1 *Frequent extrachromosomal oncogene amplification drives*
2 *aggressive tumors*

3 **Authors:** Hoon Kim[1$], Nam Nguyen[2$], Kristen Turner[3], Sihan Wu[3], Jihe Liu[1],Viraj
4 Deshpande[2], Sandeep Namburi[1], Howard Y. Chang[4,5], Christine Beck[1], Paul
5 Mischel[3,6,7*], Vineet Bafna[2*], Roel Verhaak[1*]

6
7 **Affiliations:**

8 1. The Jackson Laboratory for Genomic Medicine, Farmington, Connecticut 06032, USA

9 2. Department of Computer Science and Engineering, University of California at San
10 Diego, La Jolla, California 92093, USA

11 3. Ludwig Institute for Cancer Research, University of California at San Diego, La Jolla,
12 California 92093, USA

13 4. Center for Personal Dynamic Regulomes, Stanford University, Stanford, California
14 94305, USA

15 5. Howard Hughes Medical Institute, Stanford University, Stanford, CA 94305, USA

16 6. Moores Cancer Center, University of California at San Diego, La Jolla, California
17 92093, USA

18 7. Department of Pathology, University of California at San Diego, La Jolla, California
19 92093, USA

20

21 $-co-first authors

22 *co-corresponding authors

23

24 **Extrachromosomal DNA (ecDNA) amplification promotes high oncogene copy**
25 **number, intratumoral genetic heterogeneity, and accelerated tumor evolution[1-3],**
26 **but its frequency and clinical impact are not well understood. Here we show, using**
27 **computational analysis of whole-genome sequencing data from 1,979 cancer**
28 **patients, that ecDNA amplification occurs in at least 26% of human cancers, of a**
29 **wide variety of histological types, but not in whole blood or normal tissue. We**
30 **demonstrate a highly significant enrichment for oncogenes on amplified ecDNA**
31 **and that the most common recurrent oncogene amplifications arise on ecDNA.**
32 **EcDNA amplifications resulted in higher levels of oncogene transcription**
33 **compared to copy number matched linear DNA, coupled with enhanced chromatin**
34 **accessibility. Patients whose tumors have ecDNA-based oncogene amplification**

35  **showed increase of cell proliferation signature activity, greater likelihood of lymph**
36  **node spread at initial diagnosis, and significantly shorter survival, even when**
37  **controlled for tissue type, than do patients whose cancers are not driven by ecDNA-**
38  **based oncogene amplification. The results presented here demonstrate that**
39  **ecDNA-based oncogene amplification plays a central role in driving the poor**
40  **outcome for patients with some of the most aggressive forms of cancers.**

41  Somatic gain of function alterations in growth controlling genes, especially driver
42  oncogenes, plays a central role in the development of cancer[4-6]. Oncogene amplification
43  is one of the most common gain of function alterations in cancer, enabling tumor cells to
44  circumvent the checks and balances that are in place during homeostasis and providing
45  selective and autonomous advantage to drive tumor growth. EcDNA-based amplification
46  has long been recognized as a way for cells to increase the copy number of specific
47  genes[7,8], but their frequency appears to be vastly underestimated[2,9]. EcDNA amplification
48  has recently emerged as a powerful mechanism for enabling tumors to concomitantly
49  reach high copy of growth promoting genes, while still maintaining intratumoral genetic
50  heterogeneity through its non-chromosomal mechanism of inheritance[1-3]. To date,
51  cytogenetic methods requiring live cells in metaphase have been used to infer
52  intranuclear localization of DNA amplifications and extrachromosomal status[10].
53  Consequently, it has been challenging to accurately assess the frequency, distribution,
54  and clinical impact of ecDNA-based amplification. More recently computational analyses
55  of whole-genome sequencing data have suggested a relatively high frequency of ecDNA
56  in some cancer types[11,12]. Here we set out to perform a global survey of the frequency of
57  ecDNA-based oncogene amplification, while investigating its contents and determining
58  its clinical context.

59  EcDNA are characterized by two distinguishing properties: 1. ecDNAs are highly and
60  focally amplified and  2. they are circular. These properties provide a basis for the
61  AmpliconArchitect tool, that enables detection and characterization of ecDNA from whole-
62  genome sequencing data (**Fig. 1A**)[11]. We applied AmpliconArchitect[11] to whole-genome
63  sequencing data from The Cancer Genome Atlas (TCGA), to quantify and characterize
64  the architecture of amplified regions that are larger than 10kb  and have more than 4

2

65    copies (CN>4) above median sample ploidy (**Supplementary Table 1**). Amplicons were

66    classified as 'Circular' (**Extended Data Fig. 1A**) representing amplicons residing

67    extrachromosomally or ecDNA structures that reintegrated into non-native chromosomal

68    locations as homogenously staining regions (HSRs), 'Amplified-noncircular' for linear

69    amplifications, or as 'heavily rearranged', for non-circular amplicons containing segments

70    from different chromosomes, or regions that were very far apart on chromosomes (>1Mb)

71    regions. Sample lacking amplifications were labeled 'no copy number amplification (CNA)

72    detected'.

73    To evaluate the accuracy of the computational predictions, we similarly analyzed whole

74    genome sequencing data from a panel of 34 cancer cell lines[1,2], for which tumor cells in

75    metaphase could be examined. We used 15 unique fluorescence in-situ hybridization

76    (FISH) probes in combination with matched centromeric probes (60 distinct "cell-line,

77    probe" combinations) to determine the chromosomal or extrachromosomal location of a

78    set of amplicons. We observed that 100% of amplicons characterized as 'Circular' by

79    whole genome sequencing profile demonstrated extrachromosomal fluorescent signal

80    (**Extended Data Fig. 1B**). Circular amplicons had a median count of 14.5 ecDNA per cell,

81    in contrast with the 'Amplified-noncircular' category, which had a median count of 0.0

82    ecDNA per cell. However, ecDNAs may be undercounted in amplicons with low copy

83    number**.** 'Heavily rearranged amplicons' showed at least one ecDNA per cell in two of five

84    cases, suggesting that this category consists of a mixture of chromosomal and

85    extrachromosomal amplifications. We excluded the more ambiguous category of 'heavily

86    rearranged' amplicons from futher comparisons, confining our analysis to 1,695 TCGA

87    samples. The analytic results of the 256 samples containing the more ambiguous 'heavily

88    rearranged amplicons' are presented in the supplement (**Extended Data Fig. 2**).

89        We found that 436 (26%) of the 1,695 tumor samples carried one or more Circular

90    amplicons, suggesting that ecDNA-based amplification is a common event in human

91    cancer (**Fig. 1B**). In contrast, Circular amplifications were found in <0.5% of matched

92    whole blood or normal tissue samples, suggesting that extrachromosomal amplification

93    is a mechanism that is used primarily by cancer cells (**Fig. 1B**). Of note, our analysis does

94    not reflect the presence of circulating cell free DNA in blood, or of small (<1 kb), circular,

95      non-amplified DNAs, that have been shown to be common in non-neoplastic and tumor

96      tissues [13-15]. EcDNA-based Circular amplicons were found in all cancer types except

97      acute myeloid leukemia and thryroid carcinoma, including at high frequency in many

98      cancers that are considered to be amongst the most aggressive histological types. The

99      distribution of Circular amplicon frequencies across the samples are consistent with

100     earlier results on cancer models, showing that ecDNA driven amplifications were a

101     defining feature of multiple cancer sub-types, but not normal cells[2].

102     The chromosomal distribution of the 627 Circular amplicons was highly non-

103     random (**Fig. 1C**), more so when compared to the Amplified-noncircular regions

104     (**Extended Data Fig. 3A**). We found that 41% of the 24 most recurrent amplified

105     oncogenes were most frequently present on Circular amplicons, with frequencies ranging

106     from 25% of samples for *PAX8* to 91% for *CDK4* (**Fig. 1D**). The result carried over to a

107     larger list of 707 genes that were amplified in at least five samples, with 41% of those

108     oncogenes most frequently being amplified on circular structures (**Extended Data Fig.

109     3B**). We found that oncogenes amplified on circular amplicons achieved higher copy

110     numbers than the same oncogenes amplified on Amplified-noncircular structures

111     (**Extended Data Fig. 3C**).  We further observed that the association between ecDNA

112     structures and oncogene amplification did not extend to breakpoints. For 24 frequently

113     amplified oncogenes, the frequency of observing a specific number of breakpoints in a

114     unit interval decayed exponentially, consistent with random occurrence around the

115     oncogene (**Fig. 1E; Extended Data Fig. 3D, Extended Data Fig. 3E**).   These results

116     suggest that ecDNA are formed through a random process, where selection for higher

117     copies of growth promoting driver oncogenes leads to rapid oncogene amplification

118     during cancer development and progression, in a way that also retains intratumoral

119     genetic heterogeneity, due to its mechanism of uneven inheritance[3,16].

120     Circular amplicons also differed from Amplified-noncircular amplifications in other

121     notable ways. Circular and Amplified-noncircular amplifications showed similar likelihood

122     of occurring in samples with chromosome-arm level aneuploidy (**Extended Data Fig. 4A**)

123     and whole-genome duplication,   which might arise as a result of chromosome

124     missegregation [17] or other mitotic errors [18] (**Extended Data Fig. 4B**). Smaller and more

4

125 focal genomic gains and losses result from different mutagenic processes, associated
126 with genomic instability. We observed an increase in the number of DNA segments in
127 samples marked by Circular amplicons, compared to other categories (**Fig. 2A**). The
128 frequency of copy number losses was comparable between Circular and Amplified-
129 noncircular amplicon samples (**Extended Data Fig. 4C**), but genomic segment gains
130 were more frequently detected in samples with circular amplification ( Wilcoxon rank sum
131 test: p-val < 1e-14) (**Extended Data Fig. 4D**). This observation coincided with a threefold
132 increase in gene fusion events inferred from matching RNAseq profiles (**Fig. 2B;** Binomial
133 test: p-value <1e-138) compared to Non-circular amplification. Clustered mutations, also
134 referered to as kataegis, were significantly more frequently detected in Circular amplicons
135 relative to Amplified-noncircular amplicons, suggesting increased incidence of kataegis
136 (Hypergeometric test: p-value $\cong 0$)(**Extended Data Fig. 4E**). The majority of Circular
137 amplicon breakpoints showed no or minimal sequence homology (<5 bp), raising the
138 possibility that non-homologous end joining could be involved in ecDNA formation. In
139 contrast, Amplified-noncircular amplicon breakpoints showed significantly more micro-
140 homologies than were seen on circular amplicons (**Extended Data Fig. 4F,** p-
141 value<0.0005; two-sided Fisher's exact test).

142      We sought to examine the transcriptional consequences of circular ecDNA
143 amplification at the population level. We detected a highly significant correlation between
144 DNA copy number and gene expression level in all categories of DNA amplification,
145 Circular and Non-circular. However, at comparable DNA copy number, oncogenes on
146 Circular amplicons were significantly more highly expressed than those on Amplified-
147 noncircular amplicons (p-value < 0.003; Wilcoxon rank sum test; **Fig. 3A; Extended Data**
148 **Fig. 5**), showing a higher transcriptional rate (2.6X higher compared to Amplified-
149 noncircular, 8.3X higher compared to oncogenes on non-amplified regions). To test if the
150 epigenetic mechanisms governing gene expression were different between Circular
151 amplifications and Amplified-noncircular regions, we analyzed the overlapping ATAC-seq
152 profiles available for 24 samples[19]. The results (**Fig. 3B**) showed that chromatin of
153 Circular amplicons was significantly more accessible compared to Amplified-noncircular
154 categories (1.3 times higher ATAC-seq signal; Wilcoxon rank sum test; p-value < 0.003),

5

155  suggesting that increased accessibility plays a role in dysregulation and higher
156  expression of oncogenes on circular amplifications (ecDNAs).

157  Having developed a way to stratify tumors based upon amplification architecture,
158  we examined the impact of ecDNA-based amplification on two hallmarks of cancer,
159  immune evasion and cell proliferation. We used previously developed gene expression
160  signatures[20] to evaluate  the distribution of immune infiltrate and cell proliferation scores
161  by amplicon grouping. The cellular proliferation but not immune infiltration pathway
162  scores were significantly higher (**Fig. 4A**, p-val < 1e-7; Wilcoxon Rank Sum Test;
163  **Extended Data Fig. 6**) in the Circular amplification category compared to the other two
164  groups. We did not observe difference in activity of the immune signature score
165  between groups (**Fig. 4A**, p-val < 0.03; Wilcoxon Rank Sum Test).  The increased
166  activity of the cell proliferation gene signature suggested a higher rate of proliferation
167  and tumors that behave more aggressively.

168  To determine whether cancers that have ecDNA amplification were associated
169  with tumor progression, we examined the impact of circular amplification on lymph node
170  status at initial presentation, and overall survival. We found that the proportion of cases
171  in which the tumor had spread to a lymph node at the initial time of diagnosis was
172  significantly increased in tumor samples that had either circular or non-circular
173  amplification (**Fig. 4B**; p-value < 0.02 no- CNA vs Amplified-noncircular, p-value < 1.0e-
174  05 Circular-amplicon vs Amplified-noncircular). Additionally, we found a significant
175  difference in overall survival of patients stratified by amplification category. Patients
176  whose tumors contained circular amplification associated with significantly worse overall
177  outcomes compared to patients whose tumors harbored either non-circular amplifications
178  or no amplifications (**Fig. 4C**; p-val < 1e-15 versus no-CNA detected; p-val < 0.07 against
179  Amplified-noncircular; Log-rank test). To account for the possibility that differences in
180  survival rate are being influenced by the disease subtype, as circular amplicons are much
181  more prevalent in aggressive cancers such as glioblastoma, we fit the data to a Cox
182  Hazard model that tested survival after controlling for disease subtype.  The model
183  showed that patients with circular amplicons had significantly higher hazard rates (**Fig.**
184  **4D**; 28% increase in hazard rate relative to no-CNA, p-val < 0.03).

6

185    Cancer genomics is itself evolving from reading out the "code" to unraveling its function.

186    The 3D organization of the genome plays a critical role in determining how that genome

187    functions, or malfunctions, as occurs in cancer. The data presented here demonstrate

188    that ecDNA play a critical role in cancer, providing a mechanism for achieving and

189    maintaining high copy oncogene amplification and diversity. This mechanism of

190    amplification is operant in a large fraction of human cancers, and contributes to the poor

191    outcomes for patients. The potential to leverage the presence of ecDNAs in a quarter of

192    human cancers for diagnostics or therapeutics provides a link between cancer genomics

193    and broad utility for patient populations.

194

195    METHODS

196    **AmpliconArchitect**

197    We used AmpliconArchitect[11] infer the architecture of the `amplicons' --- large (>10kb)

198    rearrangements with high copy numbers (CN>4) that are inferred to have co-amplified

199    as a structure. AmpliconArchitect takes as an input aligned WGS sequences and seed

200    intervals of the amplicon.  AmpliconArchitect then searches for other regions that belong

201    to the amplicon by exploring the seed intervals, and extends beyond the intervals if it

202    encounters copy number changes or discordant edges that support a breakpoint. The

203    collection of intervals and breakpoints are combined to form a fine network with nodes

204    representing segments and edges representing rearrangements, which we call the

205    breakpoint graph. This breakpoint graph is can be further decomposed into simple

206    cycles to identify any circular paths within the amplicon structure, which is indicative of

207    ecDNA presence. The detected amplicons were annoted with the Ensembl Release 75

208    gene database (GRCh37).

209

210    **Amplicon and sample classification**

211    As a perquisite, amplicons must contain $\geq$ 10kb of genomic segments amplified to at

212    least four copies above median ploidy in order to be considered a valid amplicon. We

213    then use the AmpliconArchitect derived breakpoint graph to classify amplicons into

214    three categories: 1. Circular amplification; 2. Heavily rearranged amplification; and, 3.

    7

215    Amplified-noncircular (**Extended Data Fig. 1A)**. Amplicons were denoted as Circular

216    amplification if the segments form a cycle in the graph of total size at least 10kb and has

217    at least a copy count of four. Non-circular amplicons were denoted as heavily

218    rearranged if the breakpoints connect segments from different chromosomes, or distal

219    (>1Mb) regions (**Extended Data Fig. 1A)**. Non-circular, non-distal amplicons were

220    denoted as locally rearranged. All other regions that were not part of any amplicon

221    structure were classified as not-amplified. While an amplicon may fit the requirements

222    for several categories (i.e., a circular amplicon may also comprise heavily rearranged

223    amplifications), priority was given to the circular amplification category, followed by

224    heavily rearranged and finally amplified-noncircular. Similarly, samples were classified

225    based upon what amplicons are present within the sample, giving precedence to the

226    presence of amplicon with highest priority. For example, a sample with both circular

227    and heavily rearranged amplification would be classified as circular amplification.

228    Samples without any amplicons are classified as `No CNA detected'.

229    **Cell line validation**

230    We ran AmpliconArchitect on the whole-genome sequencing data from 11 glioblastoma

231    neurosphere cultures from deCarvalho et al [1]. with FISH images and the 23 cell lines

232    from different cancer types from the Turner et al [2] study. Fluorescence in-situ

233    hybridization (FISH) was performed in parallel, as described. The seed interval for each

234    cell line included the probe region. For each probe, we reported whether it landed in an

235    amplicon (inferred from AmpliconArchitect), and if so what was the amplicon

236    classification. The distribution of the average ecDNA per cell was computed as the

237    average number FISH probes that co-localized on ecDNA across all the images for that

238    particular cell line+FISH probe combination (**Extended Data Fig. 1B**). Wilcoxon Rank

239    Sum test was used to detect significant differences in average ecDNA counts per cell

240    across the amplicon classes.

241    **TCGA processing**

242    We processed TCGA whole genome sequencing BAMs through the Institute for

243    Systems Biology Cancer Genomics Cloud (https://isb-cgc.appspot.com/) that provides a

244    cloud-based platform for TCGA data analysis. We used genome-wide snp6 copy

     8

245    number segments with copy number log ratio equal to 1 as seed interval(s) of interest

246    that are required for the input to AmpliconArchitect[11]. Default parameters and reference

247    files were used for all other settings. Details on how to run AmpliconArchitect have been

248    described in the corresponding manuscript[11] and its source code depository.

249    We ran AmpliconArchitect on tumor and normal WGS samples from 1979 patients.

250    Samples were classified based upon the amplicon with highest precedence present in

251    the sample, or classified as `No CNA detected' if no amplicons are present in the

252    sample.  Samples classified as highly-rearranged are removed from further analysis due

253    to the ambiguity of ecDNA status of the sample.

**TCGA processed data**

255    The processed data (hg19) and clinical data were found at the GDC

256    (https://portal.gdc.cancer.gov/ legacy-archive/search/f) and the PancanAtlas publication

257    page (https://gdc.cancer.gov/about-data/publications/pancanatlas). Somatic muations

258    for TCGA whole genome sequencing were downloaded from the ICGC PCAWG portal

259    (https://dcc.icgc.org/pcawg)[21].

**Oncogene analysis**

261    We examined the enrichment of the 24 recurrent oncogenes known to be activated by

262    amplification by counting the total number of base pairs from the amplicon classes from

263    all the tumor samples that overlap these oncogenes.  We then simulated 10,000

264    replicates by sampling random regions of the same size of the amplicons and computed

265    an empirical expected distribution of base pairs covering the oncogenes if the amplicons

266    were randomly sampled across the genome.  We report the z-score between the

267    empirical distribution and observed value for the amplicon classes.  We also report the

268    average copy count, estimated from AmpliconArchitect.  For each of these oncogenes

269    on an amplicon structure, we reported the position of breakpoint detected within a 1 MB

270    region flanking the oncogene using the breakpoint graph to infer breakpoints.  We

271    partitioned the region into 1000 bp windows and counted the total number of

272    breakpoints that landed in each window, and display a histogram of these counts.  We

273    modeled the histograms using an exponential distribution and show that under the

9

274   assumption that the breakpoints are distributed randomly, the histograms closely follow

275   the exponential distribution.

276   We used *allOnco* (http://www.bushmanlab.org/links/genelists), a set of 2,579 cancer

277   genes generated from curated collections cancer genes from many different

278   publications.  We identified all amplicons that overlapped with the oncogenes and report

279   the proportion amplified oncogenes that are circular.

## Genomic instability analyses

281   We computed total copy number gains/losses as the number of snp6 copy number

282   segments with copy number >=4 or <= 1.  Wilcoxon Rank Sum test was used to test for

283   a significant difference between the two distributions.  We used the data from a previous

284   study[22] to estimate the genome doubling status and chromosomal arm duplication and

285   loss for each sample.  Wilcoxon Rank Sum test was used to test significance between

286   the distribution of gains and losses and Chi-squared test was used to test significance

287   between the distribution of whole genome doublings.

288   We used the data from the TCGA fusion database (https://tumorfusions.org/)[23] to

289   identify fusions events that occur on an amplicon.  For each fusion in the database, we

290   consider it valid if both ends of the fusion breakpoint junction occur on the same

291   amplicon. In total, 710 amplified fusions were detected.  We computed the average

292   fusion events per 10 Mb as the total number of fusions that landed within an amplicon

293   class divided by the sum of all the base pairs of the amplicon class multiplied by 10e7.

294   To test whether circular amplicons were enriched fusion events, we computed the p-

295   value of observing at least the number of fusion events on circular amplicon under a

296   binomial distribution where the probability *p* was estimated using the total number of

297   fusion events on the amplified-noncircular  divided by the total base pairs of the

298   amplified non-circular event and the number of trials *n* as the total base pairs of the

299   circular amplicons.

## RNAseq and ATACseq analyses

301   Of the 1,695 tumor samples, 1,608 had RNA-seq data in the format of FPKM-UQ

302   expression data.  For each gene within each disease cohort, we computed a baseline

303   FPKM-UQ as the average FPKM-UQ of all samples for which the gene was not found

10

304    on an amplicon (i.e., average expression of the unamplified gene).  We then computed

305    the fold-change in expression of each gene on each amplicon as the FPKM-UQ of the

306    amplified gene divided by the average FPKM-UQ of the unamplified samples, and

307    report the distribution of fold-changes versus the copy number.  Tukey's range test was

308    used to test significance between slope of the FPKMs for circular and amplified-

309    noncircular.

310    ATAC-seq profiles were available for 24 samples.  For each amplicon in each sample,

311    fold-change in ATAC-seq signal was computed as the average ATAC-seq signal across

312    the amplicon region divided by the average ATAC-seq signal for the same region in the

313    unamplified samples of the same cancer type.  Wilcoxon Rank Sum test was used to

314    test significance between the two distributions.

315    **Kataegis**

316    Localized mutation clusters (kataegis loci) were defined as having 6 or more

317    consecutive mutations with an inter-mutation distance of < 1kb in a similar way to a

318    previously used approach [24].

319    **Inferring breakpoint homologies**

320    For each breakpoint, sequencing reads around +/- 1000 bps of the breakpoint were

321    locally reassembled with SvABA[25] to produce a contiguous consensus sequence of

322    each breakpoint, precise breakpoint positions, and the level of homology at breakpoints.

323    **Statistical analysis**

324    Survival curves were estimated with the Kaplan–Meier method, and comparison of

325    survival curves between groups was performed with the log-rank test in R survival

326    package. Hazard ratios were estimated with the Cox proportional hazards regression

327    model in the survival R package.

328    **Competing Interests**

329    P.S.M., H.Y.C., V.B. and R.G.W.V. are co-founders of Boundless Bio, Inc. (BB), and

330    serve as consultants. V.B. is a co-founder, and has equity interest in Digital Proteomics,

331    LLC (DP), and receives income from DP. The terms of this arrangement have been

332    reviewed and approved by the University of California, San Diego in accordance with its

11

333 conflict of interest policies. BB and DP were not involved in the research presented

334 here.

335 **ACKNOWLEDGEMENTS**

350

351 **REFERENCES**

352 1  deCarvalho, A. C. *et al.* Discordant inheritance of chromosomal and extrachromosomal DNA
353   elements contributes to dynamic disease evolution in glioblastoma. *Nat Genet* **50**, 708-717,
354   doi:10.1038/s41588-018-0105-0 (2018).
355 2  Turner, K. M. *et al.* Extrachromosomal oncogene amplification drives tumor evolution and
356   genetic heterogeneity. *Nature* **543**, 122-125, doi:10.1038/nature21356 (2017).
357 3  Verhaak, R. G. W., Bafna, V. & Mischel, P. S. Extrachromosomal oncogene amplification in tumor
358   pathogenesis and evolution. *Nat Rev Cancer*, doi:10.1038/s41568-019-0128-6 (2019).
359 4  Weischenfeldt, J. *et al.* Pan-cancer analysis of somatic copy-number alterations implicates IRS4
360   and IGF2 in enhancer hijacking. *Nat Genet* **49**, 65-74, doi:10.1038/ng.3722 (2017).
361 5  Zack, T. I. *et al.* Pan-cancer patterns of somatic copy number alteration. *Nat Genet* **45**, 1134-
362   1140, doi:10.1038/ng.2760 (2013).
363 6  Beroukhim, R. *et al.* The landscape of somatic copy-number alteration across human cancers.
364   *Nature* **463**, 899-905, doi:10.1038/nature08822 (2010).
365 7  Alt, F. W., Kellems, R. E., Bertino, J. R. & Schimke, R. T. Selective multiplication of dihydrofolate
366   reductase genes in methotrexate-resistant variants of cultured murine cells. *J Biol Chem* **253**,
367   1357-1370 (1978).

368  8      Kohl, N. E. *et al.* Transposition and amplification of oncogene-related sequences in human
369          neuroblastomas. *Cell* **35**, 359-367 (1983).
370  9      Mitelman, F., Johansson, B. & Mertens, F. The impact of translocations and gene fusions on
371          cancer causation. *Nat Rev Cancer* **7**, 233-245, doi:10.1038/nrc2091 (2007).
372  10     Trask, B. J. Fluorescence in situ hybridization: applications in cytogenetics and gene mapping.
373          *Trends Genet* **7**, 149-154, doi:10.1016/0168-9525(91)90378-4 (1991).
374  11     Deshpande, V. *et al.* Exploring the landscape of focal amplifications in cancer using
375          AmpliconArchitect. *Nat Commun* **10**, 392, doi:10.1038/s41467-018-08200-y (2019).
376  12     Xu, K. *et al.* Structure and evolution of double minutes in diagnosis and relapse brain tumors.
377          *Acta Neuropathol*, doi:10.1007/s00401-018-1912-1 (2018).
378  13     Moller, H. D. *et al.* Circular DNA elements of chromosomal origin are common in healthy human
379          somatic tissue. *Nat Commun* **9**, 1069, doi:10.1038/s41467-018-03369-8 (2018).
380  14     Kumar, P. *et al.* Normal and Cancerous Tissues Release Extrachromosomal Circular DNA
381          (eccDNA) into the Circulation. *Mol Cancer Res* **15**, 1197-1205, doi:10.1158/1541-7786.MCR-17-
382          0095 (2017).
383  15     Shibata, Y. *et al.* Extrachromosomal microDNAs and chromosomal microdeletions in normal
384          tissues. *Science* **336**, 82-86, doi:10.1126/science.1213307 (2012).
385  16     Turner, K. M. *et al.* Extrachromosomal oncogene amplification drives tumor evolution and
386          genetic heterogeneity. *Nature*, doi:10.1038/nature21356 (2017).
387  17     Davoli, T. & de Lange, T. The causes and consequences of polyploidy in normal development and
388          cancer. *Annu Rev Cell Dev Biol* **27**, 585-610, doi:10.1146/annurev-cellbio-092910-154234 (2011).
389  18     Bielski, C. M. *et al.* Genome doubling shapes the evolution and prognosis of advanced cancers.
390          *Nat Genet* **50**, 1189-1195, doi:10.1038/s41588-018-0165-1 (2018).
391  19     Corces, M. R. *et al.* The chromatin accessibility landscape of primary human cancers. *Science*
392          **362**, doi:10.1126/science.aav1898 (2018).
393  20     Davoli, T., Uno, H., Wooten, E. C. & Elledge, S. J. Tumor aneuploidy correlates with markers of
394          immune evasion and with reduced response to immunotherapy. *Science* **355**,
395          doi:10.1126/science.aaf8399 (2017).
396  21     Rheinbay, E. *et al.* Discovery and characterization of coding and non-coding driver mutations in
397          more than 2,500 whole cancer genomes. *bioRxiv*, 237313, doi:10.1101/237313 (2017).
398  22     Taylor, A. M. *et al.* Genomic and Functional Approaches to Understanding Cancer Aneuploidy.
399          *Cancer Cell* **33**, 676-689 e673, doi:10.1016/j.ccell.2018.03.007 (2018).
400  23     Hu, X. *et al.* TumorFusions: an integrative resource for cancer-associated transcript fusions.
401          *Nucleic Acids Res* **46**, D1144-D1149, doi:10.1093/nar/gkx1018 (2018).
402  24     Alexandrov, L. B., Nik-Zainal, S., Wedge, D. C., Campbell, P. J. & Stratton, M. R. Deciphering
403          signatures of mutational processes operative in human cancer. *Cell Rep* **3**, 246-259,
404          doi:10.1016/j.celrep.2012.12.008 (2013).
405  25     Wala, J. A. *et al.* SvABA: genome-wide detection of structural variants and indels by local
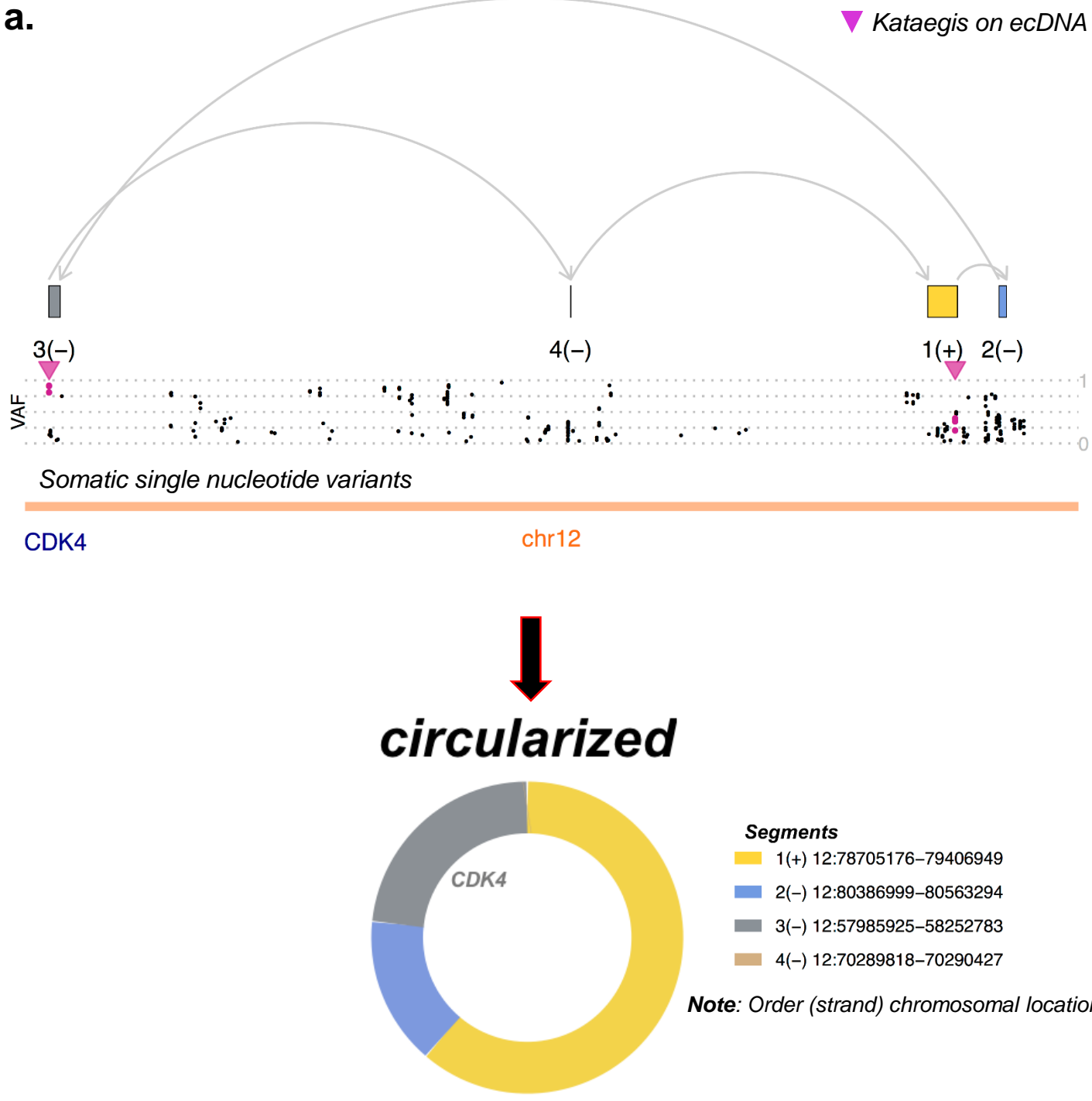406          assembly. *Genome Res* **28**, 581-591, doi:10.1101/gr.221028.117 (2018).
407
408

**a.**

▼ *Kataegis on ecDNA*

3(−)    4(−)    1(+)  2(−)

VAF

*Somatic single nucleotide variants*

CDK4    chr12

*circularized*

CDK4

**Segments**
- ■ 1(+) 12:78705176–79406949
- ■ 2(−) 12:80386999–80563294
- ■ 3(−) 12:57985925–58252783
- ■ 4(−) 12:70289818–70290427

***Note***: *Order (strand) chromosomal location*

**Fig. 1 | Frequency of circular amplification across tumor and non-tumor tissues. A**. Representative example of a Circular DNA amplification.
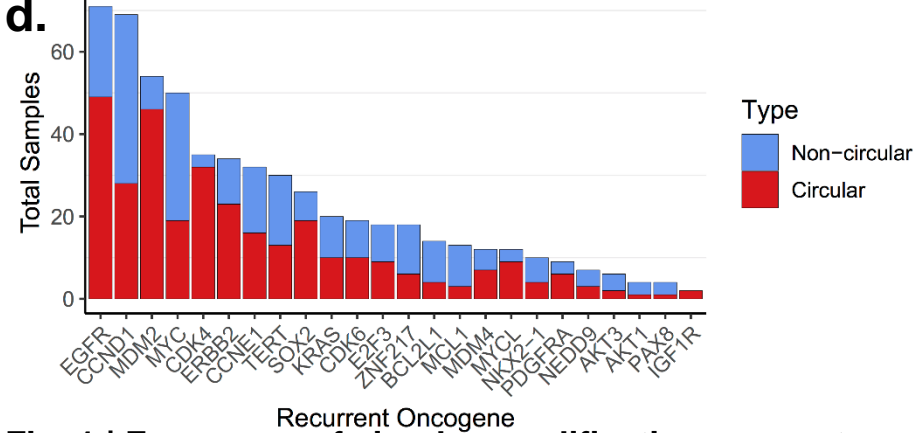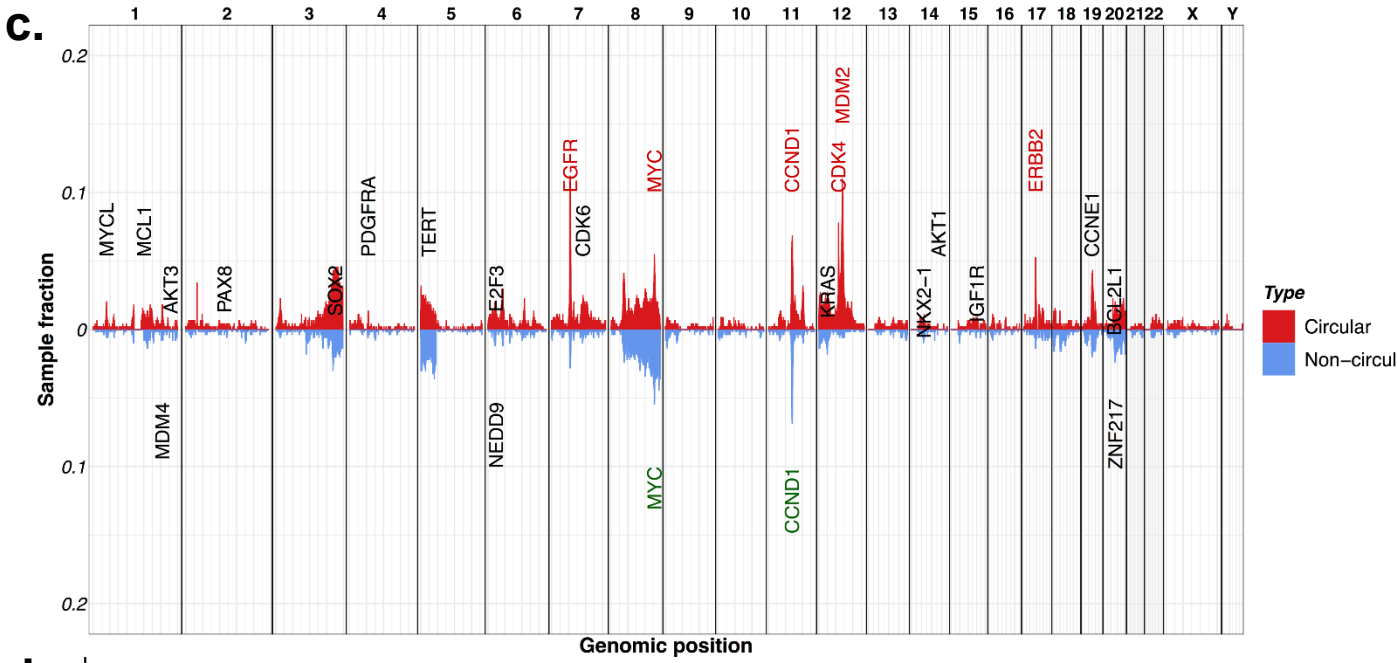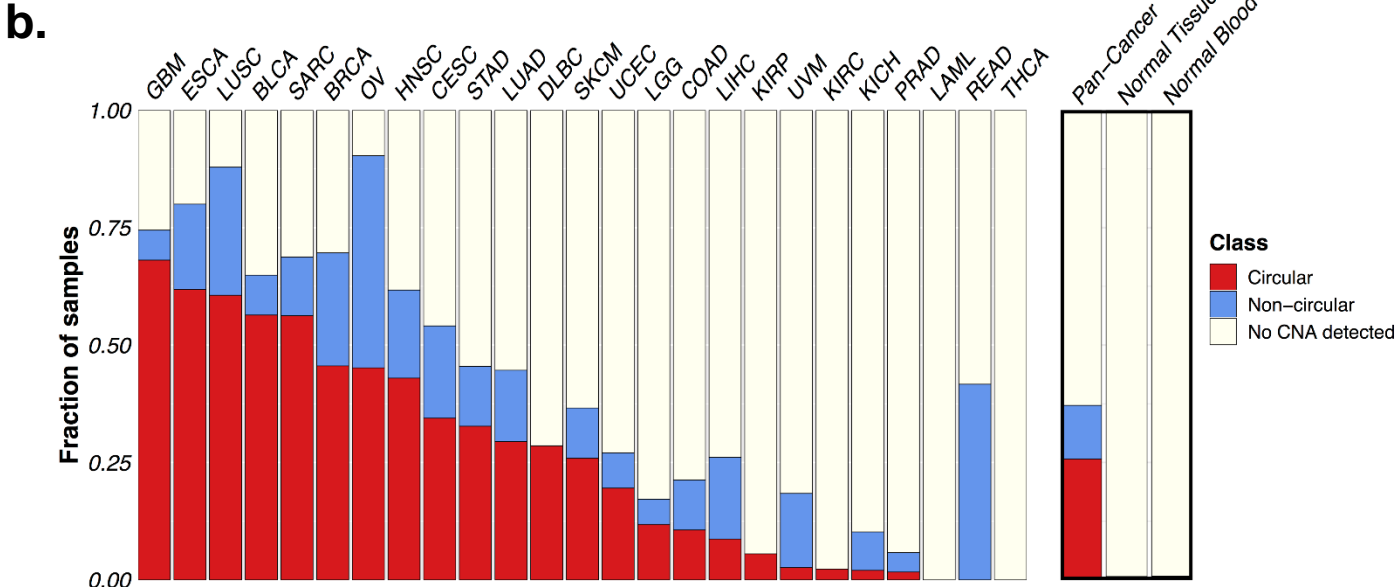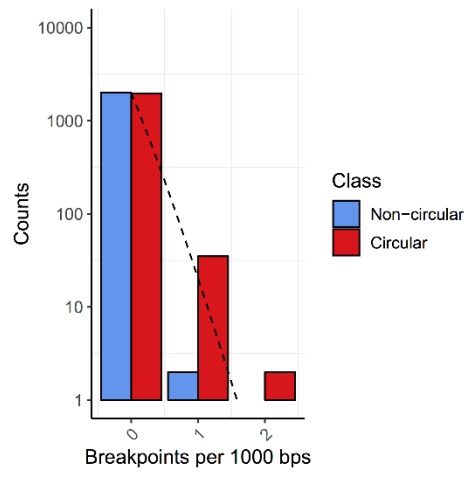
**Fig. 1 | Frequency of circular amplification across tumor and non-tumor tissues. B.** Distribution of circular, non-circular, and no copy number alteration (no CNA) detected categories by tumor and normal tissue. **C**. Genome-wide distribution of circular (red) and non-circular (blue) amplification peaks. **D**. Classification of circular vs non-circular amplification status by gene. Shown are 24 most frequently amplified oncogenes.

**Fig. 1 | Frequency of circular amplification across tumor and non-tumor tissues. E**. Breakpoint locations and distribution of breakpoints across all samples with amplified EGFR (top), CDK4 (middle), and MYC (bottom).

**Fig. 2 | Total number of copy number segments and transcript fusions are increased in Circular amplicon tumor samples**. **A.** TCGA copy number array data was used to count the total number of DNA segments within a sample. Circular samples contained statistically significantly more DNA segments than non-circular and no CNA detected (p-val < 1e-5 and 1e-128, respectively; Wilcox Rank Sum Test). **B.** Circular structures expressed significantly more gene fusions compared to non-circular amplicons, after size normalization.

**Fig. 3 | Gene expression and chromatin accessibility of amplicon classes. A.** Copy number of the gene versus its fold-change in FPKM for all genes with a copy count greater than 4 and less than 100, for each gene on each amplicon. The fold-change in FPKM is computed as the gene's (FPKM-UQ+1) divided by the average of (FPKM-UQ+1) for the same gene in all other tumor samples from the same cohort for which the gene is not on any amplicon (i.e., not amplified). Linear regression lines are shown for each classification class. Tukey's range test shows genes on circular structures are significantly different to genes on non-circular structures (p-value < 1e-15). **B.** For each amplicon in the 24 TCGA samples with ATAC-seq and AmpliconArchitect results, the log2 fold-change in ATAC-seq signal across the amplicon relative to tissue types without amplification within the same region is shown. Each point represents a separate amplicon. The distribution of fold-change for circular amplicons is statistically significantly higher than non-circular (Wilcoxon rank sum test; p-value < 0.003).
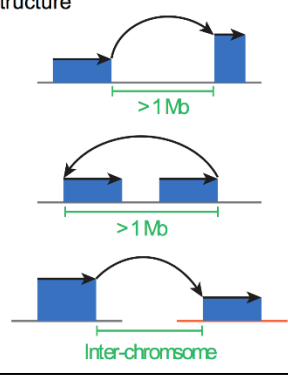
**Fig. 4 | Presence of circular amplification associates with poor outcomes**. **A.** Cell proliferation gene expression signature single sample GSEA (ssGSEA) scores by amplification category. Shown are means and 95% confidence intervals of the ssGSEA scores. Samples with circular structures showed significantly higher ssGSEA scores than samples with non-circular amplicons. **B.** Lymph node stage for primary tumors showing samples with amplification are more likely to have spread to the lymph node at time of diagnosis. **C.** Kaplan-Meier five-year survival curves by amplification category. Both amplified-noncircular and circular amplification have significantly worst outcome than No CNA samples (p-val < 1e-4 and 1e-15, respectively). Circular amplification has worse but not significant outcome compared to amplified-noncircular (p-val < 0.07). **D.** Cox-Hazard model, incorporating disease and patient cohorts as parameters showing circular amplification results in significantly higher hazard rates.

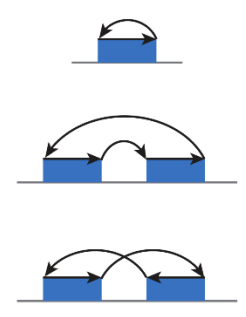**Extended Data Fig. 1 | Amplicon classification. A**. Schematic representation of the three classification categories. Amplicons are classified using a hierarchical scheme based upon the genomic reconstruction of the amplified regions (i.e., any region with a copy count of 4 or greater) and the presence or absence of discordant edges between these regions. Amplicons must have a minimum 10kb of amplified regions in order to be considered a valid amplicon. The first category is circular amplicon, which is an amplicon that contains one or more amplified segments forming a cyclic path of at least 10kb bps in length and has an average amplification of four copies. The second category is heavily-rearranged amplicon, which is an amplicon that contains amplified segments that are connected by discordant read pairs, and at least one breakpoint junctions is inter-chromosomal or greater than 1Mb is size. The third category is non-circular amplicon, which is any amplicon that contains amplified segments with no discordant edges or with discordant edges, but all breakpoint junctions are less than 1 Mb in size. All other regions are considered not amplified. As the classification scheme is hierarchical, each amplicon can only have one class, and the highest rank class has precedent (i.e., an amplicon that is both circular and heavily-rearranged will be classified only as circular). As samples can have multiple amplicons, the sample is classified as the amplicon with highest precedent (i.e., a sample with 1 circular amplicon and 3 heavily-rearranged amplicons would be classified as circular).

**b.**

**Extended Data Fig. 1 | Amplicon classification. B**. Validation on cell line data. Validation of the classification scheme on cell line data with FISH experiments for detecting ecDNA from the Turner et al. and deCarvalho et al. studies. FISH probes were designed for selected oncogenes and DAPI staining was performed to determine whether the FISH probe landed on chromosomal DNA or ecDNA. For each cell (represented as an image of the cell in metaphase), the number of positive ecDNA probes were counted, and for each cell line, the average positive ecDNA per cell was reported. For each probe, we report whether it landed in an amplicon (inferred from AA), and if so, what was the amplicon's classification. The distribution for the average ecDNA per cell between the circular and heavily rearranged (p-value < 0.05; Wilcoxon rank sum test) and No CNA detected/heavily rearranged were statistically significantly different (p-value < 0.001; Wilcoxon rank sum test).

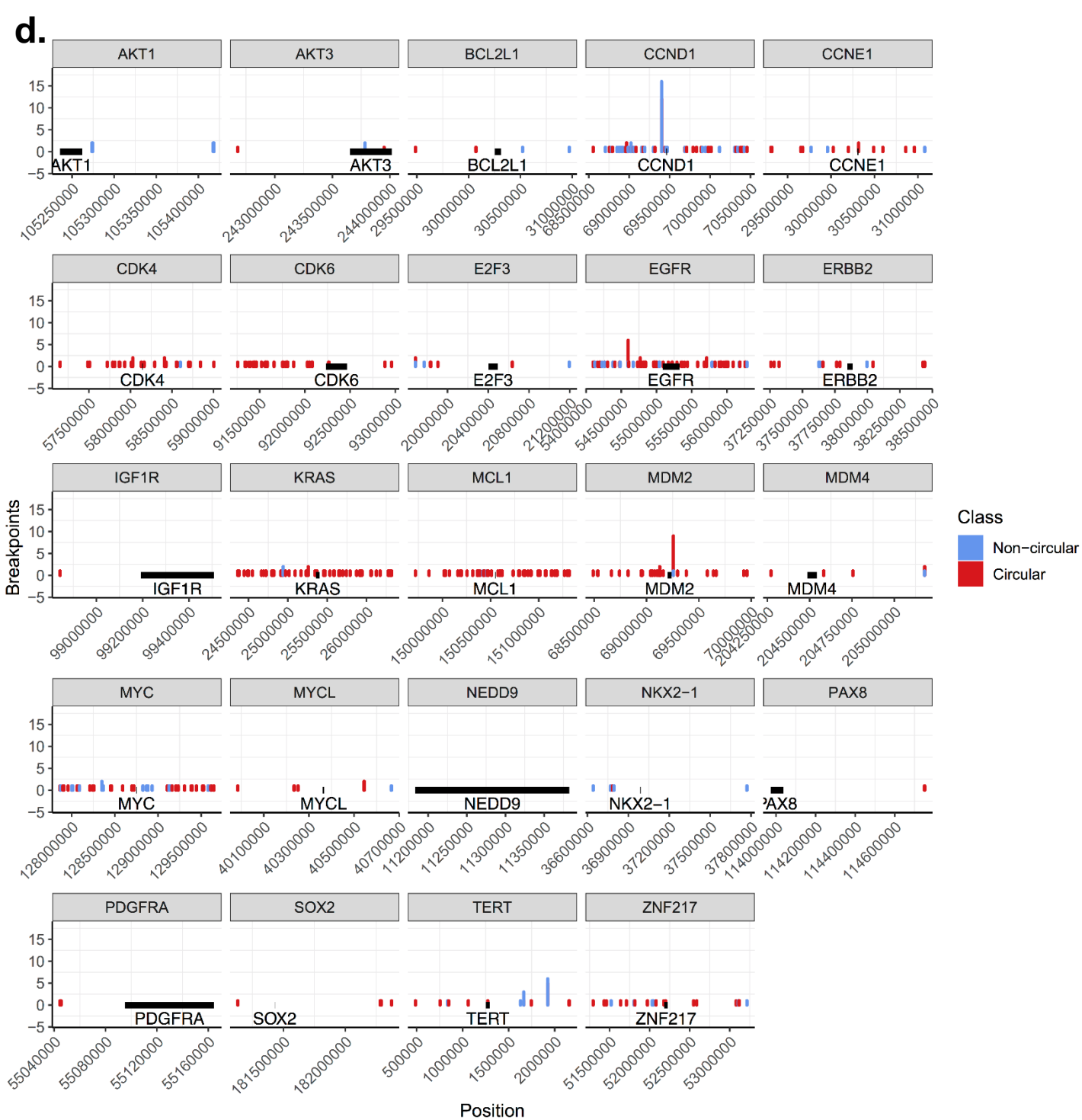**Extended Data Fig. 2 | A.** summary result including heavily-rearranged category. **B.** Genome doubling events by amplification class. **C.** Total number of copy number segments by amplification class. **D.** Kataegis frequency differences between amplification categories. **E.** Breakpoint homology by amplification class.

**Extended Data Fig. 3. Circular vs amplified non-circular amplification comparisons. A**. 24 recurrently amplified oncogenes significantly overlap circular regions (z-score 10.9), especially compared to amplified non-circular (z-score 4.0). **B**. For all oncogenes with copy number >= 4 (defined from the DNA copy number array data) and present in at least 5 samples, we show the class distribution of that oncogene. The oncogenes are ordered by proportion on circular amplification. **C**. For the 24 recurrent oncogenes known to be activated via amplification (**Zack et al. Nat Gen. 2013**), we report the average copy number for the oncogenes for circular amplification versus amplified-noncircular amplification.
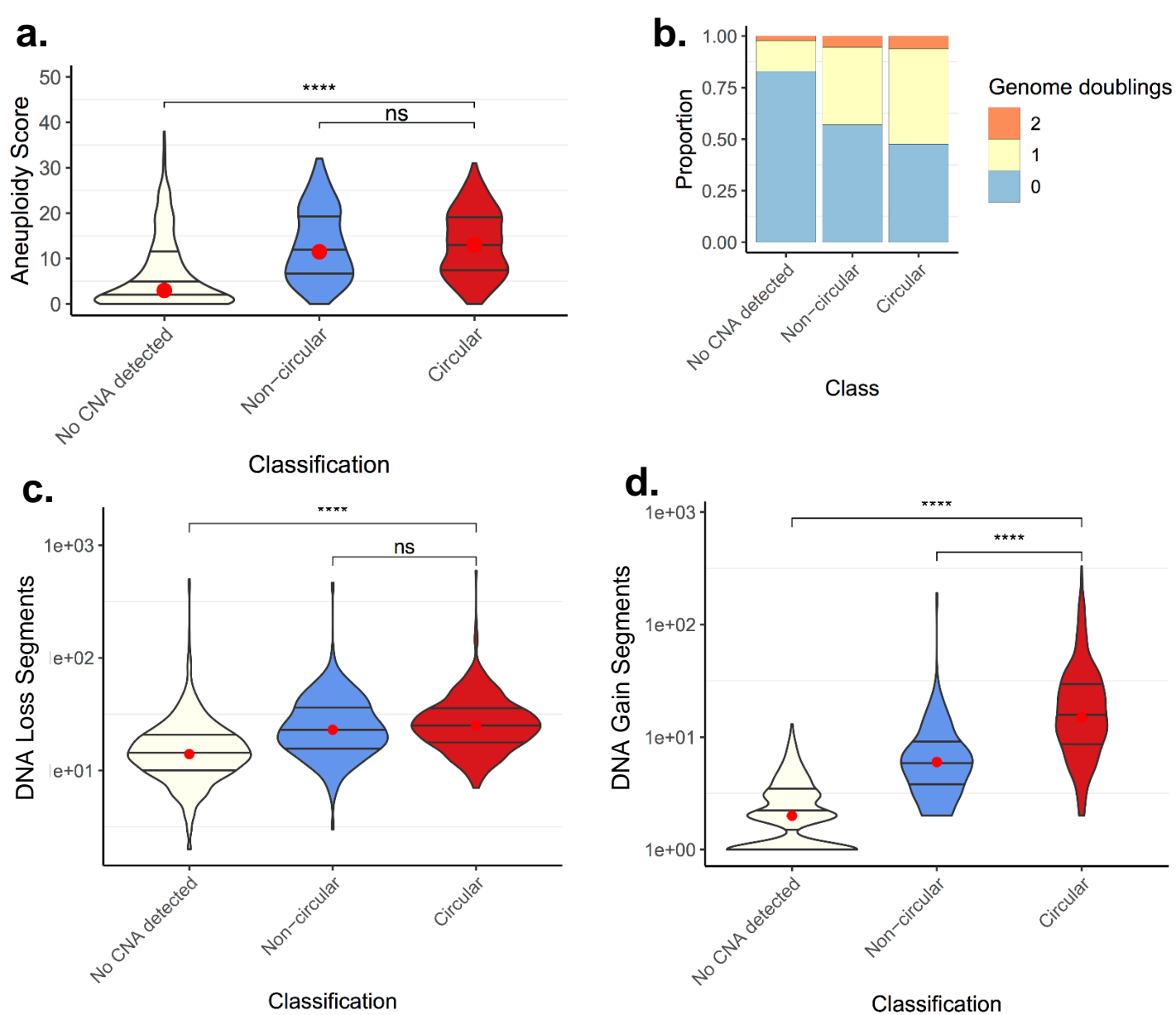
**Extended Data Fig. 3. Circular vs amplified non-circular amplification comparisons. D.** Breakpoint locations across the 24 recurrent oncogenes activated by amplification. Outliers in CCND1 and MDM2 were results of mapping bias due to short ALU repeats near the oncogene region.
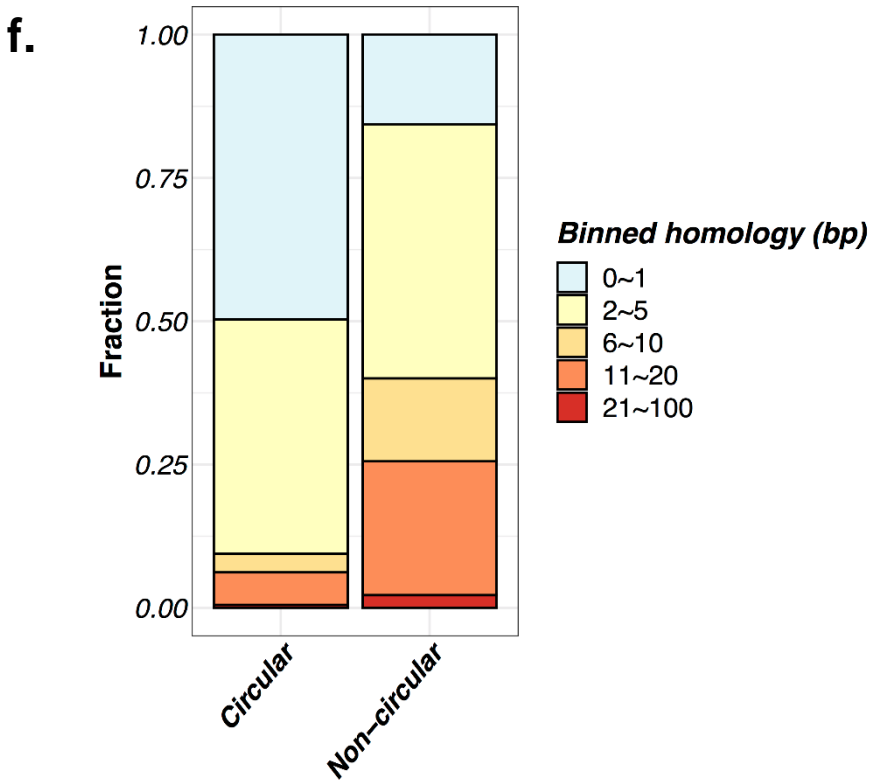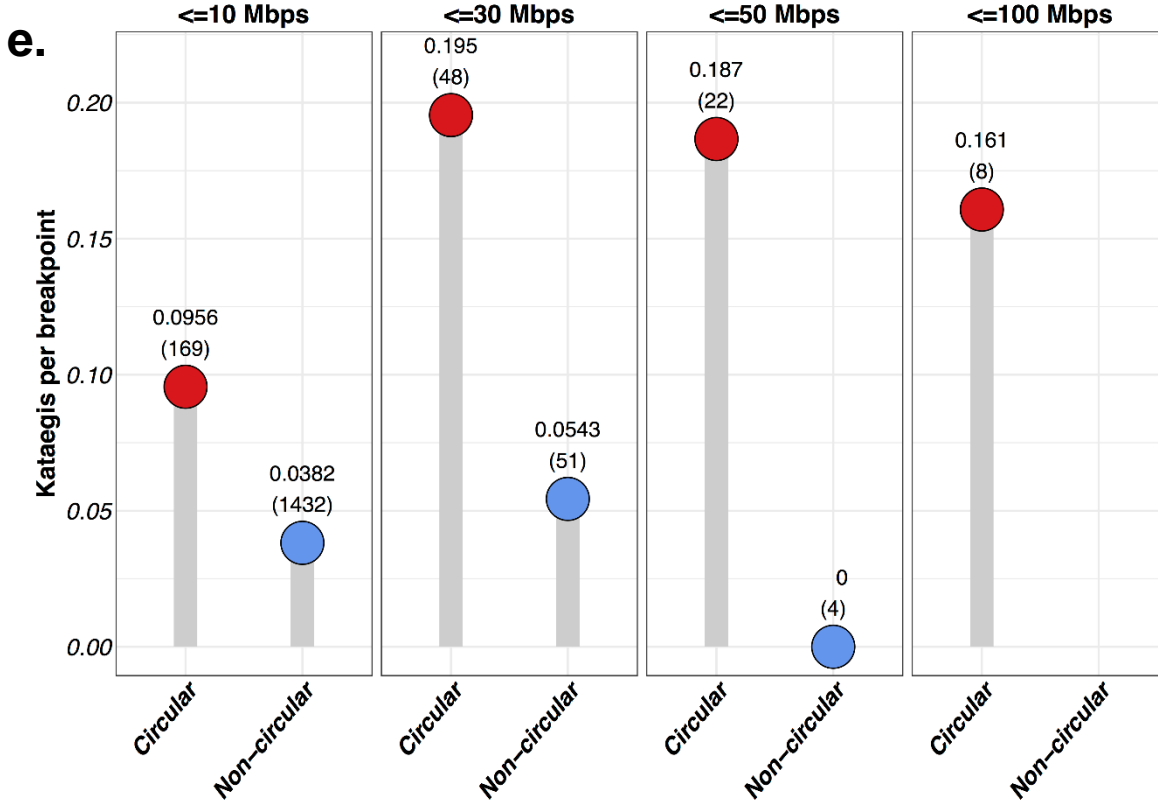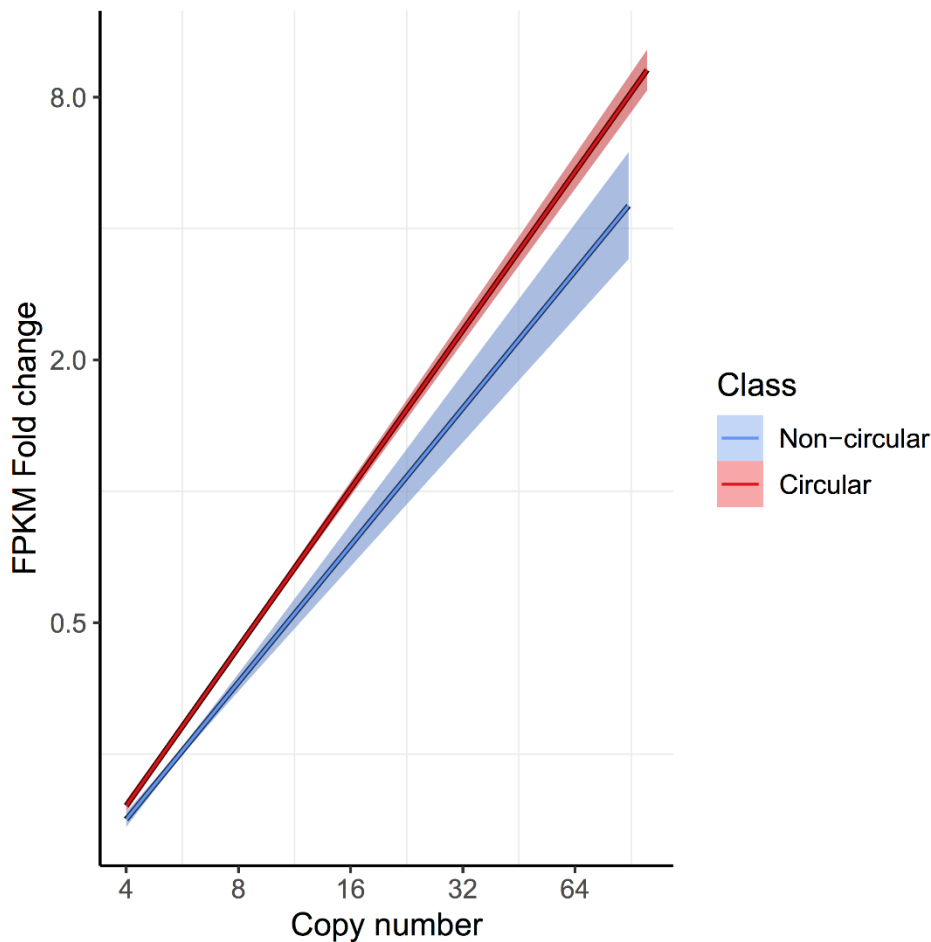
**e.**



**Extended Data Fig. 3. Circular vs amplified non-circular amplification comparisons. E.** Breakpoint locations across the 24 recurrent oncogenes activated by amplification.
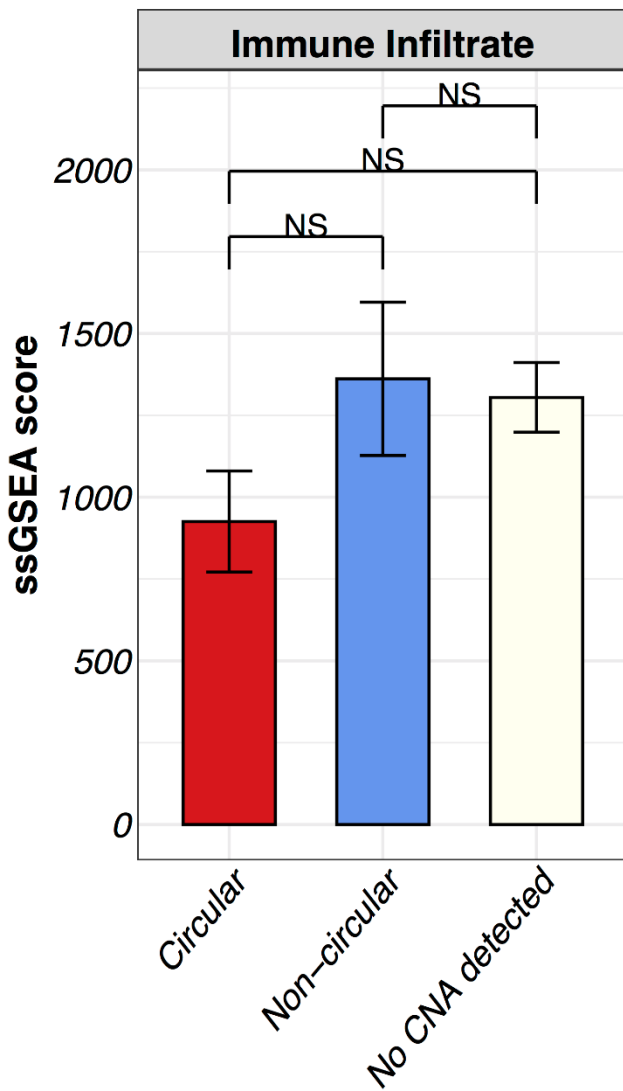
**Extended Data Fig. 4 | Aneuploidy and genomic instability events by amplification class. A.** Chromosome arm aneuploidy scores showing no difference in chromosomal arm level events between amplified-noncircular and circular amplification. **B.** Genome doublings distribution across classes showing no difference in distribution between amplified-noncircular and circular amplification. Circular amplification and non-amplified are different (Chi-square test; p-val < 1e-12). Circular amplification and amplified-noncircular are not different (Chi-square test; p-val < 0.10). **C.** Distribution for total DNA loss segments by amplification class. TCGA CNV array data was used to count the total number of DNA losses within a sample. A DNA loss was defined as a segment with CN <= 1. **D.** Same as C, but for gain segments (CN >=4). Circular samples contain statistically significantly more DNA gains than non-circular and no-CNA detected (p-val < 1e-14 and 1e-127, respectively; Wilcox Rank Sum Test). Non-circular contain statistically significantly more DNA gains than no-CNA detected (p-val < 1e-35).

**Extended Data Fig. 4 | Aneuploidy and genomic instability events by amplification class. E.** Kataegis frequency differences between amplification categories. Amplicons were grouped into Amplicon-size bins, and # kataegis frequency was normalized for the number of DNA breakpoints, demonstrating a higher occurrence of kataegis in Circular compared to Non-circular amplifications. The number of amplicons used is shown in parentheses. **F.** Breakpoint homology by amplification class. Note that inserted sequences were excluded.

**Extended Data Fig. 5 | FPKM fold-change versus copy number.** For each gene on each amplicon, we report the copy number of the gene versus its fold-change in FPKM for all genes with a copy count greater than 4 and less than 100. The fold-change in FPKM is computed as the gene's (FPKM-UQ+1) divided by the average of (FPKM-UQ+1) for the same gene in all other tumor samples from the same cohort for which the gene is not on any amplicon (i.e., is not amplified). Linear regression lines are shown for each classification class. Tukey's range test shows genes on circular structures are significantly different to genes on non-circular structures (p-value < 1e-15).

**Extended Data Fig. 6 | Circular amplification associates with worse outcomes.**
Immune gene expression signature single sample GSEA (ssGSEA) scores by
amplification category. Shown are means and 95% confidence intervals of the
ssGSEA scores. No significant difference was observed between classes.