

1 **Extrachromosomal DNA is associated with chromothripsis events and diverse**
2 **prognoses in gastric cardia adenocarcinoma**

3
4
5 **Authors:** Xue-Ke Zhao^{1,3}, Pengwei Xing^{2,3}, Xin Song^{1,3}, Miao Zhao², Linxuan Zhao²,
6 Yonglong Dang², Ling-Ling Lei¹, Rui-Hua Xu¹, Wen-Li Han¹, Pan-Pan Wang¹, Miao-
7 Miao Yang¹, Jing-Feng Hu¹, Kan Zhong¹, Fu-You Zhou¹, Xue-Na Han¹, Chao-Long
8 Meng¹, Jia-Jia Ji¹, Xingqi Chen^{2,4*}, Li-Dong Wang^{1,4*}

9
10
11 **Affiliations:**

12 ¹ State Key Laboratory of Esophageal Cancer Prevention & Treatment and Henan Key
13 Laboratory for Esophageal Cancer Research of The First Affiliated Hospital, Zhengzhou
14 University, Zhengzhou, Henan Province, 450052, PR China

15
16 ² Department of Immunology, Genetics and Pathology, Uppsala University, 75108,
17 Uppsala, Sweden

18
19
20 ³ Contributed equally

21
22 ⁴ Lead author

23 * Correspondence to: L.D.W, ldwang2007@126.com or X.C, xingqi.chen@igp.uu.se
24

25 **Abstract:**

26 Extrachromosomal DNA plays an important role in oncogene amplification in tumour cells and
27 poor outcomes across multiple cancers. However, the function of extrachromosomal DNA in
28 gastric cardia adenocarcinoma (GCA) is very limited. Here, we investigated the availability and
29 function of extrachromosomal DNA in GCA from a Chinese cohort of GCA using whole-
30 genome sequencing (WGS), whole-exome sequencing (WES), and immunohistochemistry.
31 For the first time, we identified the ecDNA amplicons present in most GCA patients, and found
32 that some oncogenes are present as ecDNA amplicons in these patients. We found that
33 oncogene ecDNA amplicons in the GCA cohort were associated with the chromothripsis
34 process and may be induced by accumulated DNA damage due to local dietary habits in the
35 geographic region. Strikingly, we observed diverse correlations between the presence of
36 ecDNA oncogene amplicons and prognosis, where *ERBB2* ecDNA amplicons correlated with
37 good prognosis, *EGFR* ecDNA amplicons correlated with poor prognosis, and *CCNE1* ecDNA

38 amplicons did not correlate with prognosis. Large-scale ERBB2 immunohistochemistry results
39 from 1668 GCA patients revealed that there was a positive correlation between the presence
40 of ERBB2 and prognosis in 2-7-year survival; however, there was a negative correlation
41 between the presence of ERBB2 and prognosis in 0-2-year survival. Our observations indicate
42 that the presence of *ERBB2* ecDNA in GCA patients may represent a good prognosis marker.

43

44

45 **Introduction**

46 Extrachromosomal DNA (ecDNA) was first identified more than half a century ago¹, and has
47 been associated with genomic instability^{2,3}. With next-generation sequencing technologies and
48 high throughput imaging platforms, an increasing number of studies have shown that ecDNAs
49 are present in most tissues, and contribute to the intratumoral heterogeneity and cancer
50 progression^{2,4-8}. Using computational analysis of whole-genome sequencing (WGS) data from
51 a large-scale cancer cohort, it has been demonstrated that the presence of ecDNA is cancer-
52 type specific, and is associated with oncogene amplification and poor outcomes across
53 multiple cancers⁷.

54

55 The cardia is located between the esophagus and the stomach. Gastric cardia
56 adenocarcinoma (GCA) and esophageal squamous cell carcinoma (ESCC) occur together in
57 the Taihang Mountains of north central China at high rates⁹⁻¹¹. Gastric cancer in this area
58 occurs primarily in the uppermost portion of the stomach and is referred to as GCA, and those
59 in the remainder of the stomach are called gastric noncardia adenocarcinoma (GNCA)¹².
60 Adenocarcinomas from junction of esophagogastric junction are usually classified as Siewert
61 type II of esophagogastric junction adenocarcinoma in western countries¹³⁻¹⁷, where Barrett's
62 esophagus is very common and has been considered as an important precancerous lesion of
63 adenocarcinoma at esophagogastric junction¹⁸. However, GCA from a Chinese population in
64 this area has distinct features compared to Western countries^{11,18,19}, and very low frequency
65 of Barrett's esophagus is observed¹⁸. Instead, GCA in this area shares similar features with
66 that of esophageal squamous cell carcinoma^{11,18}. A previous study reported that oncogene
67 amplification and gene rearrangements drive the progression and poor prognosis of GCA²⁰.
68 However, it is still unclear whether ecDNA is present in GCA, and what role it plays in the GCA
69 progression or whether it is correlated with patient prognosis. Therefore, we investigated the
70 availability and function of ecDNA in GCA in a Chinese cohort of GCA using whole-genome
71 sequencing (WGS), whole-exome sequencing (WES), and immunohistochemistry, and
72 explored the relationship between the presence of oncogene ecDNA amplicons and prognosis
73 in GCA.

74

75 **Results**

76 **Characterization of ecDNA amplicons in GCA**

77 Since ecDNA can be identified from WGS data using amplification region reconstruction tool,
78 AmpliconArchitect (AA)^{2,4-7,21}, we first performed WGS of 36 pairs of GCA tumour and tumour-
79 adjacent normal tissue from a high incidence GCA rate region in the northern region of China,
80 Henan Province (see **Methods**). All of our WGS data in 36 pairs of samples had sufficient
81 sequencing coverage and a high mapping rate (>95% mapping rate) (**Supplementary Fig. 1a**,
82 **Supplementary Table 1**). In addition, we performed single-nucleotide variant (SNV) analysis
83 in the 36 GCA patients and found that the top ranking mutated gene (81% mutation rate) was
84 TP53 (**Supplementary Fig. 1b**), which agrees with previous gene mutation studies in GCA
85 patients^{12,18,20,22}. Then, we applied AA to these 36 pairs of whole genome sequencing (WGS)
86 data pertaining to GCA tumor and tumor-adjacent normal tissue (**Fig. 1a**). Following the AA
87 pipeline, we treated the tumour-adjacent normal tissue as the background to call the somatic
88 copy number alteration(CNA) and identified ecDNA amplicons in our GCA cohort. Using this
89 strategy, ecDNA amplicons were identified in 28 of 36 GCA patients(**Fig. 1b**), and the
90 frequency (77.8%) of ecDNA amplicons observed in our GCA cohort is similar to that of
91 esophageal cancer (~80%) but higher than that of gastric cancer (~50%) in a previous report⁷.
92 Moreover, the number of ecDNAs identified from individual patients showed the high
93 heterogeneity across the GCA cohort (**Fig. 1b**), with a range of ecDNA amplicons from 0 to 24.
94 For most patients, the number of ecDNA amplicons was less than 10, and only five patients
95 had more than 10 ecDNA amplicons (**Fig. 1b**). In our GCA cohort, ecDNA amplicons were
96 further classified into five categories⁷ (**Fig. 1b, Supplementary Fig. 1c-e, Supplementary**
97 **Table 2**): Circular (n = 45), Complex (n = 21), Linear (n = 50), breakage-fusion-bridge (BFB)
98 (n = 4) and Invalid (n = 31), which occurred heterogeneously across the GCA patient cohort
99 (**Fig. 1b**). We further validated the circular feature of circular ecDNA amplicons identified from
100 AA software using another in silico method, Circle-finder, which identifies circular DNA from
101 paired-end high-throughput sequencing data²³⁻²⁵. By checking the sequencing read orientation
102 and junction points of circular ecDNA using Circle-finder, we found that 89.94-100% of circular
103 ecDNA amplicons identified from AA contained the same junctional reads detected by Circle-
104 finder (**Supplementary Fig. 1f-h**). The high proportion of overlapping circular ecDNA
105 amplicons from Circle-finder and AA results convinced us that the ecDNA amplicons identified
106 with AA are reliable.

107
108 Next, we analyzed the size of ecDNA amplicons in our GCA cohort. The size of ecDNA
109 amplicons from GCA ranged from 100 Kbp to 22.6 Mbp, with a median size of 350 Kb
110 (**Supplementary Fig. 2a**), where 75% of ecDNA amplicons were between 1-2 Mbp, and only
111 1% of ecDNA amplicons were larger than 20 Mbp (**Supplementary Fig. 2b**). Some large

112 ecDNA amplicons (> 20 Mbp) could be deconvoluted into multiple potential combinations of
113 amplicons using AA software (**Supplementary Figure 2c**). Since deconvolution is performed
114 using a computational prediction, there is still the possibility that multiple structures from these
115 large ecDNA amplicons are independent from circular amplicons. We also investigated the
116 frequency of ecDNA amplicons in different chromosomes. We found ecDNA amplicons of
117 different lengths in all chromosomes (**Supplementary Fig. 2d, 2e**) and the number of ecDNA
118 amplicons in the different chromosomes was independent of the length of the chromosome
119 (**Supplementary Fig. 2d**). We concluded that ecDNA amplicons occur heterogeneously
120 across GCA patients (**Fig. 1b, Supplementary Fig. 2e**).

121
122 Next, we performed genomic annotation for all ecDNA amplicons (**Fig. 1c, Supplementary**
123 **Fig. 2f, 2h**). We found that ecDNA amplicons occurred in different parts of the genome,
124 including 2452 sites in protein coding regions and 579 sites in long intergenic non-protein
125 coding RNA (lincRNA) (**Fig. 1c**). However, the frequency of ecDNA amplicons observed in
126 coding regions (6.28%) was higher than the proportion of coding regions in the whole genome
127 (3.48%) (**Supplementary Fig. 2f**). Furthermore, the proportion of ecDNA amplicons detected
128 in the exons (14.5%) is higher than that of exons in the entire genome (9.2%) (**Supplementary**
129 **Fig. 2g**). These ecDNA amplicons are also identified at regions of small RNAs (**Fig. 1c**),
130 including miRNAs (302 sites), SnRNAs (130 sites), SnoRNAs (63 sites), and rRNAs (37 sites).
131 Interestingly, we found that 82 ecDNA amplicons were from oncogenes and tumor suppressor
132 genes (TSGs) (**Fig. 1c**). Next, we focused on the analysis of oncogene and TSG ecDNA
133 amplicons in our GCA cohort (**Fig. 1d**). The oncogene and TSG ecDNA amplicons across the
134 GCA cohort exhibited a high heterogeneity, and the number of such oncogene and TSG
135 ecDNA amplicons varied from 1 to 11 (**Fig. 1d, 1e**). Amplification of the cyclin-E1 (*CCNE1*)
136 in the GCA was observed in a previous report²⁶. Specifically, we found that *CCNE1* ecDNA
137 amplicons occurred in 11 patients in our cohort (**Fig. 1d**). *ERBB2* is a member of the human
138 epidermal growth factor receptor (*EGF* family), and it has been reported that *ERBB2*
139 amplification plays an important role in GCA progression²⁶. We found that four patients had
140 *ERBB2* ecDNA amplicons (**Fig. 1d**). The, *CDK12*, *EGFR* and *MYC*, oncogenes and TSGs
141 were also found in the ecDNA format in more than three patients in the cohort (**Fig. 1d**). The
142 other name for *ERBB2* is *HER2*, and *EGFR* is also called *HER1* or *ERBB1*²⁷. Both *HER1* and
143 *HER2* are members of the *EGF* family. The identification of *HER1* ecDNA and *HER2* ecDNA
144 in GCA reflects the role of the *EGF* family in GCA progression²⁸. However, we did not observe
145 codetection of *HER1* ecDNA amplicons and *HER2* ecDNA amplicons in the same GCA patient
146 (**Fig. 1d**), which likely indicates the heterogeneous features in our GCA cohort. The frequent
147 detection of ecDNA amplicons in The Cancer Genome Atlas (TCGA) reflects the presence of
148 cancer specific oncogene ecDNA amplicons in each cancer type⁷, where the ecDNA amplicons

149 from gastric cancer and esophagus cancer are investigated. Since the cardia is located at the
150 junction of esophageal and stomach, we next investigated whether the list of ecDNA amplicons
151 from GCA was similar to that of gastric cancer or esophageal cancer using the TCGA report.
152 We found that GCA shares some common oncogene ecDNA amplicons with both gastric
153 cancer and esophageal cancer including *CCNE1*, *EGFR*, and *MYC* (**Supplementary Fig. 3**).
154 The top two ranking ecDNA amplicons, *ERBB2* and *CCNE1*, were the same in both gastric
155 cancer and GCAs. However, the top ranking list of oncogene ecDNA amplicons was different
156 between esophageal cancer and GCAs (**Supplementary Fig. 3**), where *CCND1* and *EGFR*
157 were the top two ranking oncogene amplicons in the esophageal cancer. Our results indicate
158 that the top oncogene ecDNA amplicons from GCAs is more similar to those from gastric
159 cancer. In addition, we observed that several oncogenes and TSG ecDNA amplicons appear
160 in the same GCA patient (**Fig. 1d**). The cyclization of oncogene ecDNA is highly amplified due
161 to its rolling-circle replication mechanism, and the circular ecDNA could contain different
162 oncogenes from different regions of the genome². Thus, we examined whether these different
163 oncogenes and TSGs in the same patient were located in the same ecDNA amplicon. We first
164 divided the highly amplified regions into segments, recombined them together by read
165 orientation and read junctions, and further reconstructed circular ecDNA containing multiple
166 oncogenes and/or TSG ecDNA amplicons (**Fig. 1d-f**, **Supplementary Fig. 4a-d**,
167 **Supplementary Table 3**). We referred multiple (two or more than two) oncogenes and/or
168 TSGs in the same ecDNA amplicon as oncogene ecDNA co-amplification (**Fig. 1d**), and
169 investigated the frequency of such occurrences (**Fig. 1d, 1e**). We found i) co-amplification of
170 oncogenes occurred in 50% of patients (18 of 36 patients) (**Fig. 1e**, **Supplementary Fig. 4a**);
171 ii) the frequency of oncogene ecDNA co-amplification varied from 50% to 100% of all
172 oncogene amplifications in different patients (**Fig. 1e**); and iii) some pairs of oncogene ecDNA
173 co-amplifications were observed in more than one patient (**Supplementary Table 3**), where
174 oncogene and TSG ecDNA pairs of *ERBB2* and *CDK12*, *RARA* and *SMARCE1*, and *CBLC*
175 and *BCL3* occurred in 3 patients; oncogene ecDNA pairs of *EGFR* and *IRF4*, *PPARG* and
176 *RAF1*; and pairs of *CDK12*, *ERBB2* and *RARA* occurred in 2 patients. Interestingly, *EGFR* and
177 *CDK6* with a physical distance of 40 Mbp, are located in the same circular ecDNA (**Fig. 1f**).
178 Using the normal genome copy number as the background, we found that the *EGFR* and *CDK6*
179 circular ecDNAs were amplified forty times compared to other parts of the genome (**Fig. 1f**).
180 The coamplification of *EGFR* and *CDK6* in the same ecDNA amplicons indicates that different
181 genes could work together during the progression of GCAs.

182

183 **Validation of ecDNA amplicons using Circle-Seq**

184 To further evaluate the accuracy of ecDNA amplicon prediction from the AmpliconArchitect
185 prediction, we chose 10 pairs of GCAs from our cohort to perform ecDNA sequencing with

186 Circle-seq²⁹ (See **Methods, Supplementary Fig. 5a**). We performed ecDNA peak calling from
187 Circle-seq using adjacent normal tissue as the control³⁰. Among 10 pairs of these selected
188 GCA patients for Circle-Seq, seven of them were ecDNA amplicon positive by WGS prediction
189 (**Fig. 1b**), and ecDNA amplicons (ranging from 491 to 39020) were identified in all of them
190 using Circle-Seq (**Supplementary Fig. 5b**). Then, we checked the overlapping ecDNA
191 segments from Circle-seq and predicated ecDNA amplicons from the WGS in the seven pairs
192 of GCAs. We found that most ecDNA amplicons identified in the WGS appeared in the Circle-
193 seq peak, where 100% WGS ecDNA in four GCAs, more than 80% WGS ecDNA in two GCAs,
194 and 50% WGS ecDNA in one GCA were confirmed by Circle-seq (**Fig. 2a**). Since *CCNE1* was
195 the most dominant detected ecDNA amplicon across the cohort, we determined the detailed
196 structure of *CCNE1* in Circle-seq (**Supplementary Fig. 5c**). We found that there was a clear
197 enrichment of *CCNE1* in two GCAs from both Circle-seq and WGS, and that both had a similar
198 tendency for amplification (**Supplementary Fig. 5c**). However, there was no *CCNE*
199 amplification in the normal samples, in either WGS or Circle-seq, indicating that our ecDNA
200 amplicon detection, identified with AmpliconArchitect prediction from the WGS data, is reliable.
201 The AA computational tool not only predicted the ecDNA amplicon, but also provided the
202 structure of the ecDNA amplicon. Upon closer inspection comparing the fine structure of
203 ecDNA amplification between the WGS and Circle-seq, we found that the fine structure was
204 not always the same (**Fig. 2b**). The *FGFR2* ecDNA amplicon exhibited highly amplified
205 segments with fluctuations in WGS prediction but not in the Circle-seq detection (**Fig. 2b**). The
206 difference in the fine structure from WGS and Circle-seq likely reflects the technical bias of the
207 ecDNA amplicon prediction from the WGS and library preparation from the Circle-seq.

208

209 **EcDNA amplicons in GCA is associated with chromothripsis**

210 Even though ecDNA amplicons are widely detected in different types of cancer, the sources of
211 ecDNA amplification remain unknown. It has been reported that chromothripsis contributes to
212 cancer progression and drives ecDNA amplification in cancer^{3,31,32}, and that some ecDNA
213 amplicons are generated during chromothripsis process². Next, we aimed to understand the
214 relationship between chromothripsis and ecDNA amplicons in our GCA cohort. We used the
215 ShatterSeek package³³ to identify chromothripsis events across the 36 GCA patients
216 (**Supplementary Fig. 6a**). Strikingly, we found that chromothripsis occurred in 34 GCA
217 patients across our cohort (**Supplementary Fig. 6b**). We also divided the chromothripsis
218 events into fine categories with the parameters of high confidence (HC) and low confidence
219 (LC) (see **Methods**). This revealed that HC chromothripsis occurred in 61.1% of GCAs across
220 the cohort, and LC chromothripsis occurred in 88.9% of all GCA samples. We found that the
221 frequency of chromothripsis in GCA patients was quite diverse across the cohort, where the
222 range of chromothripsis was from 0 to 4 for HCs and 0 to 14 for LCs (**Supplementary Fig.**

223 **6c**). The location of the chromothripsis events in the genome was also quite heterogeneous
224 across the cohort (**Fig. 3a**). When we aligned chromothripsis events and ecDNA amplicons
225 on the genome browser, we observed a clear overlap between ecDNA amplicons and
226 chromothripsis at some of the oncogene ecDNA loci, including the *ERBB2* and *MYC* genes
227 (**Fig. 3b, Supplementary Fig. 7**). To further explore the relationship between chromothripsis
228 and ecDNA amplification, we quantified the number of ecDNA amplicons that overlapped with
229 chromothripsis (**Fig. 3c**). The results showed that 17.22% of ecDNA amplicons occurred in HC
230 chromothripsis, and 15.89% occurred in LC chromothripsis. Taken together, these results
231 indicate that 33.11% of ecDNA amplicons might be caused by chromothripsis (**Fig. 3c**). To
232 further determine the relationship between ecDNA amplicon and chromothripsis, we calculated
233 the correlation between the number of chromothripsis events and the total length of all ecDNA
234 (**Fig. 3d**). The results clearly demonstrated a positive correlation between ecDNA amplicons
235 and chromothripsis events (Pearson's correlation = 0.42). Our results indicate the ecDNA
236 amplicons in GCAs are more likely to occur due to chromothripsis, and that such events could
237 contribute to GCA progression if the chromothripsis event occurs at the oncogene site.

238

239 Comprehensive analysis of chromothripsis using large-scale samples of human cancers from
240 TCGA showed that the frequency of chromothripsis is greater than 50% in several cancer
241 types³⁴. However, the frequency of chromothripsis in our GCA cohort was 94% (**Fig. 3a**), which
242 is extremely high. Previous reports have shown that chromothripsis is associated with genomic
243 instability and DNA damage³⁵⁻³⁹. Thus, we investigated potential risk factors contributing to such
244 a high frequency of chromothripsis in our GCA cohort by analyzing genome stability and DNA
245 damage. First, we performed microsatellite instability (MSI) detection by
246 immunohistochemistry (IHC) staining of four proteins (MLH1, MSH2, MSH6 and PMS2)^{40,41}.
247 We found that only 9 of 36 samples were MSI-high samples (**Supplementary Fig. 8a, 8b,**
248 **Supplementary Table 4**), and 27 patients were MSI-low. The two chromothripsis-negative
249 samples were all in the MSI-low group (**Supplementary Fig. 8b**), and there was no correlation
250 between MSI grade and chromothripsis events (**Supplementary Fig. 8b**, $p = 1$, Fisher's exact
251 test). Thus, we concluded that the high frequency of chromothripsis is not likely due to the high
252 proportion of MSI-high samples in our cohort. Second, we calculated chromosomal instability
253 (CIN) for all 36 samples in accordance with a previous report⁴² and divided GCA patients into
254 four groups based on the genome integrity index (from low to high: 0 to 0.2, 0.2 to 0.4, 0.4 to
255 0.6, 0.6-0.8) (see **Methods**). We found only 2 samples in our GCA patients in the high-grade
256 CIN group (**Supplementary Fig. 8c, Supplementary Table 4**). The two chromothripsis-
257 negative samples were in the low-grade CIN group (**Supplementary Fig. 8c**), and there was
258 no correlation between CIN grade and chromothripsis events (**Supplementary Fig. 8c**, $p =$
259 0.381, Fisher's exact test). Thus, we concluded that the high frequency of chromothripsis is

260 not likely due to the high proportion of high-grade CIN in our cohort. Third, we performed IHC
261 staining of γ H2AX protein, a crucial biomarker for the detection of DNA double strand breaks⁴³,
262 in our GCA cohort. We found that 80.55% (29/36) of GCA patients were γ H2AX protein positive
263 (**Fig. 3e, 3f, Supplementary Table 4**). The two chromothripsis-negative samples were both
264 γ H2AX protein negative (**Fig. 3f**), and there was a significant correlation between the presence
265 of γ H2AX and chromothripsis events (**Fig. 3f**, $p = 0.033$, Fisher's exact test). We also found
266 that the total length of chromothripsis in γ H2AX protein-positive patients was significantly
267 longer than that in γ H2AX protein-negative patients (**Fig. 3g**, $p = 0.025$). Thus, we suspect that
268 the high frequency of chromothripsis is most likely due to the high degree of DNA damage that
269 has accumulated in GCA patients. All GCA patients in our study were from the high incidence
270 area for GCA in Henan Province, northern China⁹, where the intake of nitrosamine-rich foods,
271 such as pickled vegetables, has been well recognized as one of the key risk factors for GCA⁴⁴.
272 Accumulating evidence has demonstrated that nitrosamine is a very important factor for DNA
273 alkylation, synthesis disorder, high instability and even DNA double strand breaks⁴⁵⁻⁵⁰. Thus,
274 we suspected that nitrosamine exposure in our GCA cohort may accumulate DNA damage,
275 potentially inducing a high frequency of chromothripsis. As ecDNA amplicons in our GCA
276 cohort are more likely to occur due to chromothripsis, as stated above, and it was also
277 proposed that chromothripsis is a primary mechanism that accelerates genomic DNA
278 rearrangement and amplification into ecDNA by a recent study³, our data suggest that local
279 dietary habits from the geographic region in our cohort may contribute to ecDNA occurrence
280 in GCA patients.

281

282 **The presence of oncogene ecDNA does not increase the mutation frequency in GCA**

283 Oncogene amplification is a key factor contributing to human cancer⁵¹. A high frequency of
284 oncogene mutations has also been reported in GCA^{20,22}. Since both oncogene amplification
285 (**Fig. 1d**) and oncogene mutations (**Supplementary Fig. 1b**) were observed in our GCA cohort,
286 we investigated whether there was a high frequency of oncogene mutations in the region of
287 ecDNA oncogene amplicons. We calculated numbers of SNVs in the whole genome as well
288 as in only ecDNA amplicon present regions (**Supplementary Fig. 9a**) and found mutation
289 frequency in the ecDNA amplicon regions occur at a similar level as in the whole genome from
290 most patients, except for two GCA samples (**Supplementary Fig. 9a**). Statistical analysis
291 showed that there was no significant difference in mutation frequency between ecDNA
292 amplicon regions and the whole genome in our GCA cohort (**Supplementary Fig. 9b**, $p =$
293 0.18). We also compared the numbers of SNVs in regions of individual oncogene or TSG
294 ecDNA regions (same oncogene or TSG ecDNA observed in 2 or more patients) between
295 present and absent oncogene ecDNA patients (**Supplementary Fig. 9c**) and found that there
296 were significantly more SNVs in the ecDNA present group only with respect to the BIRC3 gene

297 **(Supplementary Fig. 9c, $p = 0.031$)** but not at other oncogenes **(Supplementary Fig. 9c)**.
298 Thus, we concluded that there may be no relationship between oncogene mutations and the
299 presence of oncogene ecDNA amplicons in GCA patients.

300

301 **The presence of oncogene ecDNA amplicons has the diverse correlation with the** 302 **prognosis of GCA**

303 It was reported that the presence of ecDNA is associated with oncogene amplification and
304 poor outcomes across multiple cancers⁷. Thus, we investigated the relationship between
305 oncogene amplification, the presence of ecDNA and patient prognosis in our GCA cohort. We
306 first explored the relationship between oncogene amplification and GCA patient prognosis by
307 focusing on the top 11 high frequency of oncogenes and TSGs ecDNA amplicons. We found
308 that most of the top 11 high frequency oncogene amplifications across the cohort with a copy
309 number (CN) greater than 5 came from ecDNA amplicons **(Supplementary Fig. 10)**. We
310 compared the gene copy numbers and patient survival time by splitting the gene amplification
311 into different groups (High, Low, Normal) **(Supplementary Fig. 10)**. As expected, the survival
312 time in some GCA patients after surgery was shorter in those with a high copy number of
313 certain oncogenes, including *EGFR*, *MYC*, and *BIRC3* **(Supplementary Fig. 10)**. Surprisingly,
314 we found that patients with a low CN amplification of *CCNE1* and *ERBB2* survived for a shorter
315 period compared to those with a normal gene CN **(Supplementary Fig. 10)**, and patients
316 survived even longer with a high CN of *CCNE1* and *ERBB2* amplification **(Supplementary Fig.**
317 **10)**. To further investigate our observation, we performed a correlation study between different
318 ranges of CN amplification and survival time from the *CCNE1*, *ERBB2*, and *EGFR* genes **(Fig.**
319 **4a)**. The results indicated that the short survival time was due to the high range of oncogene
320 amplification in *EGFR*. However, for *ERBB2* and part of the sample of *CCNE1*, the tendency
321 was completely opposite. Specifically, we found that four samples with a high CN of *CCNE1*,
322 caused by ecDNA amplicons, exhibited an average survival time of 5.08 years, and all samples
323 with a high CN of *ERBB2* had an average survival time of 6.59 years **(Fig. 4a)**.

324

325 Furthermore, we focused on investigating the relationship between prognosis and CN of three
326 oncogenes: *CCNE1*, *ERBB2*, and *EGFR*. *EGFR* followed the tendency that those with high-
327 range oncogene amplification had a decreased survival time than those with low-range
328 amplification ($p = 0.0013$) **(Fig. 4b)**. The relationship between *EGFR* copy number and patient
329 survival time reflects oncogene function in tumorigenesis from GCAs. For both *ERBB2* and
330 *CCNE1*, we found that patients with low range amplification had the worst prognosis compared
331 to those with normal and high range amplification**(Fig. 4b)**. To our surprise, patients with high
332 range amplification from *CCNE1* and *ERBB2* had the best prognosis compared to those with
333 low and middle range amplification **(Fig. 4b)**. To further confirm the relationship between

334 oncogene amplification and patient survival, we performed the WES sequencing on another
335 independent GCA cohort with 39 GCA patients together with our 36 GCA patient cohorts
336 (**Supplementary Fig. 11a, Supplementary Table 5**). First, the copy numbers of *ERBB2* from
337 WGS in the 36 patients were very similar to the copy numbers detected in the WES data
338 (**Supplementary Fig. 11b**), which indicates that the WES data could be used to validate our
339 WGS observation of *ERBB2* gene amplification. Next, we focused on the WES data for 75
340 GCA patients, and we observed a similar tendency, namely, that the high-range *ERBB2*
341 amplification was correlated with increased survival time (**Supplementary Fig. 11c,**
342 **Supplementary Table 6**). Taken together, we concluded that our observation is independent
343 of the specific GCA cohort. This negative correlation between oncogene amplification and
344 patient prognosis has previously been reported in many independent studies, including large
345 group studies in the TCGA⁷. We found a similar tendency for some oncogenes in GCA, such
346 as *EGFR*. The negative correlation is true for the low range amplification from *ERBB2* and
347 *CCNE1* (**Fig. 4b**); however, the correlation becomes positive when these two genes undergo
348 high range amplification (**Fig. 4b**).

349
350 Next, we investigated the relationship between the presence of oncogene ecDNA amplicons
351 and patient prognosis by dividing patients into ecDNA present and absent groups (**Fig. 4c**),
352 and we found diverse correlations of present oncogene ecDNA amplicons and patient
353 survival. In brief, we found no significant difference in prognosis for the absence and
354 presence of *CCNE1* ecDNA amplicons (**Fig. 4c**, $p = 0.55$); the presence of *EGFR* ecDNA
355 amplicons had a negative correlation with patient prognosis (**Fig. 4c**, $p = 0.036$); and the
356 presence of *ERBB2* ecDNA amplicons had a positive correlation with patient prognosis (**Fig.**
357 **4c**, $p = 0.0068$). To understand whether our observation was due to clinicopathological
358 factors from GCA patients, we first investigated the relationship between clinicopathological
359 phenotypes and prognosis in GCA (**Methods, Supplementary Fig. 12, Supplementary**
360 **Table 4**). We found that UICC tumour stage was the only clinicopathological factor correlated
361 with GCA survival (**Supplementary Fig. 12i**). Next, we performed survival analysis using
362 clinicopathological variables of patients together with the presence of ecDNA amplicons
363 (*ERBB2*, *EGFR*, *CCNE1*) by dividing patients into those with and without ecDNA amplicons
364 (**Supplementary Fig. 12**). We found that the presence of *ERBB2* ecDNA amplicons may be
365 relevant to the UICC tumour stage but not to other clinicopathological variables
366 (**Supplementary Fig. 12**). However, the presence of *EGFR* and *CCNE1* ecDNA amplicons
367 was not relevant to any clinicopathological variables (**Supplementary Fig. 12**). Since both
368 UICC tumour stage (**Supplementary Fig. 12i**) and the presence of *ERBB2* ecDNA (**Fig. 4c**)
369 are contributing factors to patient survival, we assumed that there might be some connection
370 between the presence of the *ERBB2* ecDNA amplicon and GCA stage. However, our sample

371 size was too small (36 cases) to obtain further conclusions. It will be very interesting to
372 perform further studies with larger sample sizes of patients to obtain additional conclusions in
373 the future.

374

375 The positive correlation between the presence of *ERBB2* ecDNA in GCA and patient
376 prognosis is paradoxical to large-scale TCGA studies in many cancer types⁷, where the
377 presence of ecDNA amplicons was shown to be associated with poor outcomes. Since it was
378 reported that there is a paradoxical relationship between chromosomal instability and survival
379 outcomes in cancer⁴², we examined whether the positive correlation between the presence of
380 *ERBB2* ecDNA amplicons and patient prognosis is due to chromosomal instability (CIN) in
381 our GCA cohort. The survival analysis from the four groups of CIN (**Methods**,
382 **Supplementary Fig. 8c**) shows that GCA patients with stable chromosomes survived longer
383 than patients with unstable chromosomes in our cohort (**Supplementary Fig. 13a**). However,
384 we did not find that *ERBB2* ecDNA amplicons present in samples were only enriched in
385 specific CIN groups (**Supplementary Fig. 13b**), and we did not observe a significant
386 difference in CIN values between ecDNA present samples and ecDNA absent samples
387 (**Supplementary Fig. 13c**, $p = 0.33$). Thus, we concluded that the paradoxical relationship
388 between the presence of *ERBB2* ecDNA amplicons in GCA patients and survival outcome is
389 independent of CIN. A recent study showed chromatin structure of ecDNA is highly
390 accessible⁵², we assumed that the *ERBB2* gene is highly expressed in ecDNA present GCA
391 patients. It was also reported the amplification of *ERBB2* gene was followed by *ERBB2* gene
392 overexpression in the same GCA tissue^{18,53-55}. At the same time, we observed a positive
393 correlation between *ERBB2* gene expression and ERBB2 protein expression in GCA patients
394 ($n = 44$) (**Supplementary Fig. 14a**, $R = 0.79$, **Supplementary Table 7**). Thus, we
395 hypothesized that protein levels of ERBB2 were also high in *ERBB2* ecDNA present
396 patients, and that a high level of ERBB2 protein would be positively associated with GCA
397 prognosis. To test our hypothesis, we performed immunohistochemistry of the ERBB2
398 protein from 1668 GCA patients (with 0- to 7- year survival time after surgery) (see **Methods**,
399 **Supplementary Fig. 14b**, **Supplementary Table 8**). Although we did not observe a
400 significant difference in patient prognosis among all patients ($n = 1668$, **Supplementary Fig.**
401 **14c**, $p = 0.16$), there was a significant difference in patient prognosis in patients surviving
402 between 0-2 years (including 2 years) after surgery ($n = 750$, **Fig. 4d**, $p = 0.016$) and in
403 patients surviving between 2-7 years ($n = 918$, **Fig. 4d**, $p = 0.025$). We concluded that there
404 is a positive correlation between ERBB2 protein presence and patient prognosis in 2-7 year
405 survival after surgery, and there is a negative correlation between ERBB2 protein presence
406 and patient prognosis in the 0-2 year survival after surgery. It was reported ERBB2 protein
407 expression and gene amplification correlate with better survival in esophageal

408 adenocarcinoma⁵⁶, and the positive correlation between the presence of ERBB2 protein and
409 increased patient survival (2-7 years of survival) in our GCA cohort likely also reflects the
410 similarity between esophageal adenocarcinoma features and GCA. Since we assumed that
411 the protein level of ERBB2 is high in *ERBB2* ecDNA-positive patients, our observation
412 indicates that the *ERBB2* ecDNA amplicon may represent a good prognostic marker in GCA
413 patients.

414

415 **Discussions**

416 In summary, for the first time, we identified ecDNA amplicons in GCA patients using WGS data,
417 and validated these ecDNA amplicons using Circle-seq. We found that these ecDNA
418 amplicons are present in most GCA patients, and have exhibit heterogeneity in different GCA
419 patients. Additionally, for the first time, we found that several oncogenes are in the format of
420 ecDNA amplicons in GCA patients and that different oncogenes could coamplify in the same
421 ecDNA amplicon. Interestingly, we found oncogene ecDNA amplicons were associated with
422 a high frequency of chromothripsis in our GCA cohort, and such a high frequency of
423 chromothripsis in our cohort is likely due to high degree of DNA damage induced by
424 nitrosamine exposure from a local diet⁴⁵⁻⁵⁰. We propose that local dietary habits from the
425 geographic region may have contributed to ecDNA occurrence in our GCA cohort. It was
426 reported that ecDNA is a major mechanism of drug resistance in several tumour types³, thus,
427 it will be valuable to follow clinical annotation on previous exposure therapy together with
428 ecDNA detection in large-scale samples of GCA patients to design therapy strategies for GCA
429 patients in the future.

430

431 Strikingly, we found that the correlation between the present oncogene ecDNA amplicons
432 and patient prognosis was different depending on gene in GCA patients, where *ERBB2*
433 ecDNA amplicons correlated with good prognosis, *EGFR* ecDNA amplicons correlated with
434 poor prognosis and *CCNE1* ecDNA amplicons did not correlate with prognosis. The
435 relationship between presence of ecDNA and prognosis in GCA reported in this study is
436 different from a previous report indicating that oncogene ecDNA amplicons correlate with
437 poor prognosis in other cancers from TCGA⁷, and our observation likely reflects the
438 heterogeneous nature of cancers. These diverse associations of oncogene ecDNA
439 amplification and prognosis may aid in designing better personal therapy strategies for GCA
440 patients in the future. Large-scale ERBB2 immunohistochemistry results from 1668 GCA
441 patients demonstrated that there was a positive correlation between ERBB2 protein
442 presence and patient prognosis in 2-7-year survival after surgery; however, there was a
443 negative correlation between ERBB2 protein presence and patient prognosis in 0-2-year
444 survival after surgery. This paradoxical relationship between ERBB2 protein presence and

445 prognosis is similar to a previous report on the relationship between ERBB2 protein
446 expression and improved survival in esophageal adenocarcinoma⁵⁶, which likely reflects the
447 similarity in features between esophageal adenocarcinoma and GCA. Since we assumed
448 that the protein levels of ERBB2 are high in *ERBB2* ecDNA-positive patients, our observation
449 indicates that the *ERBB2* ecDNA amplicon may represent a good prognostic marker in GCA
450 patients.

451

452 **Acknowledgements:** This work is supported by the National Key R&D Program of China
453 (2016YFC0901403 to L.D.W.), the National natural science foundation of China (81872032,
454 U1804262 to L.D.W.), the Swedish Research Council (VR-2016-06794, VR-2017-02074 to
455 X.C.), Beijer Foundation (to X.C.), Jeassons Foundation (to X.C.), Petrus och Augusta
456 Hedlunds Stiftelse (to X.C.), Göran Gustafsson's prize for younger researchers (to X.C.),
457 Vleugel Foundation (to X.C.), and Uppsala University (to X.C.).

458

459 **Author contributions**

460 L.D.W and X.C. conceived and designed the study; X.K.Z., X.S., L.L.L., R.H.X., W.L.H., P.P.W.,
461 and F.Y.Z. contributed to the collection of the patient materials and clinical information; X.K.Z.,
462 X.S., M.M.Y., J.F.H., and K.Z. prepared the WGS and ecDNA sequencing of GCA; P.X.
463 performed all the sequencing data mining; L. Z., Y. D., L.L.L., X.N.H., C.L.M., and J.J.J. were
464 responsible for the protein expressions of ERBB2, γ H2AX, and MSI staining in the GCA and
465 analysis of the relationship with the GCA survival; M.Z., X.K.Z., P.X., L.D.W., and X.C. wrote
466 the manuscript together with input from all authors; and L.D.W. and X.C. supervised all aspects
467 of this work.

468

469 **Competing Financial Interests statements**

470 The authors declare no competing financial interests.

471

472 **Availability of the data**

473 All raw data is deposited in the China National Center for Bioinformation with access number
474 of HRA00081

475

476 **References**

- 477 1 Hotta, Y. & Bassel, A. Molecular Size and Circularity of DNA in Cells of Mammals and
478 Higher Plants. *P Natl Acad Sci USA* **53**, 356-&, doi:DOI 10.1073/pnas.53.2.356 (1965).
479 2 Verhaak, R. G. W., Bafna, V. & Mischel, P. S. Extrachromosomal oncogene amplification
480 in tumour pathogenesis and evolution. *Nat Rev Cancer* **19**, 283-288,
481 doi:10.1038/s41568-019-0128-6 (2019).

- 482 3 Shoshani, O. *et al.* Chromothripsis drives the evolution of gene amplification in cancer.
483 *Nature* **591**, 137-141, doi:10.1038/s41586-020-03064-z (2021).
- 484 4 Koche, R. P. *et al.* Extrachromosomal circular DNA drives oncogenic genome
485 remodeling in neuroblastoma (vol 52, pg 29, 2019). *Nat Genet* **52**, 464-464,
486 doi:10.1038/s41588-020-0598-1 (2020).
- 487 5 Turner, K. M. *et al.* Extrachromosomal oncogene amplification drives tumour evolution
488 and genetic heterogeneity. *Nature* **543**, 122-+, doi:10.1038/nature21356 (2017).
- 489 6 Wu, S. H. *et al.* Circular ecDNA promotes accessible chromatin and high oncogene
490 expression. *Nature* **575**, 699-+, doi:10.1038/s41586-019-1763-5 (2019).
- 491 7 Kim, H. *et al.* Extrachromosomal DNA is associated with oncogene amplification and
492 poor outcome across multiple cancers. *Nat Genet* **52**, 891-+, doi:10.1038/s41588-020-
493 0678-2 (2020).
- 494 8 deCarvalho, A. C. *et al.* Discordant inheritance of chromosomal and extrachromosomal
495 DNA elements contributes to dynamic disease evolution in glioblastoma. *Nat Genet* **50**,
496 708-+, doi:10.1038/s41588-018-0105-0 (2018).
- 497 9 Li, K. Mortality and incidence trends from esophagus cancer in selected geographic
498 areas of china circa 1970-90. *Int J Cancer* **102**, 271-274, doi:10.1002/jhc.10706 (2002).
- 499 10 Wang, L. D., Zhou, Q. & Yang, C. S. Esophageal and gastric cardia epithelial cell
500 proliferation in northern Chinese subjects living in a high-incidence area. *J Cell Biochem*
501 *Suppl* **28-29**, 159-165 (1997).
- 502 11 Wang, L. D. *et al.* Genome-wide association study of esophageal squamous cell
503 carcinoma in Chinese subjects identifies susceptibility loci at PLCE1 and C20orf54. *Nat*
504 *Genet* **42**, 759-U746, doi:10.1038/ng.648 (2010).
- 505 12 Hu, N. *et al.* Genomic Landscape of Somatic Alterations in Esophageal Squamous Cell
506 Carcinoma and Gastric Cancer. *Cancer Res* **76**, 1714-1723, doi:10.1158/0008-
507 5472.Can-15-0338 (2016).
- 508 13 James D. Brierley, M. K. G., Christian Wittekind. *TNM Classification of Malignant*
509 *Tumours, 8th Edition.* (Wiley-Blackwell, 2016).
- 510 14 Li, Y., Li, J. & Li, J. Two updates on oesophagogastric junction adenocarcinoma from
511 the fifth WHO classification: Alteration of definition and emphasis on HER2 test. *Histol*
512 *Histopathol* **36**, 339-346, doi:10.14670/HH-18-296 (2021).
- 513 15 Maric, R. & Cheng, K. K. Classification of adenocarcinoma of the oesophagogastric
514 junction. *Br J Surg* **86**, 1098-1099, doi:10.1046/j.1365-2168.1999.01197-15.x (1999).
- 515 16 Moureau-Zabotto, L. *et al.* Impact of the Siewert Classification on the Outcome of
516 Patients Treated by Preoperative Chemoradiotherapy for a Nonmetastatic
517 Adenocarcinoma of the Oesophagogastric Junction. *Gastroenterol Res Pract* **2015**,
518 404203, doi:10.1155/2015/404203 (2015).
- 519 17 Siewert, J. R. & Stein, H. J. Classification of adenocarcinoma of the oesophagogastric
520 junction. *Br J Surg* **85**, 1457-1459, doi:10.1046/j.1365-2168.1998.00940.x (1998).
- 521 18 Wang, L. D., Zheng, S., Zheng, Z. Y. & Casson, A. G. Primary adenocarcinomas of lower
522 esophagus, esophagogastric junction and gastric cardia: in special reference to China.
523 *World J Gastroenterol* **9**, 1156-1164, doi:10.3748/wjg.v9.i6.1156 (2003).
- 524 19 Y. Guanrei, a. Q. S. Incidence Rate of Adenocarcinoma of the Gastric Cardia, and
525 Endoscopic Classification of Early Cardial Carcinoma in Henan Province, the People's
526 Republic of China. *Endoscopy* **19**, 7-10 (1987).

- 527 20 Frankell, A. M. *et al.* The landscape of selection in 551 esophageal adenocarcinomas
528 defines genomic biomarkers for the clinic. *Nat Genet* **51**, 506-516,
529 doi:10.1038/s41588-018-0331-5 (2019).
- 530 21 Deshpande, V. *et al.* Exploring the landscape of focal amplifications in cancer using
531 AmpliconArchitect. *Nat Commun* **10**, doi:ARTN 392
532 10.1038/s41467-018-08200-y (2019).
- 533 22 Suh, Y. S. *et al.* Comprehensive Molecular Characterization of Adenocarcinoma of the
534 Gastroesophageal Junction Between Esophageal and Gastric Adenocarcinomas. *Ann*
535 *Surg*, doi:10.1097/SLA.0000000000004303 (2020).
- 536 23 Kumar, P. *et al.* Normal and Cancerous Tissues Release Extrachromosomal Circular
537 DNA (eccDNA) into the Circulation. *Mol Cancer Res* **15**, 1197-1205, doi:10.1158/1541-
538 7786.Mcr-17-0095 (2017).
- 539 24 Dillon, L. W. *et al.* Production of Extrachromosomal MicroDNAs Is Linked to Mismatch
540 Repair Pathways and Transcriptional Activity. *Cell Rep* **11**, 1749-1759,
541 doi:10.1016/j.celrep.2015.05.020 (2015).
- 542 25 Shibata, Y. *et al.* Extrachromosomal MicroDNAs and Chromosomal Microdeletions in
543 Normal Tissues. *Science* **336**, 82-86, doi:10.1126/science.1213307 (2012).
- 544 26 Wong, S. S. *et al.* Genomic landscape and genetic heterogeneity in gastric
545 adenocarcinoma revealed by whole-genome sequencing. *Nat Commun* **5**, 5477,
546 doi:10.1038/ncomms6477 (2014).
- 547 27 Sergina, N. V. & Moasser, M. M. The HER family and cancer: emerging molecular
548 mechanisms and therapeutic targets. *Trends Mol Med* **13**, 527-534,
549 doi:10.1016/j.molmed.2007.10.002 (2007).
- 550 28 Nielsen, T. O., Friis-Hansen, L., Poulsen, S. S., Federspiel, B. & Sorensen, B. S. Expression
551 of the EGF Family in Gastric Cancer: Downregulation of HER4 and Its Activating Ligand
552 NRG4. *Plos One* **9**, doi:ARTN e94606
553 10.1371/journal.pone.0094606 (2014).
- 554 29 Moller, H. D. Circle-Seq: Isolation and Sequencing of Chromosome-Derived Circular
555 DNA Elements in Cells. *Methods Mol Biol* **2119**, 165-181, doi:10.1007/978-1-0716-
556 0323-9_15 (2020).
- 557 30 Duttke, S. H., Chang, M. W., Heinz, S. & Benner, C. Identification and dynamic
558 quantification of regulatory elements using total RNA. *Genome Res* **29**, 1836-1846,
559 doi:10.1101/gr.253492.119 (2019).
- 560 31 Forment, J. V., Kaidi, A. & Jackson, S. P. Chromothripsis and cancer: causes and
561 consequences of chromosome shattering. *Nat Rev Cancer* **12**, 663-670,
562 doi:10.1038/nrc3352 (2012).
- 563 32 Voronina, N. *et al.* The landscape of chromothripsis across adult cancer types. *Nat*
564 *Commun* **11**, 2320, doi:10.1038/s41467-020-16134-7 (2020).
- 565 33 Cortes-Ciriano, I. *et al.* Comprehensive analysis of chromothripsis in 2,658 human
566 cancers using whole-genome sequencing. *Nat Genet* **52**, 331+, doi:10.1038/s41588-
567 019-0576-7 (2020).
- 568 34 Cortes-Ciriano, I. *et al.* Comprehensive analysis of chromothripsis in 2,658 human
569 cancers using whole-genome sequencing. *Cancer Res* **78**, doi:10.1158/1538-
570 7445.Am2018-Lb-378 (2018).
- 571 35 Zhang, C. Z. *et al.* Chromothripsis from DNA damage in micronuclei. *Nature* **522**, 179-
572 184, doi:10.1038/nature14493 (2015).

- 573 36 Przybytkowski, E. *et al.* Chromosome-breakage genomic instability and chromothripsis
574 in breast cancer. *BMC Genomics* **15**, 579, doi:10.1186/1471-2164-15-579 (2014).
- 575 37 Zhang, C. Z., Leibowitz, M. L. & Pellman, D. Chromothripsis and beyond: rapid genome
576 evolution from complex chromosomal rearrangements. *Genes Dev* **27**, 2513-2530,
577 doi:10.1101/gad.229559.113 (2013).
- 578 38 Ratnaparkhe, M. *et al.* Defective DNA damage repair leads to frequent catastrophic
579 genomic events in murine and human tumors. *Nat Commun* **9**, 4760,
580 doi:10.1038/s41467-018-06925-4 (2018).
- 581 39 Umbreit, N. T. *et al.* Mechanisms generating cancer genome complexity from a single
582 cell division error. *Science* **368**, doi:10.1126/science.aba0712 (2020).
- 583 40 Lindor, N. M. *et al.* Immunohistochemistry versus microsatellite instability testing in
584 phenotyping colorectal tumors. *J Clin Oncol* **20**, 1043-1048,
585 doi:10.1200/JCO.2002.20.4.1043 (2002).
- 586 41 Chen, L. Z., Chen, G., Zheng, X. W. & Chen, Y. Expression status of four mismatch repair
587 proteins in patients with colorectal cancer: clinical significance in 1238 cases. *Int J Clin*
588 *Exp Patho* **12**, 3685-3699 (2019).
- 589 42 Birkbak, N. J. *et al.* Paradoxical Relationship between Chromosomal Instability and
590 Survival Outcome in Cancer. *Cancer Res* **71**, 3447-3452, doi:10.1158/0008-5472.Can-
591 10-3667 (2011).
- 592 43 Turinetto, V. & Giachino, C. Multiple facets of histone variant H2AX: a DNA double-
593 strand-break marker with several biological functions. *Nucleic Acids Res* **43**, 2489-2498,
594 doi:10.1093/nar/gkv061 (2015).
- 595 44 Taylor, P. R. *et al.* Prevention of Esophageal Cancer - the Nutrition Intervention Trials
596 in Linxian, China. *Cancer Res* **54**, S2029-S2031 (1994).
- 597 45 Weitberg, A. B. & Corvese, D. Effect of vitamin E and beta-carotene on DNA strand
598 breakage induced by tobacco-specific nitrosamines and stimulated human phagocytes.
599 *J Exp Clin Cancer Res* **16**, 11-14 (1997).
- 600 46 Wang, L. *et al.* Mutations of O6-methylguanine-DNA methyltransferase gene in
601 esophageal cancer tissues from Northern China. *Int J Cancer* **71**, 719-723,
602 doi:10.1002/(sici)1097-0215(19970529)71:5<719::aid-ijc5>3.0.co;2-u (1997).
- 603 47 Deng, C. *et al.* Genetic polymorphism of human O6-alkylguanine-DNA alkyltransferase:
604 identification of a missense variation in the active site region. *Pharmacogenetics* **9**, 81-
605 87, doi:10.1097/00008571-199902000-00011 (1999).
- 606 48 Groenen, P. J. & Busink, E. Alkylating Activity in Food-Products - Especially Sauerkraut
607 and Sour Fermented Dairy-Products after Incubation with Nitrite under Quasi-Gastric
608 Conditions. *Food Chem Toxicol* **26**, 215-225, doi:Doi 10.1016/0278-6915(88)90122-6
609 (1988).
- 610 49 Duell, E. J. *et al.* Vitamin C transporter gene (SLC23A1 and SLC23A2) polymorphisms,
611 plasma vitamin C levels, and gastric cancer risk in the EPIC cohort. *Genes Nutr* **8**, 549-
612 560, doi:10.1007/s12263-013-0346-6 (2013).
- 613 50 Hodgson, R. M., Wiessler, M. & Kleihues, P. Preferential methylation of target organ
614 DNA by the oesophageal carcinogen N-nitrosomethylbenzylamine. *Carcinogenesis* **1**,
615 861-866, doi:10.1093/carcin/1.10.861 (1980).
- 616 51 Verhaak, R. G. W., Bafna, V. & Mischel, P. S. Extrachromosomal oncogene amplification
617 in tumour pathogenesis and evolution. *Nat Rev Cancer* **19**, 283-288,
618 doi:10.1038/s41568-019-0128-6 (2019).

- 619 52 Wu, S. *et al.* Circular ecDNA promotes accessible chromatin and high oncogene
620 expression. *Nature* **575**, 699-703, doi:10.1038/s41586-019-1763-5 (2019).
- 621 53 Park, J. B., Rhim, J. S., Park, S. C., Kimm, S. W. & Kraus, M. H. Amplification,
622 Overexpression, and Rearrangement of the ErbB-2 Protooncogene in Primary Human
623 Stomach Carcinomas. *Cancer Res* **49**, 6605-6609 (1989).
- 624 54 Huang, J. X. *et al.* HER2 gene amplification in esophageal squamous cell carcinoma is
625 less than in gastroesophageal junction and gastric adenocarcinoma. *Oncol Lett* **6**, 13-
626 18, doi:10.3892/ol.2013.1348 (2013).
- 627 55 Houldsworth, J., Cordoncardo, C., Ladanyi, M., Kelsen, D. P. & Chaganti, R. S. K. Gene
628 Amplification in Gastric and Esophageal Adenocarcinomas. *Cancer Res* **50**, 6417-6422
629 (1990).
- 630 56 Plum, P. S. *et al.* HER2/neu (ERBB2) expression and gene amplification correlates with
631 better survival in esophageal adenocarcinoma. *BMC Cancer* **19**, 38,
632 doi:10.1186/s12885-018-5242-4 (2019).
- 633

634 **Supplementary Figures 1-14 are in the separated file.**

635 **Supplementary Table 1-8 is in the separated file.**

636 **Materials and Methods**

637 **GCA samples collection and follow-up visiting of patients**

638 All clinical samples were collected following the ethic permit from the local hospitals located at
639 high-incidence areas of GCA in the Taihang Mountains of north central China. All patients in
640 our study were not received radiotherapy or chemotherapy before the surgery. 1668 GCA
641 patients for ERBB2 immunohistochemistry (IHC) staining are from the Esophageal Cancer
642 database (from years of 1973-2020) which established and maintained by Henan Key
643 Laboratory for Esophageal Cancer Research of the First Affiliated Hospital, Zhengzhou
644 University, China¹⁻⁴. In our Esophageal cancer database, Clinical GCA tumors and matched
645 normal tissues are both preserved with snap freezing in liquid nitrogen and archived in
646 formalin-fixed paraffin-embedded (FFPE) tissue block for each GCA patient. In the studies of
647 whole genome sequencing (WGS), whole exome sequencing (WES), RNA-Seq and protein
648 expression measurement with mass spectrometry, snap freezing samples were used. In the
649 study of IHC staining, FFPE samples were applied. The diagnosis of GCA patients were
650 always identified by two well-trained pathologists in the pathology department of the local
651 hospital, where the hematoxylin and eosin (HE) staining was used to quantify the content of
652 tumor cell in tissue section and only GCA samples with more than 80% tumor cells are used
653 for our study. The matched normal tissue samples were selected from the adjacent epithelial
654 tissue which is 5-10 cm away from the edge of tumor. Both of 36 pairs of GCA tumor and
655 matched adjacent normal tissue for whole-genome sequencing (WGS) and 75 pairs of GCA
656 tumor and matched adjacent normal tissue for whole-exon sequencing (WES) are scanned
657 and confirmed with two well-trained pathologists in the same procedure. The complete
658 clinicopathological information of all patients was recorded and included in our study. All
659 patients are included in regular follow-up visiting plan with following frequency: once every
660 three months during the first year, once each 6 months during the second year, and once per
661 year after the third year. The definition of overall survival time for dead patients is a period from
662 diagnosis to death, and the definition of overall survival time for alive patients is a period from
663 diagnosis to last follow-up visit (Jan 2021).

664

665 **WGS library preparation and sequencing**

666 WGS sequencing libraries were prepare following the previous report with slight modifications⁵.
667 In brief, genomic DNA was extracted from snap freezing GCA tumor and matched normal
668 tissue with DNeasy Blood & Tissue Kit (69504, QIAGEN) following manufacturer instruction.
669 DNA concentration was measured by Qubit DNA Assay Kit in Qubit 2.0 Fluorometer (Invitrogen).
670 A total amount of 0.4µg DNA per sample was fragmented to an average size of ~350bp with

671 hydrodynamic shearing system (Covaris, Massachusetts, USA) and subjected to DNA library
672 preparation with Illumina TruSeq DNA sample preparation kit (15026486, Illumina).
673 Sequencing was carried out on Illumina NovaSeq 6000 with 150bp paired end mode according
674 to the manufacturer instruction.

675

676 **WES library preparation and sequencing**

677 WES sequencing libraries were prepared following the previous report with slight modifications⁶.
678 In brief, genomic DNA was extracted from snap freezing GCA tumor or matched normal tissue
679 using DNeasy Blood & Tissue Kit (69504, QIAGEN) according to the manufacturer's instruction.
680 DNA degradation and contamination were monitored on 1% agarose gels. DNA concentration
681 was measured by Qubit DNA Assay Kit in Qubit 2.0 Fluorometer (Invitrogen). A total amount of
682 0.6 µg genomic DNA per sample was fragmented to an average size of 180~280bp and
683 subjected to DNA library preparation using Illumina TruSeq DNA sample preparation kit. The
684 Agilent SureSelect Human All ExonV5 Kit (5190-6209, Agilent Technologies) was used for
685 exome capture according to the manufacturer's instruction. In brief, DNA libraries were
686 hybridized with liquid phase with biotin labeled probes from the Agilent SureSelect Human All
687 ExonV5 Kit, then magnetic streptavidin beads were used to capture the exons of genes.
688 Captured DNA fragments were enriched in a PCR reaction with index barcodes for sequencing.
689 Final libraries were purified using AMPure XP beads (A63880, Beckman Coulter) and
690 quantified using the Agilent high sensitivity DNA kit (5067-4626, Agilent Technologies). WES
691 libraries were sequenced on Illumina Novaseq 6000 (Illumina) with 150bp paired end mode
692 according to the manufacturer instruction.

693

694 **Circle-Seq library preparation and sequencing**

695 EcDNA sequencing Service was provided by CloudSeq Biotech Inc. (Shanghai, China) by
696 following the published procedures with slight modification⁷. Circle-Seq was performed on 10
697 pairs of snap freezing GCA tumors and matched normal tissues. In brief, 6 mg of snap freezing
698 GCA tumors or matched normal tissues tissue were suspended in L1 solution (A&A
699 Biotechnology, 010-50) and supplemented with 15 µl proteinase K (ThermoFisher, E00491)
700 before incubation overnight at 50 °C with agitation. After Lysis, samples were alkaline treated,
701 followed by precipitation of proteins and separation of chromosomal DNA from circular DNA
702 through an ion exchange membrane column (Plasmid Mini AX; A&A Biotechnology, 010-50).
703 Column-purified DNA was treated with FastDigest MssI (ER1341, Thermo Scientific,) to
704 remove mitochondrial circular DNA and incubated at 37 °C for 16 h. Remaining linear DNA
705 was removed by exonuclease (E3101K, Plasmid-Safe ATP-dependent DNase, Epicentre,) at
706 37 °C in a heating block and enzyme reaction was carried out continuously for 1 week, adding
707 additional ATP and DNase every 24 h (30 units per day) according to the manufacturer's

708 protocol (E3101K, Plasmid-Safe ATP-dependent DNase, Epicentre,). ecDNA-enriched
709 samples were used as template for phi29 polymerase amplification reactions (150043, REPLI-
710 g Midi Kit) amplifying ecDNA at 30 °C for 2 days (46–48 h). Phi29-amplified DNA was sheared
711 by sonication (Bioruptor), and the fragmented DNA was subjected to library preparation with
712 NEBNext® Ultra II DNA Library Prep Kit for Illumina (E7645S, New England Biolabs).
713 Sequencing was carried out on Illumina NovaSeq 6000 with 150bp paired end mode.

714

715 ***ERBB2* RNA expression measurement and *ERBB2* protein expression measurement in** 716 **GCA patients**

717 *ERBB2* RNA expression measurement and *ERBB2* protein expression measurement in 44
718 GCA patients from our Esophageal Cancer database (from years of 1973-2020), where
719 *ERBB2* RNA expression (Normalized value with RPKM (Reads Per Kilobase Million)) was
720 extracted from RNA-seq data, and *ERBB2* protein expression was extracted from mass
721 spectrometry. For same GCA patient, both library for RNA-seq and library for mass
722 spectrometry (MS) are prepared. The procedures of libraries preparation are briefly
723 described as below. For RNA-seq library preparation: First, 100mg of each snap freezing
724 GCA tumor tissue was used for total RNA isolation with TRIzol® Reagent (15596026,
725 Thermo Fisher Scientific). RNA purity was checked using the NanoPhotometer®
726 spectrophotometer (IMPLEN). RNA concentration was measured using Qubit® RNA Assay
727 Kit in Qubit® 2.0 Fluorometer (Life Technologies). RNA integrity was assessed using the
728 Bioanalyzer 2100 system (Agilent Technologies). Then, two RNA-seq libraries were prepared
729 for each GCA patients with technical replicates. 50ng total RNA was used as input for each
730 RNA library preparation. The RNA-Seq libraries were prepared with NEBNext® Ultra™ RNA
731 Library Prep Kit for Illumina (E7530L, NEB) by following manufacturer's instruction. RNA-seq
732 libraries were purified with AMPure XP beads (A63880, Beckman Coulter) to select 150~200
733 bp cDNA fragments. Sequencing library was quantified on the Bioanalyzer 2100 system
734 (Agilent Technologies). The libraries were sequenced on an Illumina Novaseq 6000 platform
735 with 150 bp paired-end reads. The RNA-seq sequencing libraries were aligned to the
736 genome using STAR⁸ with default parameter to reference genome (hg19). After the
737 alignment, the *ERBB2* RNA expression are extracted, and normalized with RPKM. The final
738 expression data for individual patient used to compare with protein expression is the average
739 value of two technical replicates. For mass spectrometry library preparation: First, 10 mg of
740 snap freezing GCA tumor tissues were grinded with liquid nitrogen into powder and then
741 transferred to a 5-mL centrifuge tube. After that, four volumes of lysis buffer (1% Triton X-
742 100, 1% protease inhibitor cocktail, 1% phosphatase inhibitor) was added to the cell powder,
743 followed by sonication three times on ice using a high intensity ultrasonic processor
744 (Scientz). The remaining debris was removed by centrifugation at 12,000 g at 4 °C for 10

745 min. After centrifugation, the supernatant was collected and the protein concentration was
746 measured with Pierce™ BCA protein kit (23227, Thermo Fisher Scientific) according to the
747 manufacturer's instruction. Then, the 100 µg of protein from each sample was taken for
748 protein digestion, and the volume was adjusted to the same with lysate. The sample was
749 slowly added to the final concentration of 20% v/v trichloroacetic acid (TCA) to precipitate
750 protein, then vortexed to mix and incubated for 2hs at 4 °C. The precipitated protein was
751 collected by centrifugation at 4500 g for 5 min at 4 °C. The precipitated protein was washed
752 with pre-cooled acetone for 3 times to remove traces of TCA and finally acetone was
753 removed by drying in a fume cupboard. The protein sample was then added 100 mM
754 Triethylammonium bicarbonate (TEAB) and ultrasonically dispersed. Trypsin was added at
755 1:50 trypsin-to-protein mass ratio for the first digestion overnight. The sample was reduced
756 with 5 mM dithiothreitol for 30 min at 56 °C and alkylated with 11 mM iodoacetamide for 15
757 min at room temperature in darkness. Next, 50 µg of tryptic peptides were firstly dissolved in
758 0.5 M TEAB. Each channel of peptide was labeled with their respective TMT reagent, and
759 incubated for 2 hours at room temperature. Five microliters of each sample were pooled,
760 desalted and analyzed by MS to check labeling efficiency. After labeling efficiency check,
761 samples were quenched by adding 5% hydroxylamine. The pooled samples were then
762 desalted with Strata X C18 SPE column (Phenomenex) and dried by vacuum centrifugation.
763 Then, the dried tryptic peptides were dissolved in solvent A (0.1% formic acid, 2%
764 acetonitrile/ in water), directly loaded onto a home-made reversed-phase analytical column
765 (25-cm length, 100 µm i.d.). Peptides were separated with a gradient from 8% to 10% solvent
766 B (0.1% formic acid in 90% acetonitrile) over 2 min, 10% to 23% solvent B over 38 min, 23%
767 to 33% in 14 min and climbing to 80% in 3 min then holding at 80% for the last 3 min, all at a
768 constant flowrate of 450 nL/min on an EASY-nLC 1200 UPLC system (Thermo Fisher
769 Scientific). The separated peptides were analyzed in Q Exactive™ HF-X (Thermo Fisher
770 Scientific) with a nano-electrospray ion source. The electrospray voltage applied was 2.2 kV.
771 The full MS scan resolution was set to 120,000 for a scan range of 400–1500 m/z. Up to 20
772 most abundant precursors were then selected for further MS/MS analyses with 30 s dynamic
773 exclusion. The HCD fragmentation was performed at a normalized collision energy (NCE) of
774 28%. The fragments were detected in the Orbitrap at a resolution of 45,000. Fixed first mass
775 was set as 100 m/z. Automatic gain control (AGC) target was set at 5E4, with an intensity
776 threshold of 5.8E4 and a maximum injection time of 86 ms. The resulting MS/MS data were
777 processed using MaxQuant search engine (v.1.6.10.43). Tandem mass spectra were
778 searched against the human SwissProt database (20366 entries) concatenated with reverse
779 decoy database. Trypsin/P was specified as cleavage enzyme allowing up to 2 missing
780 cleavages. The mass tolerance for precursor ions was set as 10 ppm in First search and 5
781 ppm in Main search, and the mass tolerance for fragment ions was set as 0.02 Da.

782 Carbamidomethyl on Cys was specified as fixed modification. Acetylation on protein N-
783 terminal, oxidation on Met and deamidation (NQ) were specified as variable modifications.
784 TMT-11plex quantification was performed. FDR was adjusted to < 1% and minimum score
785 for peptides was set > 40. The ERBB2 protein expression level for each patient was
786 extracted from protein lists of MS result.

787

788 **Data analysis of WGS data, WES data, copy number alteration (CNA) and ecDNA** 789 **amplicons**

790 All detailed scripts were deposited in following link: [https://github.com/chenlab2019/ecDNA-](https://github.com/chenlab2019/ecDNA-on-GCA)
791 [on-GCA](https://github.com/chenlab2019/ecDNA-on-GCA). The WGS data of 36 samples were aligned to the reference genome (hg19) using
792 BWA-MEM v.0.7.17⁹ with the default parameter and were sorted by SAMtools v.1.9¹⁰. PCR
793 duplicates were removed from aligned BAM files by Sambamba v.0.7.0¹¹. By taking matched
794 normal samples as background, tumor-specific CNAs were called by copyCat package
795 (<https://github.com/chrisamiller/copyCat>) which is loosely based on readDepth¹². During the
796 process of CNA calling, bam-window tools (<https://github.com/genome-vendor/bam-window>)
797 was used to count reads in 10Kbp window size. AmpliconArchitect (AA) was applied to filter
798 CNAs with copy number greater than 4x and size greater than 100Kbp. The adjacent CNAs
799 were merged into a single interval. These intervals were fed into AmpliconArchitect software¹³
800 as seeds to detect ecDNA amplicons¹⁴. The oncogene annotation of ecDNA amplicons was
801 based on the genome intervals of amplicons following AA pipeline¹³. The genomic annotation
802 of ecDNA amplicons was performed with intersection between regions of ecDNA amplicons
803 and genomic annotation of reference genome (hg19) with bedtools¹⁵. In brief, regions of the
804 ecDNA amplicons were extracted from the output of AA software. The intersection between
805 genomic annotation of reference genome (hg19) and ecDNA regions was performed with
806 bedtools first¹⁵, then the length of overlapping regions between genomic elements from
807 reference genome and ecDNA regions was extracted. Genomic elements were annotated to
808 ecDNA amplicons if there was one bp or longer overlapping. The occupancy of coding regions
809 and exons regions in ecDNA amplicons were calculated with following formulas:

810

$$811 \quad \text{occupancy of coding regions in ecDNA (\%)} = \frac{\text{total length of coding regions} \\ \text{in all ecDNA amplicons}}{\text{total length of} \\ \text{all ecDNA amplicons}} \times 100$$

812

813

$$814 \quad \text{occupancy of exon regions in ecDNA (\%)} = \frac{\text{total length of exon regions} \\ \text{in all ecDNA amplicons}}{\text{total length of} \\ \text{all ecDNA amplicons}} \times 100$$

815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851

EcDNA amplicons were further classified into different categories (linear, complex, circular, breakage-fusion-bridge (BFB) and invalid) with AA software (<https://github.com/jluebeck/AmpliconClassifier>) by following the previous report¹⁶. Circle-finder¹⁷⁻¹⁹ was used to confirm the circular structure of ecDNA amplicons by following the instruction, where circular junction points were detected with sequencing reads orientation. The length of overlapping region between circular ecDNA predicted from AA and circular ecDNA detected with Circle-finder was calculated with bedtools. When the length of overlapping region is longer than 1bp, circular ecDNA amplicons from AA were labelled as overlapping with results of Circle-finder.

For WES data analysis from 75 pairs of GCA tumor samples and matched adjacent normal tissues: sequencing reads containing adaptors and low-quality reads were removed and aligned to human reference genome (hg19) using BWA-MEM v.0.7.17⁹ with the default parameter and sorted by SAMtools v.1.9¹⁰. All non-primary alignments were filtered by SAMtools. PCR duplicates were marked using Picard. CNAs from tumor was called by using matched adjacent normal tissues by CNVkit²⁰. The numbers of CNAs on ERBB2 gene from each GCA patient are extracted for further analysis.

Data mining of Circle-seq

All reads were aligned to human genome hg19 using BWA-MEM v.0.7.17⁹ with default parameters. PCR duplicates were removed from the BAM file with Sambamba v.0.7.0⁹. By taking normal samples as background, peak calling on tumor samples was performed using variable-width windows of Homer v.4.11 with command *findPeaks tumor -i normal -style histone -fdr 0.001* (<http://homer.ucsd.edu/>)²³. The tumor-specific enriched peaks were considered as the fragments of circular DNA. Overlaps between enriched peaks from Circle-Seq and ecDNA amplicons from AA were calculated, and circular ecDNA amplicon from AA is labelled as validated when the overlapping regions is 1bp or longer than 1bp. For the visualization of the peak of Circle-Seq, BAM file was converted into bigwig file using deeptools bamCoverage with normalization of counts per million (CPM)²⁴.

Detection of chromothripsis events

All detailed scripts were deposited in following link: <https://github.com/chenlab2019/ecDNA-on-GCA>. Chromothripsis events from 36 pairs of GCA tumor samples were detected with ShatterSeek software v.0.4 using copy number alterations (CNAs) and structural variants (SVs) following the previous report²⁵. SVs were identified on tumor samples using the Delly²⁶ and

852 novoBreak²⁷ software by taking matched adjacent normal tissues as control, and final list of
853 SVs are merged lists from Delly and NovoBreak. CNAs from WGS were calculated with
854 copyCat package²⁸. All SVs and CNVs from tumor samples are used to identify chromothripsis
855 events with ShatterSeek, where SVs and CNVs from matched adjacent normal tissues are
856 treated as background. Events were considered as high confidence (termed HC) when there
857 were at least 7 oscillating CN segments, and considered as low confidence (termed LC) when
858 there were 4-6 oscillating CN segments¹¹. The chromothripsis events were labeled as within
859 regions of ecDNA amplicons when there is 1bp or longer intersection between segments from
860 chromothripsis and regions of ecDNA amplicons.

861

862 **Single-nucleotide variant (SNV) analysis**

863 All detailed scripts were deposited in following link: [https://github.com/chenlab2019/ecDNA-](https://github.com/chenlab2019/ecDNA-on-GCA)
864 [on-GCA](https://github.com/chenlab2019/ecDNA-on-GCA). All SNVs from WGS were called by GATK v.4.1.7 software²⁹ with Mutect2 parameter
865 and filtered by “GATK FilterMutectCalls”. The mutation profiles were visualized by
866 R/Bioconductor package “maftools”³⁰. The number of SNVs within region of ecDNA amplicons
867 and whole genome region were counted respectively for each sample. The average number
868 of SNVs per million nucleotides from regions of ecDNA amplicons and whole genome were
869 calculated with following equations:

870

$$871 \quad \text{SNVs of ecDNA} = \frac{\text{The number of SNVs in ecDNA amplicons}}{\text{The total length of ecDNA amplicons}} \times 1 \text{ million}$$

$$872 \quad \text{SNVs of whole genome} = \frac{\text{The number of SNVs within whole genome}}{\text{The total length of whole genome}} \times 1 \text{ million}$$

873

874 Numbers of SNVs within individual oncogene ecDNA amplicon from groups of absent and
875 present this gene ecDNA were also compared: first high frequency of oncogene ecDNA
876 amplicons (appeared at least in 2 patients) in 36 patients are selected, then the number of
877 SNVs within each selected oncogene from individual patient was calculated and numbers of
878 SNVs between groups of present and absent this oncogene ecDNA were compared.

879

880 **Oncogene ecDNA amplicon analysis**

881 All detailed scripts were deposited in following link: [https://github.com/chenlab2019/ecDNA-](https://github.com/chenlab2019/ecDNA-on-GCA)
882 [on-GCA](https://github.com/chenlab2019/ecDNA-on-GCA). The list of oncogenes and tumor suppressor genes ecDNA amplicons was extracted
883 from the report of AmpliconArchitect following AmpliconArchitect workflow¹³. The copy number
884 of each oncogene from 36 GCA samples was extracted from the report of copyCat.
885 Oncogenes and/or tumor suppressor genes are labeled as oncogene co-amplification if two or

886 more than two oncogenes and/or tumor suppressor genes are located in the same ecDNA
887 amplicon.

888

889 **Calculation of Chromosomal instability (CIN):** All detailed scripts were deposited in
890 following link: <https://github.com/chenlab2019/ecDNA-on-GCA>. The chromosomal instability
891 (CIN) was calculated following the previous report³¹, and groups of chromosomal instability
892 (CIN) is defined with by number of genome integrity index (GII). GII was defined as the fraction
893 of the genome that was altered based on the common regions of alteration. CIN of GCA
894 patients was divided into four groups based on GII (0-0.2, 0.2-0.4, 0.4-0.6, 0.6-0.8), and 36
895 GCA patients were assigned into different groups of CIN.

896

897 **Prognoses and statistical analysis**

898 All computational codes and scripts are deposited in following
899 link: <https://github.com/chenlab2019/ecDNA-on-GCA>. R package “survival” with Kaplan-Meier
900 method was used³² to calculate and compare patient prognosis between different groups of
901 GCA patients. The statistic methods used in prognosis analysis with clinicopathological factor
902 are as follows: Fisher’s exact test for sex, family history cigarette smoking, alcohol consuming
903 and tumor stage, and Wilcoxon signed-rank test for age. All analyses were performed on R
904 v.3.6.2, Python v.2.7.16 and Python 3.7.4. The visualization of survival curve was conducted
905 by ggplot2³³, karyoploteR³⁴, pheatmap R packages and Circos³⁷, IGV software³⁸.

906

907 **Immunohistochemistry (IHC) staining of ERBB2 protein**

908 IHC was performed by following the previous report³⁹ with slightly modifications. In brief, 5- μ m
909 thick formalin fixed paraffin-embedded GCA tissue sections were first deparaffined with
910 xylene 15mins for 3 times, then were dehydrated through 100% alcohol, 85% alcohol and 75%
911 alcohol for 5mins each, followed by distilled water rinsing for 5 mins. The epitope retrieval is
912 performed in the microwave by putting the tissue into citrate buffer (pH 6.0). After the epitope
913 retrieval, the tissue section is rinsed in Phosphate-Buffered Saline buffer (PBS, PH7.4). After
914 blocked with 3% bovine serum albumin (BSA) 30mins at room temperature, the tissues were
915 incubated with ERBB2 antibody (1:100 dilution, SAB5700151, Sigma-Aldrich) overnight at 4°C.
916 In the next day, the washing is performed with PBS buffer for 3 times, 15mins each. The
917 secondary antibody (Horseradish Peroxidase, HRP marked, PV-9000, ZSGB-BIO) was
918 incubated for 50 mins at room temperature. After the secondary antibody incubation, the
919 washing is performed with PBS buffer 3 times on shaker, 15 mins each. The tissue is stained
920 with the Harris Hematoxylin for 3 mins. At last, the tissue section was mounted and imaged.
921 Sections with no signal in any cell were defined as negative groups; sections with 5 or more
922 cells with ERBB2 positive signal were defined as positive groups.

923 **Immunohistochemistry (IHC) staining of γ H2AX:** The staining protocol is same as ERBB2
924 staining. The primary antibody of γ H2AX (SAB5700329, Sigma-Aldrich) was with 1:200 dilution.
925 The staining of γ H2AX was categories into positive and negative groups with following
926 parameters: Section with no γ H2AX signal in any cell was defined as γ H2AX negative groups;
927 section with 5 or more cells with γ H2AX positive signal was defined as γ H2AX positive groups.

928
929 **MSI detection with Immunohistochemistry (IHC) staining:** IHC staining of four mismatch
930 repair (NMR) proteins: MLH1 (1: 100, PA5-32497, Thermo Fisher Scientific), MSH2 (1: 500,
931 MA5-15740, Thermo Fisher Scientific), MSH6 (1: 100, MA5-32040, Thermo Fisher Scientific)
932 and PMS2 (1: 150, MA5-26269, Thermo Fisher Scientific), were performed on 5- μ m thick
933 FFPE tumor sections from 36 GCA patient with same protocol as stated as above in ERBB2
934 IHC staining. The patient was labeled as microsatellite instability (MSI-high) if one of NMR
935 proteins was negative stained, otherwise the patient is labeled as MSI-low.

936

937

938 **Methods only References:**

- 939 1 Wang, L. *et al.* Mutations of O6-methylguanine-DNA methyltransferase gene in
940 esophageal cancer tissues from Northern China. *Int J Cancer* **71**, 719-723,
941 doi:10.1002/(sici)1097-0215(19970529)71:5<719::aid-ijc5>3.0.co;2-u (1997).
- 942 2 Wang, L. D., Zheng, S., Zheng, Z. Y. & Casson, A. G. Primary adenocarcinomas of lower
943 esophagus, esophagogastric junction and gastric cardia: in special reference to China.
944 *World J Gastroenterol* **9**, 1156-1164, doi:10.3748/wjg.v9.i6.1156 (2003).
- 945 3 Wang, L. D. *et al.* Genome-wide association study of esophageal squamous cell
946 carcinoma in Chinese subjects identifies susceptibility loci at PLCE1 and C20orf54. *Nat*
947 *Genet* **42**, 759-U746, doi:10.1038/ng.648 (2010).
- 948 4 Wang, L. D., Zhou, Q. & Yang, C. S. Esophageal and gastric cardia epithelial cell
949 proliferation in northern Chinese subjects living in a high-incidence area. *J Cell Biochem*
950 *Suppl* **28-29**, 159-165 (1997).
- 951 5 Meyer, M. & Kircher, M. Illumina sequencing library preparation for highly multiplexed
952 target capture and sequencing. *Cold Spring Harb Protoc* **2010**, pdb prot5448,
953 doi:10.1101/pdb.prot5448 (2010).
- 954 6 Sulonen, A. M. *et al.* Comparison of solution-based exome capture methods for next
955 generation sequencing. *Genome Biol* **12**, R94, doi:10.1186/gb-2011-12-9-r94 (2011).
- 956 7 Moller, H. D. *et al.* Circular DNA elements of chromosomal origin are common in
957 healthy human somatic tissue. *Nat Commun* **9**, 1069, doi:10.1038/s41467-018-03369-
958 8 (2018).
- 959 8 Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21,
960 doi:10.1093/bioinformatics/bts635 (2013).
- 961 9 Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler
962 transform. *Bioinformatics* **25**, 1754-1760, doi:10.1093/bioinformatics/btp324 (2009).
- 963 10 Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,
964 2078-2079, doi:10.1093/bioinformatics/btp352 (2009).

- 965 11 Tarasov, A., Vilella, A. J., Cuppen, E., Nijman, I. J. & Prins, P. Sambamba: fast processing
966 of NGS alignment formats. *Bioinformatics* **31**, 2032-2034,
967 doi:10.1093/bioinformatics/btv098 (2015).
- 968 12 Miller, C. A., Hampton, O., Coarfa, C. & Milosavljevic, A. ReadDepth: a parallel R
969 package for detecting copy number alterations from short sequencing reads. *Plos One*
970 **6**, e16327, doi:10.1371/journal.pone.0016327 (2011).
- 971 13 Luebeck, J. *et al.* AmpliconReconstructor integrates NGS and optical mapping to
972 resolve the complex structures of focal amplifications. *Nat Commun* **11**, 4374,
973 doi:10.1038/s41467-020-18099-z (2020).
- 974 14 Deshpande, V. *et al.* Exploring the landscape of focal amplifications in cancer using
975 AmpliconArchitect. *Nat Commun* **10**, doi:ARTN 392
976 10.1038/s41467-018-08200-y (2019).
- 977 15 Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic
978 features. *Bioinformatics* **26**, 841-842, doi:10.1093/bioinformatics/btq033 (2010).
- 979 16 Kim, H. *et al.* Extrachromosomal DNA is associated with oncogene amplification and
980 poor outcome across multiple cancers. *Nat Genet* **52**, 891-+, doi:10.1038/s41588-020-
981 0678-2 (2020).
- 982 17 Kumar, P. *et al.* Normal and Cancerous Tissues Release Extrachromosomal Circular
983 DNA (eccDNA) into the Circulation. *Mol Cancer Res* **15**, 1197-1205, doi:10.1158/1541-
984 7786.Mcr-17-0095 (2017).
- 985 18 Dillon, L. W. *et al.* Production of Extrachromosomal MicroDNAs Is Linked to Mismatch
986 Repair Pathways and Transcriptional Activity. *Cell Rep* **11**, 1749-1759,
987 doi:10.1016/j.celrep.2015.05.020 (2015).
- 988 19 Shibata, Y. *et al.* Extrachromosomal MicroDNAs and Chromosomal Microdeletions in
989 Normal Tissues. *Science* **336**, 82-86, doi:10.1126/science.1213307 (2012).
- 990 20 Talevich, E., Shain, A. H., Botton, T. & Bastian, B. C. CNVkit: Genome-Wide Copy
991 Number Detection and Visualization from Targeted DNA Sequencing. *PLoS Comput Biol*
992 **12**, e1004873, doi:10.1371/journal.pcbi.1004873 (2016).
- 993 21 Duttke, S. H., Chang, M. W., Heinz, S. & Benner, C. Identification and dynamic
994 quantification of regulatory elements using total RNA. *Genome Res* **29**, 1836-1846,
995 doi:10.1101/gr.253492.119 (2019).
- 996 22 Ramirez, F. *et al.* High-resolution TADs reveal DNA sequences underlying genome
997 organization in flies. *Nat Commun* **9**, 189, doi:10.1038/s41467-017-02525-w (2018).
- 998 23 Heinz, S. *et al.* Simple combinations of lineage-determining transcription factors prime
999 cis-regulatory elements required for macrophage and B cell identities. *Mol Cell* **38**, 576-
1000 589, doi:10.1016/j.molcel.2010.05.004 (2010).
- 1001 24 Ramirez, F., Dundar, F., Diehl, S., Gruning, B. A. & Manke, T. deepTools: a flexible
1002 platform for exploring deep-sequencing data. *Nucleic Acids Res* **42**, W187-191,
1003 doi:10.1093/nar/gku365 (2014).
- 1004 25 Cortes-Ciriano, I. *et al.* Comprehensive analysis of chromothripsis in 2,658 human
1005 cancers using whole-genome sequencing. *Nat Genet* **52**, 331-+, doi:10.1038/s41588-
1006 019-0576-7 (2020).
- 1007 26 Rausch, T. *et al.* DELLY: structural variant discovery by integrated paired-end and split-
1008 read analysis. *Bioinformatics* **28**, i333-i339, doi:10.1093/bioinformatics/bts378 (2012).
- 1009 27 Chong, Z. *et al.* novoBreak: local assembly for breakpoint detection in cancer genomes.
1010 *Nat Methods* **14**, 65-67, doi:10.1038/nmeth.4084 (2017).

- 1011 28 Gao, R. *et al.* Delineating copy number and clonal substructure in human tumors from
1012 single-cell transcriptomes. *Nat Biotechnol* **39**, 599-608, doi:10.1038/s41587-020-
1013 00795-2 (2021).
- 1014 29 McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for
1015 analyzing next-generation DNA sequencing data. *Genome Res* **20**, 1297-1303,
1016 doi:10.1101/gr.107524.110 (2010).
- 1017 30 Mayakonda, A., Lin, D. C., Assenov, Y., Plass, C. & Koeffler, H. P. Maftools: efficient and
1018 comprehensive analysis of somatic variants in cancer. *Genome Res* **28**, 1747-1756,
1019 doi:10.1101/gr.239244.118 (2018).
- 1020 31 Birnbak, N. J. *et al.* Paradoxical Relationship between Chromosomal Instability and
1021 Survival Outcome in Cancer. *Cancer Res* **71**, 3447-3452, doi:10.1158/0008-5472.Can-
1022 10-3667 (2011).
- 1023 32 Li, J. C. A. Modeling survival data: Extending the Cox model. *Sociol Method Res* **32**, 117-
1024 120, doi:Doi 10.1177/0049124103031004005 (2003).
- 1025 33 Villanueva, R. A. M. & Chen, Z. J. ggplot2: Elegant Graphics for Data Analysis, 2nd
1026 edition. *Meas-Interdiscip Res* **17**, 160-167, doi:10.1080/15366367.2019.1565254
1027 (2019).
- 1028 34 Gel, B. & Serra, E. karyoploteR: an R/Bioconductor package to plot customizable
1029 genomes displaying arbitrary data. *Bioinformatics* **33**, 3088-3090,
1030 doi:10.1093/bioinformatics/btx346 (2017).
- 1031 35 Krzywinski, M. *et al.* Circos: An information aesthetic for comparative genomics.
1032 *Genome Res* **19**, 1639-1645, doi:10.1101/gr.092759.109 (2009).
- 1033 36 Thorvaldsdottir, H., Robinson, J. T. & Mesirov, J. P. Integrative Genomics Viewer (IGV):
1034 high-performance genomics data visualization and exploration. *Brief Bioinform* **14**,
1035 178-192, doi:10.1093/bib/bbs017 (2013).
- 1036 37 Krzywinski, M. *et al.* Circos: an information aesthetic for comparative genomics.
1037 *Genome Res* **19**, 1639-1645, doi:10.1101/gr.092759.109 (2009).
- 1038 38 Thorvaldsdottir, H., Robinson, J. T. & Mesirov, J. P. Integrative Genomics Viewer (IGV):
1039 high-performance genomics data visualization and exploration. *Brief Bioinform* **14**,
1040 178-192, doi:10.1093/bib/bbs017 (2013).
- 1041 39 Robertson, D., Savage, K., Reis-Filho, J. S. & Isacke, C. M. Multiple immunofluorescence
1042 labelling of formalin-fixed paraffin-embedded (FFPE) tissue. *BMC Cell Biol* **9**, 13,
1043 doi:10.1186/1471-2121-9-13 (2008).
- 1044

1045 **Figure Legends:**

1046 **Figure 1: Identification and characterization of ecDNA amplicons from whole-genome**
1047 **sequencing data of the GCA cohort.**

1048 **a**, Schematic of the experiment design for detecting ecDNA amplicons from WGS data of 36
1049 pairs of GCA tumour and tumour-adjacent normal tissue from a high incidence GCA rate region
1050 in the northern region of China.

1051 **b**, Detailed characterization of ecDNA amplicons from 36 GCAs, where ecDNA amplicons are
1052 further classified into circular, complex, linear, breakage-fusion-bridge (BFB) and invalid.

1053 **c**, Genomic annotation of all ecDNA amplicons, where the annotation was defined by
1054 overlapping gene regions and regions of ecDNA amplicons. TSG = tumour suppressor gene.

1055 **d**, Distribution of high-frequency oncogene and tumour suppressor gene (TSG) amplicons
1056 across all 36 samples.

1057 **e**, The summary of oncogene ecDNA co-amplification in our cohort, where co-amplification is
1058 defined when two or more than two oncogenes are in the same ecDNA amplicon;

1059 **f**, *EGFR* and *CDK6* are located in the same circular ecDNA amplicon, where the genome
1060 coverage on the left panel represents gene amplification of *EGFR* and *CDK6*, and the circular
1061 structure on the right panel is the reconstruction of *EGFR* and *CDK6* in the same circular
1062 ecDNA.

1063

1064 **Figure 2: Validation of the ecDNA amplicons using Circle-seq.**

1065 **a**, Summary of ecDNA overlapping lists from the prediction of AmpliconArchitect (AA) and
1066 identification using Circle-seq. The y-axis is the ecDNA amplicon number from WGS prediction.
1067 Overlap: the ecDNA amplicons were identified using both AA software from WGS and Circle-
1068 Seq. None: the ecDNA amplicons were only identified using AA software but not using Circle-
1069 Seq.

1070 **b**, The genome browser track at the *FGFR2* gene locus from whole-genome sequencing (WGS)
1071 and Circle-seq. The connection lines on the top represent the potential structure combination
1072 in ecDNA amplicons predicted by AA software. N = normal tissue, T = tumour tissue.

1073

1074 **Figure 3: EcDNA amplicon and chromothripsis in GCA patients.**

1075 **a**, Summary of chromothripsis events across the whole genome in our GCA cohort. HC = high
1076 confidence chromothripsis; LC = low confidence chromothripsis.

1077 **b**, *ERBB2* ecDNA amplicon in the event of chromothripsis from one GCA patient. The different
1078 connection lines on the top represent the potential different formats of chromothripsis events
1079 at the *ERBB2* gene. CN = copy number.

1080 **c**, Summary of overlapping frequency between ecDNA amplicon and chromothripsis in the
1081 GCA cohort. HC = high confidence chromothripsis; LC = low confidence chromothripsis.

1082 **d**, The correlation between total length of ecDNA amplicons and the frequency of
1083 chromothripsis in GCA patients, where each dot represents one sample.

1084 **e**, Representative images of γ H2AX immunohistochemistry (IHC) staining in our GCA cohort.

1085 **f**, Presence and absence of chromothripsis in γ H2AX-positive and γ H2AX-negative groups of
1086 GCA patients. The numbers on the bars are patient numbers.

1087 **g**, Comparisons of the total length of chromothripsis in γ H2AX-positive and γ H2AX-negative
1088 GCA patients, where each dot represents one patient, and the length of chromothripsis is the
1089 total length of all chromosomes in each sample. The p-value was calculated using the
1090 Wilcoxon signed-rank test.

1091

1092 **Figure 4: Oncogene amplification, ecDNA amplicon presence and prognosis of GCA**
1093 **patients**

1094 **a**, The relationship between gene copy number and survival time for *CCNE1*, *EGFR*, and
1095 *ERBB2* genes in the GCA cohort, where copy number of *CCNE1* and *ERBB2* genes were
1096 divided into three groups, High, Low and Normal, and copy number of the *EGFR* gene was
1097 divided into two groups, High and Normal. High = high copy number of gene amplification, Low
1098 = low copy number of gene amplification, Normal = no gene amplification.

1099 **b**, Survival analysis of different groups with three oncogene amplifications (*CCNE1*, *EGFR*,
1100 and *ERBB2*) in the cohort. The definition of High, Low and Normal is the same as in panel **a**,
1101 and the p-value was calculated using the Log rank test.

1102 **c**, Survival time of present and absent ecDNA amplicons of three oncogenes in ecDNA
1103 (*CCNE1*, *EGFR*, and *ERBB2*) in the cohort. The p-value was calculated using the Log rank
1104 test.

1105 **d**, Kaplan-Meier plot for the presence and absence of *ERBB2* protein expression in GCA tissue
1106 sections from 1668 GCA patients. Left panel: survival analysis of patients with 0-2 year survival
1107 after surgery (n = 750); right panel: survival analysis of patients with 2-7 year survival after
1108 surgery (n = 918). The p-value was calculated using the Log rank test.

1109

Figure 1

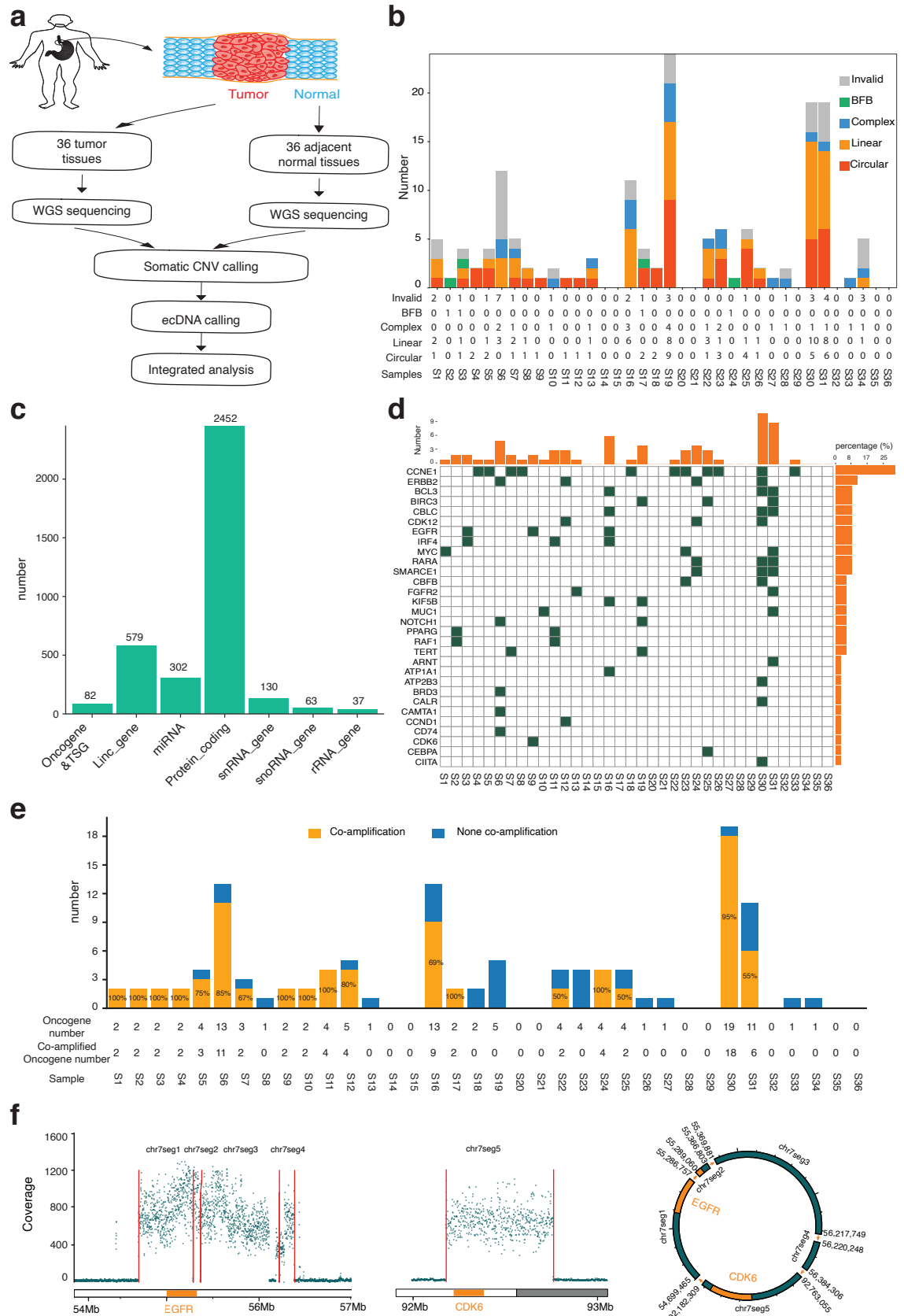
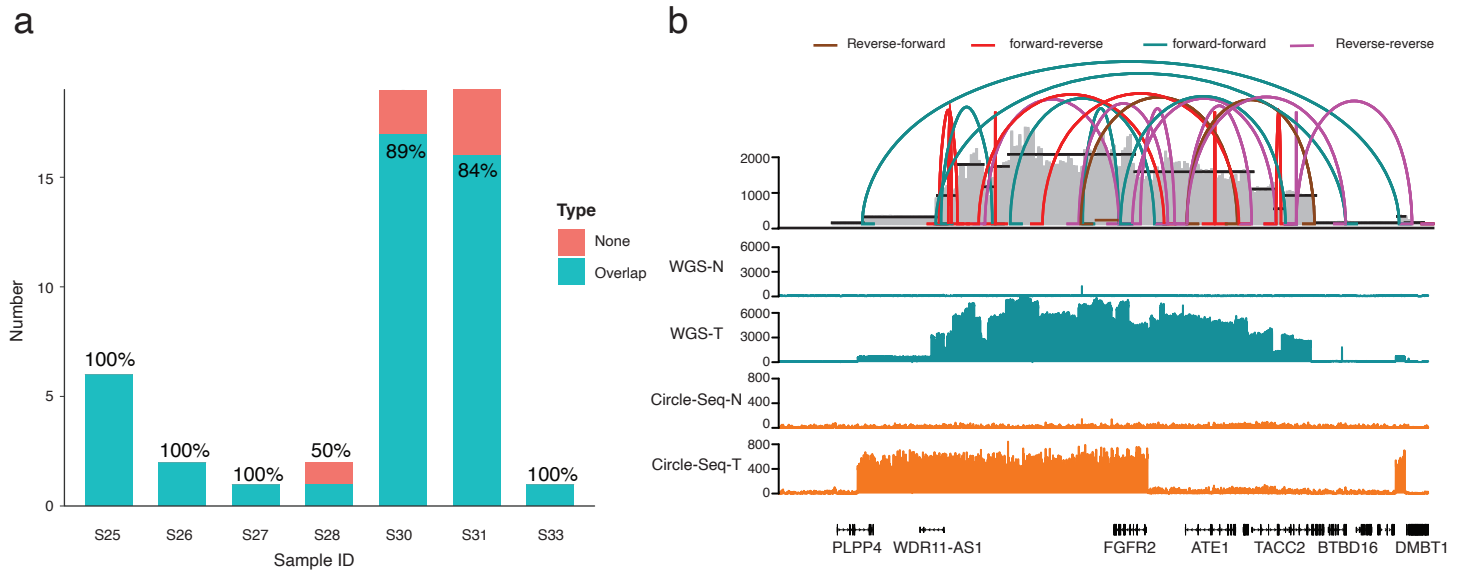


Figure 2



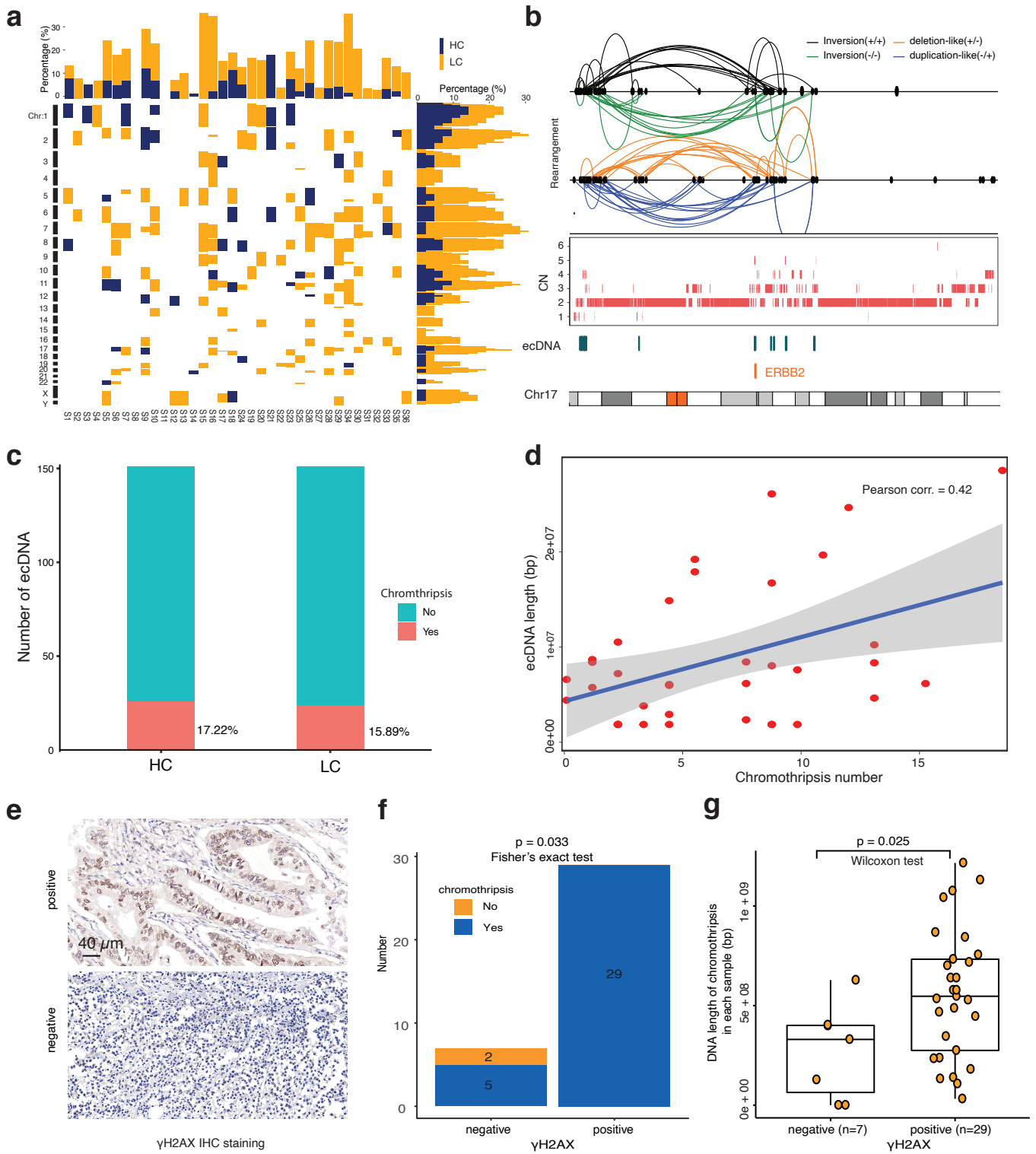


Figure 4

