Desthiobiotinylation Kit (ThermoFisher Scientific) according to the manufacturer's protocol.
450 Biotin labelling efficiency of the RNA probes was determined using the Chemiluminescent
451 Nucleic Acid Detection Module Kit (ThermoFisher Scientific) following the manufacturer's
452 protocol. Biotin labelling efficiency was normalized against the efficiency of antisense
453 transcript to determine amount of initial RNA bait used for the subsequent pulldown
454 experiment. Pulldown using these labelled RNA probes was carried out using the Pierce™
455 Magnetic RNA-Protein Pull-Down Kit (ThermoFisher Scientific) according to the
456 manufacturer's instructions. Protein lysates eluted from the pulldowns were used for
457 immunoblotting and other downstream analysis. The primer sequences with the T7-tag for
458 the PCR are provided in Supplementary Table 3.

459

460 ***In vitro* generation of sgRNA transcripts**

461 Approximately 1.4kb of genomic fragment spanning SALL4 5' UTR - exon 1 - intron 1
462 was PCR amplified (Zymo Research) and cloned into the pGEM-T Easy vector. The vector
463 was linearized with BamH1 restriction enzyme (New England Biolabs). SALL4-targeting
464 sgRNA candidates were transcribed with HiScribe™ Quick T7 High Yield RNA Synthesis
465 Kit (New England Biolabs) following manufacturer's instructions. The sgRNA target
466 sequences within SALL4 locus can be supplied upon request

467

468 ***In vitro* cleavage and selection of sgRNA transcripts**

469 *In vitro* cleavage assay was performed using purified Cas9 nuclease from *S. pyogenes*
470 (New England Biolabs) in order to select SALL4-specific sgRNA among candidates. The
471 experiment was performed according to the manufacture's protocol. The sgRNAs were
472 denatured at 95°C for 3 minutes, then Cas9 protein and sgRNAs were incubated for 10
473 minutes at 25°C to form a complex. Lastly, a linearized SALL4 DNA fragment was added to
474 the mixture and the entire reaction was incubated at 37°C for 1 hour. The reaction mixture
475 was composed of purified Cas9 protein, individual sgRNA, and linearized SALL4 genomic
476 fragment in ratio of 10: 10: 1. 1 ul of Proteinase K was added to each sample after the
477 cleavage reaction, and it was then incubated at room temperature for 10 minutes. The result
478 was analyzed with a 1% agarose gel.

479

480 **Lentiviral transduction of sgSALL4_1 and dCas9**

481 Lentiviruses expressing dCas9 or sgRNA were packaged in 293T cells the plasmids
482 psPAX2 and pMD2.G. TransIT-LT1 Transfection Reagent (Mirus) was used for transfection
483 into 293T cells. Virus was collected at 48 hours and 72 hours post-transfection. The collected

484 virus was filtered through 0.45 μm microfilters and stored at $-80\text{ }^{\circ}\text{C}$. Transduction of SNU-
485 387 cells was performed by mixing virus and 4 $\mu\text{g}/\text{mL}$ polybrene (Santa Cruz) together to
486 add to the cells seeded in T75 flasks 24 hours prior to the transduction. 24 hours after the
487 transduction, the medium was replenished with normal RPMI culture medium. Transduction
488 efficiency was determined by GFP (for sgRNA) or mCherry (dCas9) expression by FACS
489 analysis, and the positive cells were sorted by a FACS Aria machine (BD Biosciences).

490

491 **5-aza-2'-deoxycytidine(decitabine) treatment**

492 SNU387 cells were treated with 1.25 μM of 5-aza-2'-deoxycytidine (Sigma-Aldrich)
493 according to the manufacturer's instructions. Culture medium and drug were refreshed every
494 24 hours due to the drug being light-sensitive. RNA (for RT-PCR) and genomic DNA (for
495 bisulphite sequencing) were isolated after 5 days of treatment

496

497 **Digital droplet PCR**

498 Reactions for the ddPCR were prepared by harvesting 100,000 cells on each day for RNA
499 extraction and cDNA preparation. The reaction mixture was prepared with the 2x ddPCR
500 supermix for probes (Biorad, Cat #186-3026), 10-fold diluted cDNA, nuclease-free water,
501 and forward and reverse primers. Once the reaction mixture was ready, it was loaded into a
502 DG8 cartridge for the QX200 Automated Droplet Generator (Biorad, Cat#186-4003). We
503 then proceeded to the thermal cycling with a Biorad C1000 Touch Thermal Cycler with the
504 following cycle conditions: 95°C for 10minutes, 94°C for 30 seconds (40 cycles), 60°C for 2
505 minutes (40 cycles), 98°C for 10 minutes and 4°C hold. The reaction plate was loaded into
506 the QX200 Droplet Reader (Biorad, Cat#186-4003) for gene expression analysis.

507

508 **SALL4 methylation and expression correlation analysis**

509 Pearson correlations between SALL4 expression and methylation levels were performed
510 for the sites within the 5' UTR - exon 1 - intron 1 intron 1 region and distant sites in the
511 intron 1 (Fig S2). In order to show that a negative correlation is specific to primary HBV+
512 HCC patients, adjacent normal samples were used as negative controls. The data used for the
513 correlation was taken from Yang, et al (15), containing 19 pairs of primary HBV+ HCC
514 patients and their matched adjacent normal tissues.

515

516 **Materials Availability**

517 All plasmids and mouse lines generated in this study are freely available from the authors
518 upon reasonable request.

519

520 **Data and Code Availability**

521 The data used for the correlation was taken from Yang, et al, 2017 containing 19 pairs of
522 primary HBV+ HCC patients and their matched adjacent normal tissues. The authors declare
523 that all other data supporting the findings of this study are available within the paper and its
524 supplementary information files.

525

526 **References**

- 527 1. Yang JD, Roberts LR. Epidemiology and management of hepatocellular carcinoma. *Infect Dis*
528 *Clin North Am* **2010**;24(4):899-919, viii doi S0891-5520(10)00059-00.1016/j.idc.2010.07.004.
- 529 2. Keating GM, Santoro A. Sorafenib: a review of its use in advanced hepatocellular carcinoma.
530 *Drugs* **2009**;69(2):223-40 doi 10.2165/00003495-200969020-000066.
- 531 3. Yang J, Chai L, Fowles TC, Alipio Z, Xu D, Fink LM, *et al.* Genome-wide analysis reveals Sall4 to
532 be a major regulator of pluripotency in murine-embryonic stem cells. *Proc Natl Acad Sci U S*
533 *A* **2008**;105(50):19756-61 doi 10.1073/pnas.0809321105.
- 534 4. Zhang J, Tam WL, Tong GQ, Wu Q, Chan HY, Soh BS, *et al.* Sall4 modulates embryonic stem
535 cell pluripotency and early embryonic development by the transcriptional regulation of
536 Pou5f1. *Nat Cell Biol* **2006**;8(10):1114-23 doi 10.1038/ncb1481.
- 537 5. Ma Y, Cui W, Yang J, Qu J, Di C, Amin HM, *et al.* SALL4, a novel oncogene, is constitutively
538 expressed in human acute myeloid leukemia (AML) and induces AML in transgenic mice.
539 *Blood* **2006**;108(8):2726-35 doi 10.1182/blood-2006-02-001594.

- 540 6. Yue X, Xiao L, Yang Y, Liu W, Zhang K, Shi G, *et al.* High cytoplasmic expression of SALL4
541 predicts a malignant phenotype and poor prognosis of breast invasive ductal carcinoma.
542 *Neoplasma* **2015**;62(6):980-8 doi 10.4149/neo_2015_119.
- 543 7. Kobayashi D, Kuribayashi K, Tanaka M, Watanabe N. Overexpression of SALL4 in lung cancer
544 and its importance in cell proliferation. *Oncol Rep* **2011**;26(4):965-70 doi
545 10.3892/or.2011.1374.
- 546 8. Oikawa T, Kamiya A, Zeniya M, Chikada H, Hyuck AD, Yamazaki Y, *et al.* Sal-like protein 4
547 (SALL4), a stem cell biomarker in liver cancers. *Hepatology* **2013**;57(4):1469-83 doi
548 10.1002/hep.26159.
- 549 9. Yang J, Chai L, Gao C, Fowles TC, Alipio Z, Dang H, *et al.* SALL4 is a key regulator of survival
550 and apoptosis in human leukemic cells. *Blood* **2008**;112(3):805-13 doi 10.1182/blood-2007-
551 11-126326.
- 552 10. Yong KJ, Gao C, Lim JS, Yan B, Yang H, Dimitrov T, *et al.* Oncofetal gene SALL4 in aggressive
553 hepatocellular carcinoma. *N Engl J Med* **2013**;368(24):2266-76 doi
554 10.1056/NEJMoa1300297.
- 555 11. Li A, Jiao Y, Yong KJ, Wang F, Gao C, Yan B, *et al.* SALL4 is a new target in endometrial cancer.
556 *Oncogene* **2015**;34(1):63-72 doi 10.1038/onc.2013.529.
- 557 12. Liu BH, Jobichen C, Chia CSB, Chan THM, Tang JP, Chung TXY, *et al.* Targeting cancer
558 addiction for SALL4 by shifting its transcriptome with a pharmacologic peptide. *Proc Natl*
559 *Acad Sci U S A* **2018**;115(30):E7119-E28 doi 10.1073/pnas.1801253115.
- 560 13. Fan H, Cui Z, Zhang H, Mani SK, Diab A, Lefrancois L, *et al.* DNA demethylation induces SALL4
561 gene re-expression in subgroups of hepatocellular carcinoma associated with Hepatitis B or
562 C virus infection. *Oncogene* **2017**;36(17):2435-45 doi 10.1038/onc.2016.399.
- 563 14. Amabile G, Di Ruscio A, Muller F, Welner RS, Yang H, Ebralidze AK, *et al.* Dissecting the role
564 of aberrant DNA methylation in human leukaemia. *Nat Commun* **2015**;6:7091 doi
565 10.1038/ncomms8091.
- 566 15. Di Ruscio A, Ebralidze AK, Benoukraf T, Amabile G, Goff LA, Terragni J, *et al.* DNMT1-
567 interacting RNAs block gene-specific DNA methylation. *Nature* **2013**;503(7476):371-6 doi
568 10.1038/nature12598.
- 569 16. Chalei V, Sansom SN, Kong L, Lee S, Montiel JF, Vance KW, *et al.* The long non-coding RNA
570 Dali is an epigenetic regulator of neural differentiation. *Elife* **2014**;3:e04530 doi
571 10.7554/eLife.04530.
- 572 17. Wang L, Zhao Y, Bao X, Zhu X, Kwok YK, Sun K, *et al.* LncRNA Dum interacts with Dnmts to
573 regulate Dppa2 expression during myogenic differentiation and muscle regeneration. *Cell*
574 *Res* **2015**;25(3):335-50 doi 10.1038/cr.2015.21.
- 575 18. Merry CR, Forrest ME, Sabers JN, Beard L, Gao XH, Hatzoglou M, *et al.* DNMT1-associated
576 long non-coding RNAs regulate global gene expression and DNA methylation in colon cancer.
577 *Hum Mol Genet* **2015**;24(21):6240-53 doi 10.1093/hmg/ddv343.
- 578 19. Hawkins PG, Morris KV. Transcriptional regulation of Oct4 by a long non-coding RNA
579 antisense to Oct4-pseudogene 5. *Transcription* **2010**;1(3):165-75 doi
580 10.4161/trns.1.3.13332.
- 581 20. Watanabe T, Totoki Y, Toyoda A, Kaneda M, Kuramochi-Miyagawa S, Obata Y, *et al.*
582 Endogenous siRNAs from naturally formed dsRNAs regulate transcripts in mouse oocytes.
583 *Nature* **2008**;453(7194):539-43 doi 10.1038/nature06908.
- 584 21. Poliseno L, Salmena L, Zhang J, Carver B, Haveman WJ, Pandolfi PP. A coding-independent
585 function of gene and pseudogene mRNAs regulates tumour biology. *Nature*
586 **2010**;465(7301):1033-8 doi 10.1038/nature09144.
- 587 22. Han L, Yuan Y, Zheng S, Yang Y, Li J, Edgerton ME, *et al.* The Pan-Cancer analysis of
588 pseudogene expression reveals biologically and clinically relevant tumour subtypes. *Nat*
589 *Commun* **2014**;5:3963 doi 10.1038/ncomms4963.

- 590 23. Yang Y, Chen L, Gu J, Zhang H, Yuan J, Lian Q, *et al.* Recurrently deregulated lncRNAs in
591 hepatocellular carcinoma. *Nat Commun* **2017**;8:14421 doi 10.1038/ncomms14421.
592 24. Di Bisceglie AM. Hepatitis B and hepatocellular carcinoma. *Hepatology* **2009**;49(5
593 Suppl):S56-60 doi 10.1002/hep.22962.
594 25. Torrents D, Suyama M, Zdobnov E, Bork P. A genome-wide survey of human pseudogenes.
595 *Genome Res* **2003**;13(12):2559-67 doi 10.1101/gr.1455503.
596 26. Esteller M. Cancer epigenomics: DNA methylomes and histone-modification maps. *Nat Rev*
597 *Genet* **2007**;8(4):286-98 doi 10.1038/nrg2005.
598 27. Huhtanen CN, Feinberg JI, Trenchard H, Phillips JG. Acid Enhancement of Clostridium
599 botulinum Inhibition in Ham and Bacon Prepared with Potassium Sorbate and Sorbic Acid. *J*
600 *Food Prot* **1983**;46(9):807-10 doi 10.4315/0362-028X-46.9.807.
601 28. Ganem D, Prince AM. Hepatitis B virus infection--natural history and clinical consequences.
602 *N Engl J Med* **2004**;350(11):1118-29 doi 10.1056/NEJMra031087.

603

604 **Author contributions**

605 D.G.T., L.C., and Y.T. initiated the project and provided guidance throughout. D.G.T.,
606 L.C., Y.T., C.G., and J.K. designed the experiments. J.K. carried out experiments, analyzed
607 data, prepared figures, and wrote the manuscript. C.G, Y.L, and A.J carried out experiments,
608 prepared figures and edited the manuscript. M.A.B. analyzed the data, prepared figures, and
609 edited the manuscript. H.Y. and L.Y. performed the bioinformatics analysis on SALL4
610 expression and methylation. A.D.R. and J.Y. reviewed the manuscript. D.G.T, L.C, and Y.T
611 conceived of and supervised the project, designed experiments, and critically reviewed the
612 manuscript.

613

614 **Acknowledgements**

615 We thank Sudhakar Jha and Polly Chen for their insightful suggestions. We also thank the
616 Tenen, Tay, and Chai lab members for reviewing the manuscript. D.G.T. is funded by the
617 Singapore Ministry of Health's National Medical Research Council under its Singapore
618 Translational Research (STaR) Investigator Award, as well as NIH grants 1R35CA197697
619 and P01HL131477. Y.T. is funded by a Singapore National Research Foundation Fellowship
620 and National University of Singapore President's Assistant Professorship. This research is
621 supported by the National Research Foundation Singapore and the Singapore Ministry of

622 Education under its Research Centres of Excellence initiative, as well as the RNA Biology
623 Center at the Cancer Science Institute of Singapore, NUS, as part of funding under the
624 Singapore Ministry of Education's AcRF Tier 3 grants, Grant number MOE2014-T3-1-006.
625 L.C. is funded by NIH/NHLBI grant P01HL095489. A.D.R. is funded by NCI
626 R00CA188595, the Italian Association for Cancer Research (AIRC) start up grant #2014-
627 15347, and the Giovanni Armenise-Harvard Foundation.

628

629 **Figure Legends**

630 **Figure 1: SALL4 expression is negatively correlated with methylation of the**
631 **5'UTR-exon 1-intron 1 region. (A)** Schematic representation of the methylation probes. The
632 numbers refer to each CpG dinucleotide. The probe in the 5'UTR-exon 1 junction, "Probe 1",
633 assesses the methylation status of the CpG dinucleotide #11 and the intronic probe, "Probe 2"
634 assesses the CpG dinucleotide #68 in Supplementary Figure 1. **(B and C)** SALL4 expression
635 and methylation correlation analysis in 19 HBV+ patients. Compared to paired adjacent non-
636 transformed liver tissue, there is a negative correlation between SALL4 expression and
637 Probe 1 methylation, which is not observed using Probe 2. **(D)** Bisulfite sequence of the
638 5'UTR-exon 1 intron 1 region in wildtype SNU398 and SNU387 HCC cell lines. White color
639 represents hypomethylation while black represents hypermethylation of the individual CpG
640 dinucleotide. Degree of methylation was determined as a proportion of methylated cytosine
641 residue at a position out of 10 clones. Only CpG dinucleotides 1 to 35 are represented as the
642 sequencing efficiency was poor for dinucleotides 36 to 39. **(E)** Absolute quantification of
643 SALL4 mRNA expression in wildtype SNU398 and SNU387. β -actin was used as a positive
644 control for the assay. cDNA for β -actin quantification was diluted 10 times and back-
645 calculated accordingly later. The levels of β -actin were comparable between SNU398 and
646 SNU387 at about 400 to 600 copies of transcripts per cell. However, SNU398 cells expressed

647 more than 150 copies of SALL4 mRNAs while SNU387 expressed less than 10 copies on
648 average. (F) SALL4 protein levels in wildtype SNU398 and SNU387. β -actin was used as a
649 loading control for immunoblotting.

650

651 **Figure 2. CRISPR-DiR demethylates and activates SALL4.** (A) single-guide RNA design
652 for the CRISPR-DiR. The red region targets and interacts with the SALL4 5'UTR. The black
653 region interacts with dCas9. The blue regions are the two segments of ecCEBP α that interact
654 with DNMT1 (15). (B) Bisulfite sequence of CRISPR-DiR transduced SNU387 cells. The
655 data represents the methylation profile 14 days after CRISPR-DiR for SALL4. The numbers
656 indicate each of the CpG dinucleotides in the 5'UTR-exon 1 intron 1 junction. The white
657 color represents hypomethylation while the black hypermethylation of the individual CpG
658 dinucleotides. Only CpG dinucleotides 1 to 35 are represented as the sequencing efficiency
659 was poor for dinucleotides 36 to 39. (C & D) SALL4 transcript and protein levels after
660 CRISPR-DiR in SNU387. The western blot image is cropped as there were multiple lanes in
661 between "D21" and "5-aza". However, they are from the same blot and exposed for the same
662 duration. (E) Soft agar growth assay for CRISPR-DiR in SNU387. (F) Growth curve assay
663 for CRISPR-DiR for SALL4 in SNU387. 5-aza-2-deoxycytidine(decitabine) was used as a
664 positive control. NT denotes non-targeting negative control. Mean \pm SD, $n \geq 3$, *P < 0.05;
665 **P < 0.01; ***P < 0.001.

666

667 **Figure 3. SALL4P5 demethylates and activates SALL4, and associates with DNMT1.**

668 (A) Bisulfite sequencing after transient overexpression of individual SALL4 pseudogenes in
669 SNU387. The methylation status of CpG dinucleotides in SALL4 5'UTR-exon 1 intron is
670 shown. (B) Bisulfite sequencing after SALL4P5 knockdown in SNU398. ShScr denotes
671 scrambled shRNA. shSALL4P7 was used as a negative control. (C) Transcript localization in

672 SNU398. Cells were fractionated into nuclear and cytoplasmic fractions and transcript
673 expression was quantified. β -tubulin was used as a cytoplasmic fraction control, 18S rRNA
674 as the nuclear fraction control. **(D)** Biotin-labelled pull-down of DNMT1 in SNU398. Full
675 length SALL4P5 was used as a bait to pull down complexes and DNMT1 presence was
676 probed using immunoblotting. “as P5” denotes the negative control, antisense-SALL4P5 and
677 “ecCEBP α ” denotes the positive control. Full length SALL4P7 was used as a pseudogene
678 negative control as well. **(E and F)** SALL4 transcript and protein expression after pseudogene
679 overexpression in SNU387. **(G)** Soft agar growth assay for pseudogene overexpression in
680 SNU387. Mean \pm SD, $n \geq 3$, * $P < 0.05$; ** $P < 0.01$; *** $P < 0.001$.

681

682 **Figure 4. SALL4P5 is upregulated in HCC patients and during hepatitis B induction.**

683 **(A)** Relative SALL4 transcript expression in paired HCC patient samples. All expression
684 data are normalized against adjacent non-transformed tissues. Patient #1 was used as a
685 negative control with unaltered SALL4 and SALL4P5 levels. The levels of SALL4,
686 SALL4P5, and SALL4P7 were assessed for the other six patients, patients #2 to 7, as they
687 had more than 1.5-fold elevation in SALL4 levels compared to adjacent non-transformed
688 tissues. **(B)** Absolute quantification (RNA copies per cell) of Hepatitis B antigen X, HBx,
689 and SALL4 transcripts during hepatitis B induction in HepAD38B. Transcript levels of HBx
690 and SALL4 transcripts were monitored every 6-12 hours post HBV induction **(C)** Average
691 methylation in SALL4 5'UTR-exon 1 intron 1 across the 35 CpG dinucleotides in
692 HepAD38B. Blue denotes SALL4 methylation profile without HBV induction, while red
693 denotes SALL4 methylation profile after induction. Mean \pm SD, $n \geq 3$. * $P < 0.05$; ** $P < 0.01$;
694 *** $P < 0.001$.

Figure 1

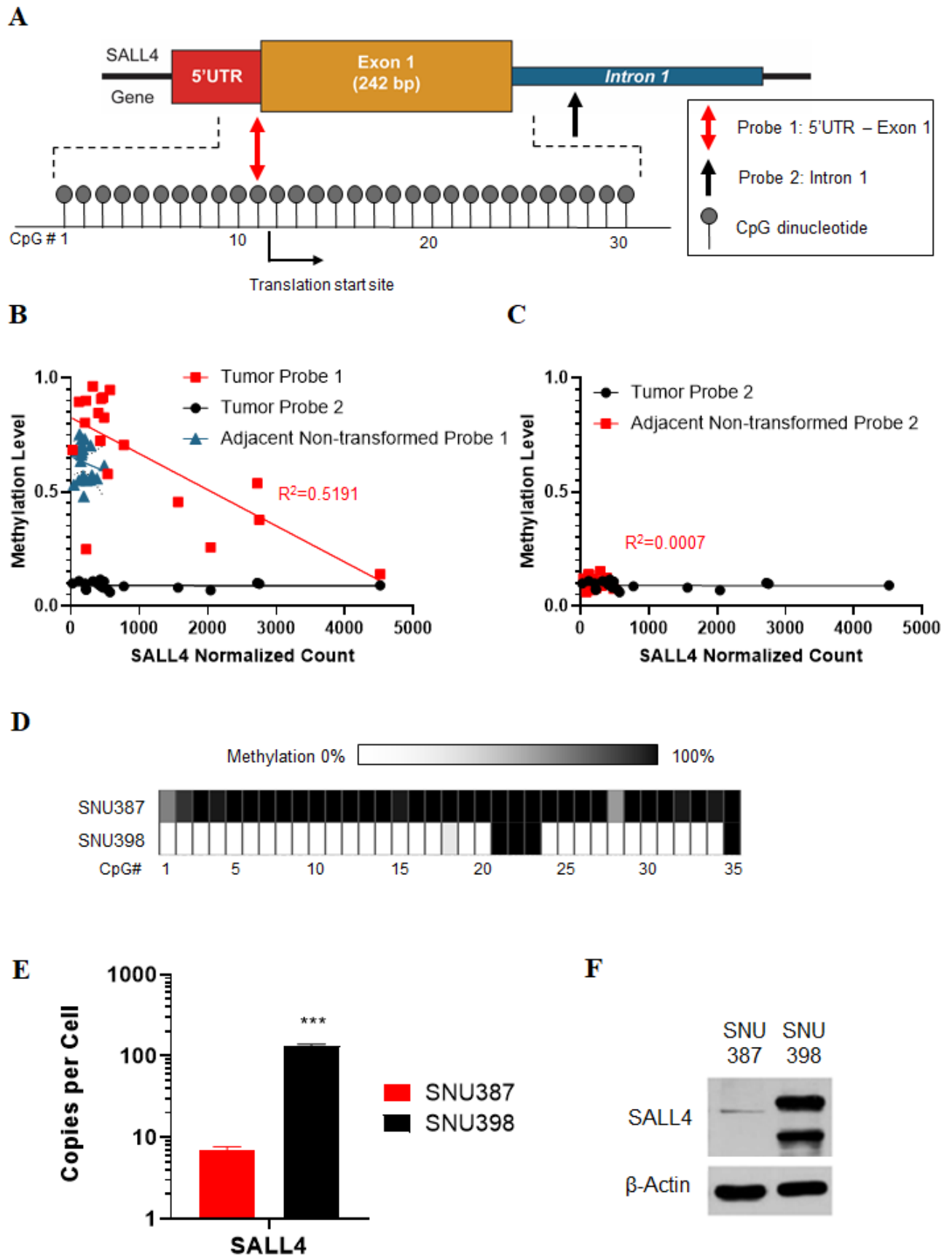
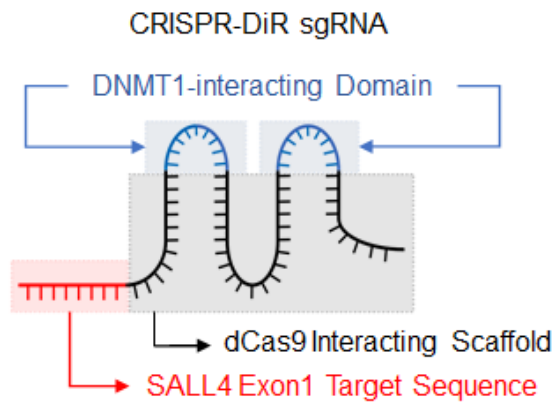
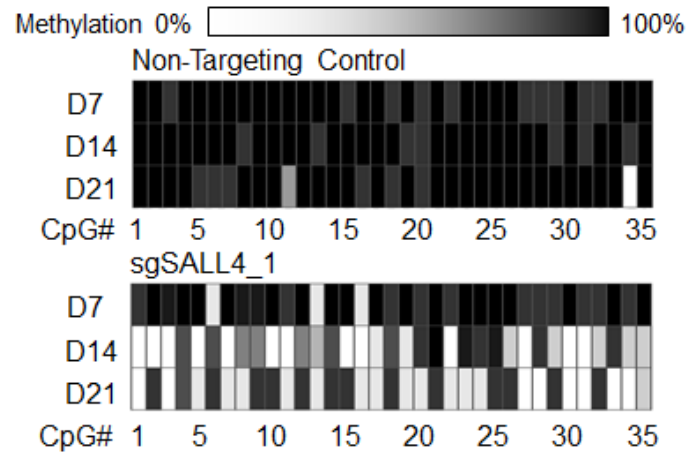


Figure 2

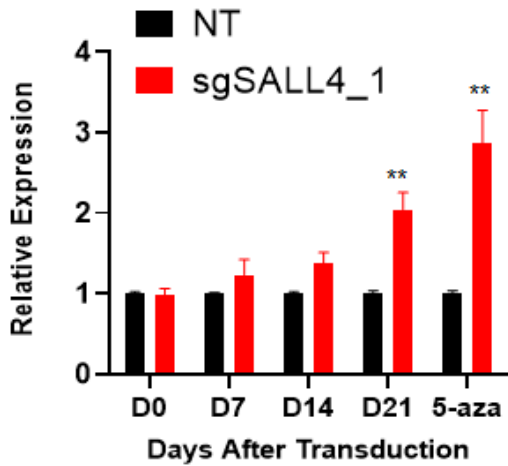
A



B



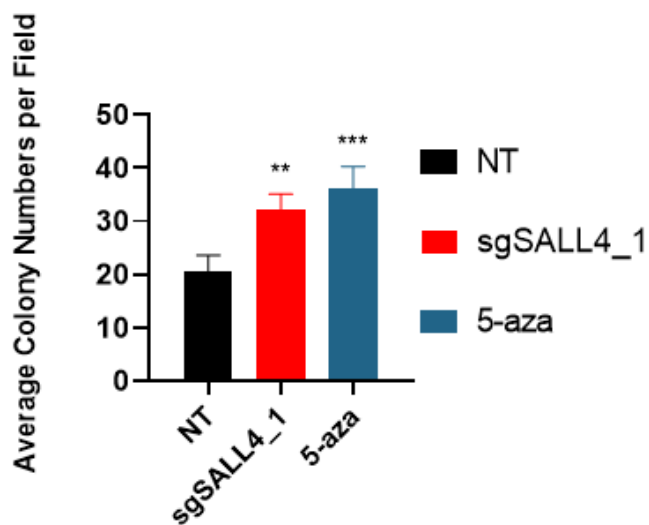
C



D



E



F

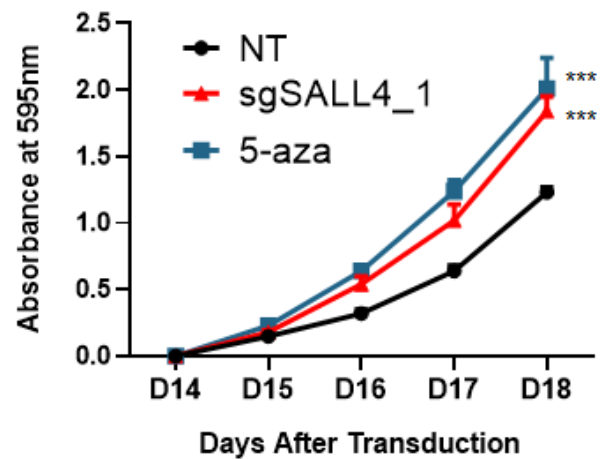


Figure 3

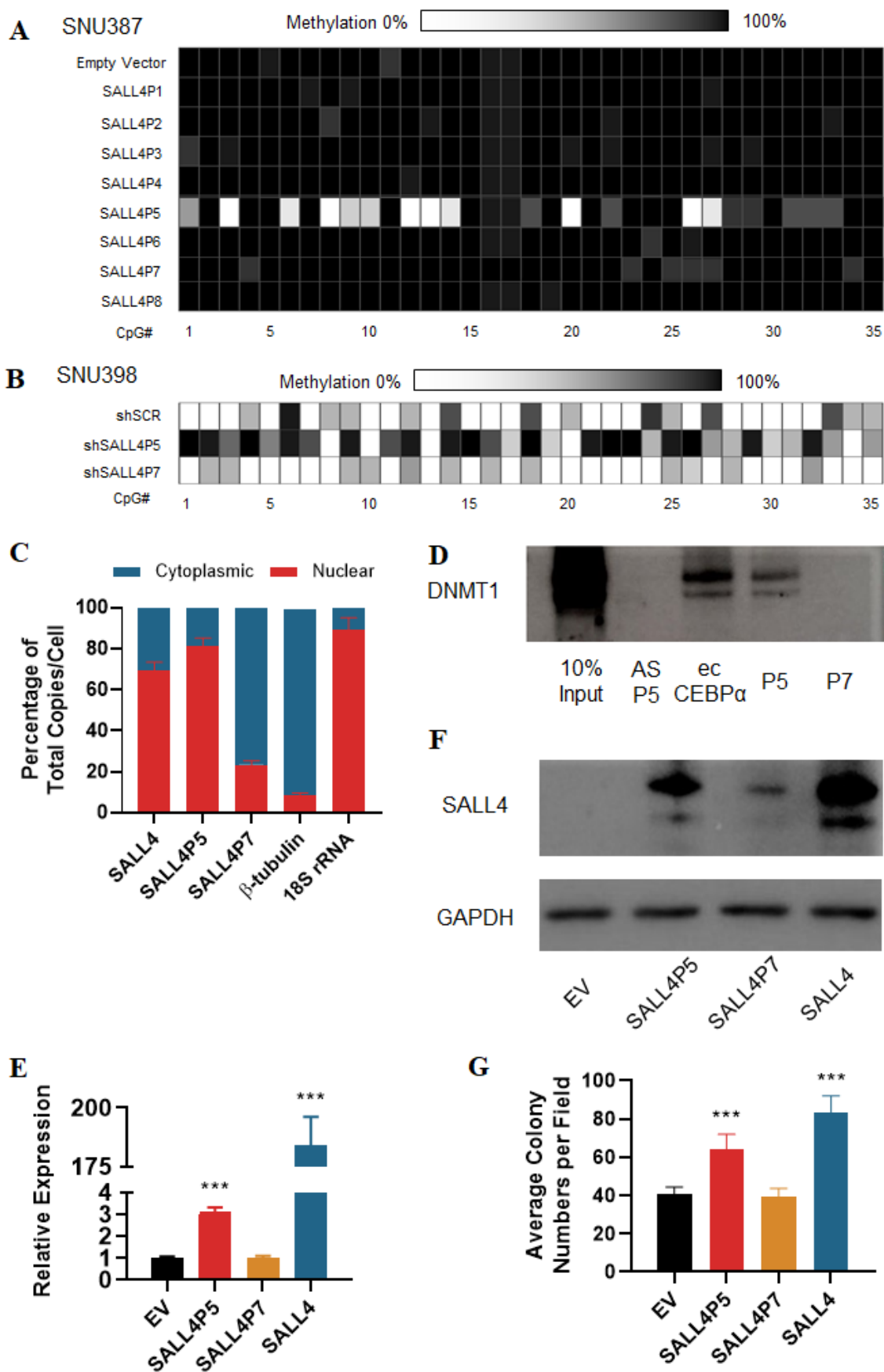


Figure 4

