

1 The landscape of selection in 551 Esophageal Adenocarcinomas defines 2 genomic biomarkers for the clinic

3
4 Frankell AM¹, Jammula S², Li X¹, Contino G¹, Killcoyne S^{1,3}, Abbas S¹, Perner J², Bower L²,
5 Devonshire G², Ococks E¹, Grehan N¹, Mok J¹, O'Donovan M⁴, MacRae S¹, Eldridge M², Tavare S²,
6 Fitzgerald RC¹ and the Oesophageal Cancer Clinical and Molecular Stratification (OCCAMS)
7 Consortium⁵

8

9 ¹ MRC cancer unit, Hutchison/MRC research centre, University of Cambridge, Cambridge, UK

10 ² CRUK Cambridge institute, University of Cambridge, Cambridge, UK.

11 ³ European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Hinxton, UK

12 ⁴ Department of Histopathology, Cambridge University Hospital NHS Trust, Cambridge, UK

13 ⁵ A full list of contributors from the OCCAMS Consortium is available at the end of the manuscript

14 **Abstract:**

15 Esophageal Adenocarcinoma (EAC) is a poor prognosis cancer type with rapidly rising incidence. Our
16 understanding of genetic events which drive EAC development is limited and there are few molecular
17 biomarkers for prognostication or therapeutics. We have accumulated a cohort of 551 genomically
18 characterised EACs (73% WGS and 27% WES) with clinical annotation and matched RNA-seq. Using a
19 variety of driver gene detection methods, we discover 77 EAC driver genes (73% novel) and 21 non-
20 coding driver elements (95% novel), and describe mutation and CNV types with specific functional
21 impact. We identify a mean of 4.4 driver events per case derived from both copy number events and
22 mutations. We compare driver mutation rates to the exome-wide mutational excess calculated using
23 Non-synonymous vs Synonymous mutation rates (dNdS). We observe mutual exclusivity or co-
24 occurrence of events within and between a number of EAC pathways (GATA factors, Core Cell cycle
25 genes, TP53 regulators and the SWI/SNF complex) suggestive of important functional relationships.
26 These driver variants correlate with tumour differentiation, sex and prognosis. Poor prognostic
27 indicators (SMAD4, GATA4) are verified in independent cohorts with significant predictive value. Over

28 50% of EACs contain sensitising events for CDK4/6 inhibitors which are highly correlated with clinically
29 relevant sensitivity in a panel EAC cell lines and organoids.

30

31 **Introduction**

32 Esophageal cancer is the eighth most common form of cancer world-wide and the sixth most
33 common cause of cancer related death¹. Esophageal Adenocarcinoma (EAC) is the predominant
34 subtype in the west, including the UK and the US. The incidence of EAC in such countries has been
35 rapidly rising, with a seven-fold increase in incidence over the last 35 years in the US². EAC is a highly
36 aggressive neoplasm, usually presenting at a late stage and is generally resistant to chemotherapy,
37 leading to five-year survival rates below 15%³. It is characterised by very high mutation rates in
38 comparison to other cancer types⁴ but also, paradoxically, there is a paucity of recurrently mutated
39 genes. EACs also display dramatic chromosomal instability and thus may be classified as a C-type
40 neoplasm which may be driven mainly by structural variation rather than mutations^{5,6}. Currently our
41 understanding of precisely which genetic events drive the development of EAC is highly limited and
42 consequentially there is a paucity of molecular biomarkers for prognosis or targeted therapeutics
43 available in the clinic.

44 Driver events undergoing positive selection during cancer evolution are a small proportion
45 of total number of genetic events that occur in each tumour⁷. Methods to differentiate driver
46 mutations from passenger mutations use features associated with known driver events to detect
47 regions of the genome, often genes, in which mutations are enriched for these features⁸. The
48 simplest of these features is the tendency of a mutation to co-occur with other mutations in the
49 same gene at a high frequency, as detected by MutsigCV⁹. MutsigCV has been applied on several
50 occasions to EAC cohorts^{6,10,11} and has identified ten known cancer genes as high confidence EAC
51 drivers (TP53, CDKN2A, SMAD4, ARID1A, ERBB2, KRAS, PIK3CA, SMARCA4, CTNNB1 and FBXW7).

52 Analysis of the non-coding genome has been performed by the PCAWG ICGC analysis and identified
53 a significantly mutated enhancer associated with TP53TG1¹². However these analyses leave most
54 EAC cases with only one known driver mutation, usually TP53, due to the low frequency at which
55 other drivers occur. Equivalent analyses in other cancer types have identified three or four drivers
56 per case^{13,14}. Similarly, detection of copy number driver events in EAC has relied on identifying
57 regions of the genome recurrently deleted or amplified, as detected by GISTIC^{10,15-18}. However,
58 GISTIC identifies relatively large regions of the genome, often containing hundreds of genes, with
59 little indication of which specific gene-copy number aberrations (CNAs) may actually confer a
60 selective advantage. There are also several non-selection based mechanisms which can cause
61 recurrent CNAs, such as fragile sites where a low density of DNA replication origins causes frequent
62 structural events at a particular loci. These have not been differentiated properly from selection
63 based recurrent CNAs¹⁹. Epigenetic events, for example methylation, may also be important sources
64 of driver events in EAC but are much more difficult to assess formally for selection.

65 Without proper annotation of the genomic variants which drive the biology of EAC tumours
66 we are left with a very large number of events, most of which are likely to be inconsequential,
67 making it extremely difficult to detect statistical associations between genomic variants and various
68 biological and clinical parameters. To address these issues, we have accumulated a cohort of 551
69 genomically characterised EACs using our esophageal ICGC project, which have high quality clinical
70 annotation, associated whole genome sequencing (WGS) and RNA-seq on cases with sufficient
71 material. We have augmented our ICGC WGS cohort with publically available whole exome²⁰ and
72 whole genome sequencing²¹ data. We have applied a number of complementary driver detection
73 tools to this cohort, using a range of driver associated features combined with analyses of RNA
74 expression to produce a comprehensive assessment and characterisation of mutations and CNAs
75 under selection in EAC. We then use these events to define functional cell processes that have been
76 selectively dysregulated in EAC and identify novel, clinically relevant biomarkers for prognostication,

77 which we have verified in independent cohorts. Finally, we have used this compendium of EAC
78 driver variants to provide an evidence base for targeted therapeutics, which we have tested *in vitro*.

79

80 **Results**

81 **A Compendium of EAC driver events and their functional effects**

82 In 551 EACs we called a total of 11,813,333 single nucleotide variants (SNVs) and small insertions or
83 deletions (Indels), with a median of 6.4 such mutations / Mb (supplementary figure 1), and 286,965
84 copy number aberrations (CNAs). We also identified 134,697 structural variants (SVs) in WGS cases.
85 Mutations or copy number variants under selection were detected using specific driver associated-
86 mutation features (Fig 1A). We use several complementary driver detection tools to detect each
87 feature, and each tool underwent quality control to ensure reliability of results (see methods). These
88 features include highly recurrent mutations within a gene (dNdScv²², ActivedriverWGS²³,
89 MutsigCV2⁹), high functional impact mutations within a gene (OncodriveFM²⁴, ActivedriverWGS²³),
90 mutation clustering (OncodriveClust²⁵, eDriver²⁶ and eDriver3D²⁷) and recurrent amplification or
91 deletion of genes (GISTIC¹⁵) undergoing concurrent over or under-expression (see methods) (Fig
92 1A)⁸.

93 These complementary methods produced highly significant agreement in calling EAC driver
94 genes, particularly within the same feature-type (supplementary figure 2) and on average more than
95 half of the genes identified by one feature were also identified by other features (Fig 1B). In total
96 seventy six EAC driver genes were discovered, 86% of which have not been detected in EAC
97 previously^{10,11,16-18,20} and 69% are known drivers in pan-cancer analyses giving confidence in our
98 methods^{22,28,29}. To detect driver elements in the non-coding genome we used ActiveDriverWGS²³ a
99 recently benchmarked³⁰ method using both function impact prediction and recurrence to determine
100 driver status (Fig 1C, supplementary figure 3). We discovered 21 non-coding driver elements using

101 this method. We have recovered several known non-coding driver elements from the pan-cancer
102 PCAWG analysis¹² including an enhancer on chr7 linked to TP53TG1, a gene required for TP53 action,
103 the only non-coding driver found in EAC in PCAWG and the promoter/5'UTR regions of PTDSS1 and
104 WRD74 which are novel in EAC but were found in other cancer types. We also identified completely
105 novel non-coding cancer driver elements including in the 5'UTR of MMP24 and promoters of two
106 related histones (HIST1H2BO and HIST1H2AM).

107 EAC is notable among cancer types for harbouring a high degree of chromosomal
108 instability²¹. Using GISTIC we identified 149 recurrently deleted or amplified loci across the genome
109 (Fig 2A). To determine which genes within these loci confer a selective advantage when they
110 undergo CNAs we use a subset of 116 cases with matched RNA-seq to detect genes within these loci
111 in which homozygous deletion or amplification causes a significant under or over-expression
112 respectively, a prerequisite for selection of CNAs. The majority of genes in these regions showed no
113 significant CN associated expression change (74%), although work in larger cohorts suggests we may
114 be underpowered to detect small expression changes³¹. We observed highly significant expression
115 changes in 17 known cancer genes within GISTIC peaks such as ERBB2, KRAS and SMAD4 which we
116 designate high-confidence EAC drivers. We also found five tumour suppressor genes where copy
117 number loss was not necessarily associated with expression modulation but tightly associated with
118 presence of mutations leading to LOH, for example ARID1A and CDH11. CDH11 was not identified by
119 our driver gene detection methods but this would suggest it may be a promising candidate for
120 further validation. To determine whether copy number changes in genes not previously associated
121 with cancer may contribute to oncogenesis we searched for genes with similar expression-CN profile
122 as most of our high-confidence drivers (see methods). We found 140 such cases which we
123 designated "candidate copy number (CN) drivers" (supplementary tables 1-4). Not all candidate
124 drivers are likely to be true CN-drivers. However, several candidate drivers such as ZNF131, YES1 and
125 PIBF1 are not accompanied by other drivers in their GISTIC peak and contain extrachromosomal-like
126 events, hence are promising candidates for further study.

127 In a subset of GISTIC loci, we observed extremely high copy number amplification,
128 commonly greater than 100 copies, and these loci were highly correlated with presence of CN-
129 drivers (Ploidy adjusted Copy number >10, Wilcox test, $p < 10^{-6}$) (supplementary figure 4). We use
130 copy number adjusted ploidy to define amplifications as it produces superior correlation with
131 expression data than absolute CN alone. Ploidy of our samples varies from 2-6 (3.5 on average) and
132 hence Ploidy adjusted copy number of >10 cut off translates into >20-60 absolute copies (on average
133 35 copies). To discern a mechanism for these ultra-high amplifications we assessed structural
134 variants (SVs) associated with these events and the copy number steps surrounding them. For many
135 of these events the extreme amplification was produced largely from a single copy number step the
136 edges of which were supported by structural variants with ultra-high read support. Two examples
137 are shown in Fig 2B and further examples in supplementary figure 5. In the first example
138 circularisation and amplification initially occurred around MYC but subsequently incorporated ERBB2
139 from an entirely different chromosome and in the second an inversion has been followed by
140 circularisation and amplification of KRAS. A pattern of extrachromosomal amplification via double
141 minutes has been previously noted in EAC²¹, and hence we refer to this amplification class with
142 ultra-high amplification (Ploidy adjusted Copy number >10) as 'extrachromosomal-like'. Several
143 deletion loci co-align with fragile sites (Fig 2A). Most deletion loci were dominated by heterozygous
144 deletions while a small subset had a far higher percentage of homozygous deletions including
145 CDKN2A and several associated with fragile site loci (Fig 2A). For some cases we may have been
146 unable to identify drivers in loci simply because the aberrations do not occur in the smaller RNA-seq
147 matched cohort.

148 We found extrachromosomal-like amplifications had an extreme and highly penetrant
149 effects on expression while moderate amplification (ploidy adjusted copy number > 2) and
150 homozygous deletion had highly significant (Wilcox test, $p < 10^{-4}$ and $p < 10^{-3}$ respectively) but less
151 dramatic effects on expression with a lower penetrance (Fig 2C). This lack of penetrance was
152 associated with low cellularity (fisher's exact test, expression cut off = 2.5 normalised FPKM, $p < 0.01$)

153 in amplified cases but also likely reflects that genetic mechanisms other than gene-dosage can
154 modulate expression in a rearranged genome. We also detected several cases of over expression or
155 complete expression loss without associated CN changes which may reflect non-genetic mechanisms
156 for driver dysregulation. For example, one case overexpressed ERBB2 at 28-fold median expression
157 however had entirely diploid CN in and surrounding ERBB2 and a second case contained almost
158 complete loss of SMAD4 expression (0.008-fold median expression) despite possessing 5 copies of
159 SMAD4.

160

161 **Landscape of driver Events in EAC**

162 The overall landscape of driver gene mutations and copy number alterations per case is depicted in
163 Fig 3A. These comprise both oncogenes and tumour suppressor genes activated or repressed via
164 different mechanisms. Occasionally different types of events are selected for in the same gene, such
165 as KRAS and ERBB2 which both harbour activating mutations and amplifications in 19% and 18% of
166 cases respectively. Passenger mutations occur by chance in most driver genes. To quantify this we
167 have used the observed:expected mutation ratios (calculated by dNdScv) to estimate the percentage
168 of driver mutations in each gene and in different mutation classes. For many genes, only specific
169 mutation classes appear to be under selection. Many tumour suppressor genes; ARID2, RNF43,
170 ARID1B for example, are only under selection for truncating mutations; ie splice site, nonsense and
171 frameshift Indel mutations, but not missense mutations which are passengers. However, oncogenes,
172 like ERBB2, only contain missense drivers which form clusters to activate gene function in a specific
173 manner. Where a mutation class is <100% driver mutations, mutational clustering can help us define
174 the driver vs passenger status of a mutation (supplementary figure 6). Clusters of mutations
175 occurring in EAC or mutations on amino acids which are mutation hotspots in other cancer types³²
176 (supplementary table 5) are indicated in Fig 3A. Novel EAC drivers of particular interest include B2M,
177 a core component of the MHC class I complex and resistance marker for Immunotherapy³³, MUC6 a

178 secreted glycoprotein involved in gastric acid resistance and ABCB1 a channel pump protein which is
179 associated with multiple instances of drug resistance³⁴. We note that several of these drivers have
180 been previously associated with gastric and colorectal cancer (supplementary table 6)^{14,35}. Lollipop
181 plots showing primary sequence distribution of mutations in these genes are provided
182 (supplementary data).

183 The identification of driver events provides a rich information about the molecular history of
184 each EAC tumour. We detect a median of five events in driver genes per tumour (IQR = 3-7, Mean =
185 5.6) and only a very small fraction of cases have no such events detected (6 cases, 1%). When we
186 remove the predicted percentage of passenger mutations using dnds ratios we find a mean of 4.4
187 true driver events per case which derive more commonly from mutations than CN events (Fig 3B).
188 Using hierarchal clustering of drivers we noted that TP53 mutant cases had significantly more CN
189 drivers (Wilcox test, $p = 0.0032$, supplementary figure 7). dNdScv, one of the driver gene detection
190 methods used, also analyses the genome-wide excess of non-synonymous mutations based on
191 expected mutation rates to assess the total number of driver mutations across the exome which is
192 calculated at 5.4 (95% CIs: 3.5-7.3) in comparison to 2.7 driver mutations which we calculate in our
193 gene-centric analysis after passenger removal. This suggests low frequency driver genes may be
194 prevalent in the EAC mutational landscape (see discussion). Further analysis suggests these missing
195 mutations are mostly missense mutations and our gene-centric analysis captures almost all
196 predicted splice and nonsense drivers (supplementary figure 8). Some of our methods use
197 enrichment of nonsense and splice mutations as a marker of driver genes and hence have a higher
198 sensitivity for these mutations.

199 To better understand the functional impact of driver mutations we analysed expression of
200 driver genes with different mutation types and compared their expression to normal tissue RNA,
201 which was sequenced alongside our tumour samples (Fig 3C). Since surrounding squamous
202 epithelium is a fundamentally different tissue, from which EAC does not directly arise, we have used

203 duodenum and gastric cardia samples as gastrointestinal phenotype controls, likely to be similar to
204 the, as yet unconfirmed, tissue of origin in EAC. A large number of driver genes have upregulated
205 expression in comparison to normal controls, for example TP53 has upregulated RNA expression in
206 WT tumour tissue and in cases with missense (see non-truncating Fig 3C) mutations but RNA
207 expression is lost upon gene truncation. In depth analysis of different TP53 mutation types reveals
208 significant heterogeneity within non-truncating mutations, for example R175H mutations correlate
209 with low RNA expression (supplementary figure 9). Normal tissue expression of CDKN2A suggests
210 that CDKN2A is generally activated in EAC, likely due to genotoxic or other cancer-associated
211 stresses³⁶ and returns to physiologically normal levels when deleted. Heterogeneous expression in
212 WT CDKN2A cases suggest a different mechanism of inhibition such as methylation in some cases.
213 Overexpression of other genes in wild type tumours, such as SIN3A, may confer a selective
214 advantage due to their oncogenic properties, in this case cooperating with MYC, which is also
215 overexpressed in EACs (Fig 3C). A smaller number of driver genes are downregulated in EAC tissue-
216 3/4 of these (GATA4, GATA6 and MUC6) are involved in the differentiated phenotype of
217 gastrointestinal tissues and may be lost with tumour de-differentiation. Driving alterations in these
218 genes have been observed in other GI cancers^{14,37,38} however their oncogenic mechanism is
219 unknown. In most genes we did not observe expression loss at the RNA level with truncation, for
220 instance ARID1A (supplementary figure 10).

221

222 **Dysregulation of specific pathways and processes in EAC**

223 It is known that selection preferentially dysregulates certain functionally related groups of genes and
224 biological pathways in cancer³⁹. This phenomenon is highly evident in EAC, as shown in Fig 4 which
225 depicts the functional relationships between EAC drivers. This provides greater functional
226 homogeneity to the landscape of driver events.

227 While TP53 is the dominant driver in EAC, 28% of cases remain TP53 wildtype. MDM2 is a E3
228 ubiquitin ligase that targets TP53 for degradation. Its selective amplification and overexpression is
229 mutually exclusive with TP53 mutation suggesting it can functionally substitute the effect of TP53
230 mutation via its degradation. Similar mutually exclusive relationships are observed between; KRAS
231 and ERBB2, GATA4 and GATA6 and Cyclin genes (CCNE1, CCND1 and CCND3). Activation of the Wnt
232 pathway occurs in 19% of cases either by mutation of phospho-residues at the N terminus of β -
233 catenin, which prevent degradation, or loss of Wnt destruction complex components like APC. Many
234 different chromatin modifying genes, often belonging to the SWI/SNF complex, are also selectively
235 mutated (31% of cases). In contrast SWI/SNF genes are co-mutated significantly more often than we
236 would expect by chance (fisher's exact test, $p < 0.01$ see methods), suggesting an increased advantage
237 to further mutations once one has been acquired. We also assessed mutual exclusivity and co-
238 occurrence in genes in different pathways and between pathways themselves (Fig 4B). Of particular
239 note are co-occurring relationships between TP53 and MYC, GATA6 and SMAD4, Wnt and Immune
240 pathways as well as mutually exclusive relationships between ARID1A and MYC, gastrointestinal (GI)
241 differentiation and RTK pathways and SWI-SNF and DNA-Damage response pathways. Wnt
242 dysregulation has been previously linked to immune escape⁴⁰ and interestingly was also associated
243 with hyper-mutated cases ($> 50,000$ SNVs or Indels, fisher's exact test, $p = 0.021$, OR= 2.4). We were
244 able to confirm some of these relationships in independent cohorts in different cancer types
245 (supplementary table 7) suggesting some of these may be pan-cancer phenomenon. As shown in Fig
246 4, all of these pathways interact to stimulate the G1 to S phase transition of the cell cycle via
247 promoting phosphorylation of Rb, although many of these pathways have multiple oncogenic or
248 tumour suppressive functions.

249 A number of other driver genes have highly related functional roles including core
250 transcriptional components (TAF1 and POLQ), drivers of immune escape (JAK1 and B2M³³), cell
251 adhesion receptors (CDH1, CHDL and PCDH17), core ribosome components (ELF3 and RPL22), core

252 RNA processing components (GPATCH8 and COIL), ion channels (KCNQ3 and TRPA1) and Ephrin
253 type-A receptors (EPHA2 and EPHA3).

254

255 **Clinical significance of driver variants**

256 Events undergoing selection during cancer evolution influence tumour biology and thus impact
257 tumour aggressiveness, response to treatment and patient prognosis as well as other clinical
258 parameters. Clinical-genomic correlations can provide useful biomarkers but also give insights into
259 the biology of these events.

260 Univariate Cox regression was performed for events in each driver gene with driver events
261 occurring in greater than 5% of EACs (ie after removal of predicted passengers, 16 genes) to detect
262 prognostic biomarkers (Fig 5A). Events in two genes conferred significantly poorer prognosis after
263 multiple hypothesis correction, GATA4 amplification (HR : 0.54 , 95% CI : 0.38 – 0.78, *P* value =
264 0.0008) and SMAD4 mutation or homozygous deletion (HR : 0.60 , 95% CI : 0.42 – 0.84, *P* value =
265 0.003). Both genes remained significant in multivariate Cox regression including pathological TNM
266 staging, resection margin, curative vs palliative treatment intent and differentiation status (GATA4 =
267 HR adjusted : 0.47, 95% CIs adjusted : 0.29 - 0.76, *P* value = 0.002 and SMAD4 = HR adjusted : 0.61,
268 95% CI adjusted : 0.40 – 0.94, *P* value = 0.026). 31% of EACs contain either SMAD4 mutation or
269 homozygous deletion or GATA4 amplification and cases with both genes altered had a poorer
270 prognosis (Fig 5B). We validated the poor prognostic impact of SMAD4 events in an independent
271 TCGA gastroesophageal cohort (HR = 0.58, 95% CI = 0.37 – 0.90, *P* value = 0.014) (Fig 5C) and we also
272 found GATA4 amplifications were prognostic in a cohort of TCGA pancreatic cancers (HR = 0.38 95%
273 CI: 0.18 – 0.80, *P* value = 0.011) (Fig 5D), the only available cohort containing a feasible number of
274 GATA4 amplifications. The prognostic impact of GATA4 has been suggested in previously published
275 independent EAC cohort¹⁷ although it did not reach statistical significance after FDR correction and
276 SMAD4 expression loss has been previously linked to poor prognosis in EAC⁴¹. We also noted stark

277 survival differences between cases with SMAD4 events and cases in which TGF β receptors were
278 mutated (Fig 5E, HR = 5.6, 95% CI : 1.7 – 18.2, *P* value = 0.005) in keeping with the biology of the
279 TGF β pathway where non-SMAD TGF β signalling is known to be oncogenic⁴².

280 In addition to survival analyses we also assessed driver gene events for correlation with
281 various other clinical factors including differentiation status, sex, age and treatment response. We
282 found Wnt pathway mutations had a strong association with well differentiated tumours (*p*=0.001,
283 OR = 2.9, fisher's test, see methods, Fig 5F). We noted interesting differences between female
284 (*n*=81) and male (*n*=470) cases. Female cases were enriched for KRAS mutation (*p* = 0.001, fisher's
285 exact test) and TP53 wildtype status (*p* = 0.006, fisher's exact test) (Fig 5G). This is of particular
286 interest given the male predominance of EAC³.

287

288 **Targeted therapeutics using EAC driver events.**

289 The biological distinctions between normal and cancer cells provided by driver events can be used to
290 derive clinical strategies for selective cancer cell killing. To investigate whether the driver events in
291 particular genes and/or pathways might sensitise EAC cells to certain targeted therapeutic agents
292 we used the Cancer Biomarkers database⁴³. We calculated the percentage of our cases which
293 contain EAC-driver biomarkers of response to each drug class in the database (summary shown Fig
294 6A, and full data supplementary table 8). Aside from TP53, which has been problematic to target
295 clinically so far, we found a number of drugs with predicted sensitivity in >10% of EACs including
296 EZH2 inhibitors for SWI/SNF mutant cancers (23%, and 33% including other SWI/SNF EAC
297 drivers), and BET inhibitors which target KRAS activated and MYC amplified cases (25%). However,
298 by far the most significantly effective drug was predicted to be CDK4/6 inhibitors where >50% of
299 cases harboured sensitivity causing events in the receptor tyrosine kinase (RTK) and core cell cycle
300 pathways (eg in CCND1, CCND3 and KRAS).

301 To verify that these driver events would also sensitise EAC tumours to such inhibitors we
302 used a panel of thirteen EAC or Barrett's HGD cell lines, which share similar genomic changes and
303 driver events^{44,45}, which have undergone whole genome sequencing⁴⁶ and assessed them for
304 presence of EAC driver events (Fig 6B). The mutational landscape of these lines was broadly
305 representative of EAC tumours. We found that the presence of cell cycle and or RTK activating driver
306 events was highly correlated with response to two FDA approved CDK4/6 inhibitors, Ribociclib and
307 Palbociclib and several cell lines were sensitive below maximum tolerated blood concentrations in
308 humans (Fig 6B, supplementary table 9, supplementary figure 11)⁴⁷. Such EAC cell lines had
309 comparable sensitivity to T47D which is derived from an ER +ve breast cancer where CDK4/6
310 inhibitors have been FDA approved. We noted three cell lines without sensitising events which were
311 highly resistant, with little drug effect even at 4000 nanomolar concentrations, similar to a known Rb
312 mutant resistant line breast cancer cell line (MDA-MB-468). Two of these three cell lines harbour
313 amplification of CCNE1 which is known to drive resistance to CDK4/6 inhibitors by bypassing CDK4/6
314 and causing Rb phosphorylation via CDK2 activation⁴⁸. To verify these effects in a more
315 representative model of EAC we treated three whole genome sequenced EAC organoid cultures⁴⁹
316 with Palbociclib and Ribociclib as well as a more recently approved CDK4/6 inhibitor, Abemaciclib. As
317 was observed in cell lines, Cell cycle and RTK driver events were present only in the more sensitive
318 organoids and CCNE1 activation in the most resistant (Fig 6C). We found Abemaciclib to be
319 significantly more potent in comparison to both other CDK4/6 inhibitors, both in organoids and cell
320 lines (supplementary figure 10). We note that the maximum tolerated blood doses of Abermaciclib
321 achieved in the clinic were also higher than the other CDK4/6 inhibitors⁵⁰, within the range of
322 sensitivity achieved in several cell lines and organoids cultures.

323

324

325

326 **Discussion**

327 We present here a detailed catalogue of coding and non-coding genomic events that have been
328 selected for during the evolution of esophageal adenocarcinoma. These events have been
329 characterised in terms of their relative impact, related functions, mutual exclusivity and co-
330 occurrence and expression in comparison to normal tissues, producing insights into EAC biology. We
331 have used this set of biologically important gene alterations to identify prognostic biomarkers and
332 actionable genomic events for personalised medicine.

333 While clinical annotation and matched RNA data is a strength of this study, in some cases we
334 may have been unable to assess selected variants for survival associations or expression changes
335 which were detected in the full 551 cohort, due to lack of representation in clinically annotated or
336 RNA matched sub cohorts. Despite rigorous analyses to detect selected events, assessment of the
337 global excess of mutations by dNdScv suggests we are unable to detect all events selected in EAC,
338 similar to many other cancer types²². All driver gene detection methods which we have used rely on
339 driver mutation re-occurrence in a gene to some degree. Many of these undetected driver
340 mutations are hence likely to be spread across a large number of genes whereby each is mutated at
341 low frequency across EAC patients. This tendency for low frequency EAC drivers may be responsible
342 for the low yield of MutsigCV in previous cohorts and may suggest that C-type cancers such as EAC,
343 are not less 'mutation-driven' than M-type cancers but rather that their mutational drivers are
344 spread across a larger number of genes⁵. The identification of these very low frequency mutations
345 will require substantially different detection techniques to those which are currently in wide spread
346 use and such methods are in development⁵¹ although they require validation. Undoubtedly many
347 copy number drivers are also left undiscovered and validation of candidates identified here is an
348 important avenue of future work.

349 While a number of previous reports have attempted to detect EAC drivers, they have had a
350 limited yield per case for a variety of reasons. The first such study²⁰ used methods which, despite

351 being well regarded at the time, were subsequently discredited⁹. Hence a number of known false
352 positive genes (EYS, SYNE1 and CNTTAP5) were erroneously reported as drivers, along with an
353 additional unknown number of genes. Since then a number of reports, including our own, on
354 medium and large cohort sizes using MutsigCV^{10,11,18} were only able to detect a small number of
355 mutational driver genes (7, 5 and 15 in each study). By using both a large cohort and more
356 comprehensive methodologies we have significantly increased this figure to 66 mutational driver
357 genes (excluding CN drivers). Detection of driver CNAs has previously relied on GISTIC to detect
358 recurrently mutated regions^{10,15-18} but no analyses have been performed to evidence which genes in
359 these large regions are true drivers. Many of the genes annotated by such papers are unlikely to be
360 CN drivers from this analysis due to their lack of expression modulation with CNAs (eg YEATS4 and
361 MCL1), the role of recurrent heterozygous losses to drive LOH in some mutational drivers (ARID1A
362 and CDH11) or their association with fragile sites (PDE4D, WWOX, FHIT). Conversely, we have been
363 able to identify novel EAC copy number drivers (eg CCND3, AXIN1, PPM1D and APC).

364 A number of discoveries made in this work require further investigation. Functional
365 characterisation of many of the driver genes described is needed to understand why they are
366 advantageous to EAC tumours and how they modify EAC biology. Particularly interesting are the GI
367 specific genes GATA4, GATA6 and MUC6 which modulate prognosis and have expression loss during
368 the transition from normal to tumour tissue. Biological pathways and processes that are selectively
369 dysregulated deserve particular attention in this regard as do the gene pairs or groups with mutually
370 exclusive or co-occurring relationships such as MYC and TP53 or SWI/SNF factors, suggestive of
371 particular functional relationships. Prospective clinical work to verify and implement SMAD4 and
372 GATA4 biomarkers in this study would be worthwhile. While EAC is a poor prognosis cancer type,
373 significant heterogeneity of survival outcome makes triaging patients in treatment groups an
374 important part of clinic practice which could be improve using better prognostication. Whole
375 genome or whole exome sequencing may be impractical for use in the clinic, however targeted NGS
376 panels to detect mutations and copy number alterations have been implemented to detect genomic

377 biomarkers in a cost effective and sensitive manner for some cancer types⁵². In EAC development of
378 a customised panel is likely to be required on the basis of this analysis. A number of targeted
379 therapeutics may provide clinic benefit to EAC cases based on their individual genomic profile. In
380 particular CDK4/6 inhibitors deserve considerable attention as an option for EAC treatment as they
381 are, by a significant margin, the treatment to which the most EACs harbour sensitivity-causing driver
382 events, excluding TP53 as an unlikely therapeutic biomarker. The in vitro validation of these
383 biomarkers for CDK4/6 inhibitors in EAC is also persuasive of possible clinical benefit using a targeted
384 approach.

385 In summary this work provides a detailed compendium of mutations and copy number
386 alterations undergoing selection in EAC which have functional and clinical impact on tumour
387 behaviour. This comprehensive study provides us with useful insights into the nature of EAC tumours
388 and should pave the way for evidence based clinical trials in this poor prognosis disease.

389

390

391 **Oesophageal Cancer Clinical and Molecular Stratification (OCCAMS) Consortium:**

392 Rebecca C. Fitzgerald¹, Ayesha Noorani¹, Paul A.W. Edwards^{1,2}, Nicola Grehan¹, Barbara Nutzinger¹,
393 Caitriona Hughes¹, Elwira Fidziukiewicz¹, Jan Bornschein¹, Shona MacRae¹, Jason Crawte¹, Alex
394 Northrop¹, Gianmarco Contino¹, Xiaodun Li¹, Rachel de la Rue¹, Maria O'Donovan^{1,3}, Ahmad
395 Miremadi^{1,3}, Shalini Malhotra^{1,3}, Monika Tripathi^{1,3}, Simon Tavaré², Andy G. Lynch², Matthew
396 Eldridge², Maria Secrier², Lawrence Bower², Ginny Devonshire², Juliane Perner², Sriganesh
397 Jammula², Jim Davies⁵, Charles Crichton⁵, Nick Carroll⁶, Peter Safranek⁶, Andrew Hindmarsh⁶,
398 Vijayendran Sujendran⁶, Stephen J. Hayes^{7,14}, Yeng Ang^{7,8,29}, Shaun R. Preston⁹, Sarah Oakes⁹, Izhar
399 Bagwan⁹, Vicki Save¹⁰, Richard J.E. Skipworth¹⁰, Ted R. Hupp¹⁰, J. Robert O'Neill^{10,23}, Olga Tucker^{11,33},
400 Andrew Beggs^{11,28}, Philippe Tanriere¹¹, Sonia Puig¹¹, Timothy J. Underwood^{12,13}, Fergus Noble¹², Jack
401 Owsley¹², Hugh Barr¹⁵, Neil Shepherd¹⁵, Oliver Old¹⁵, Jesper Lagergren^{16,25}, James Gossage^{16,24},
402 Andrew Davies^{16,24}, Fujun Chang^{16,24}, Janine Zylstra^{16,24}, Ula Mahadeva¹⁶, Vicky Goh²⁴, Francesca D.
403 Ciccarelli²⁴, Grant Sanders¹⁷, Richard Berrisford¹⁷, Catherine Harden¹⁷, Mike Lewis¹⁸, Ed Cheong¹⁸,
404 Bhaskar Kumar¹⁸, Simon L Parsons¹⁹, Irshad Soomro¹⁹, Philip Kaye¹⁹, John Saunders¹⁹, Laurence
405 Lovat²⁰, Rehan Haidry²⁰, Laszlo Igali²¹, Michael Scott²², Sharmila Sothi²⁶, Sari Suortamo²⁶, Suzy
406 Lishman²⁷, George B. Hanna³¹, Krishna Moorthy³¹, Christopher J. Peters³¹, Anna Grabowska³², Richard
407 Turkington³⁴.

408

409 ¹ Medical Research Council Cancer Unit, Hutchison/Medical Research Council Research Centre,
410 University of Cambridge, Cambridge, UK

411 ² Cancer Research UK Cambridge Institute, University of Cambridge, Cambridge, UK

412 ³ Department of Histopathology, Addenbrooke's Hospital, Cambridge, UK

413 ⁴ Oxford ComLab, University of Oxford, UK, OX1 2JD

414 ⁵ Department of Computer Science, University of Oxford, UK, OX1 3QD

415 ⁶ Cambridge University Hospitals NHS Foundation Trust, Cambridge, UK, CB2 0QQ

416 ⁷ Salford Royal NHS Foundation Trust, Salford, UK, M6 8HD

417 ⁸ Wigan and Leigh NHS Foundation Trust, Wigan, Manchester, UK, WN1 2NN

418 ⁹ Royal Surrey County Hospital NHS Foundation Trust, Guildford, UK, GU2 7XX

419 ¹⁰ Edinburgh Royal Infirmary, Edinburgh, UK, EH16 4SA

420 ¹¹ University Hospitals Birmingham NHS Foundation Trust, Birmingham, UK, B15 2GW

421 ¹² University Hospital Southampton NHS Foundation Trust, Southampton, UK, SO16 6YD

422 ¹³ Cancer Sciences Division, University of Southampton, Southampton, UK, SO17 1BJ

423 ¹⁴ Faculty of Medical and Human Sciences, University of Manchester, UK, M13 9PL

424 ¹⁵ Gloucester Royal Hospital, Gloucester, UK, GL1 3NN

425 ¹⁶ Guy's and St Thomas's NHS Foundation Trust, London, UK, SE1 7EH

426 ¹⁷ Plymouth Hospitals NHS Trust, Plymouth, UK, PL6 8DH

427 ¹⁸ Norfolk and Norwich University Hospital NHS Foundation Trust, Norwich, UK, NR4 7UY

428 ¹⁹ Nottingham University Hospitals NHS Trust, Nottingham, UK, NG7 2UH

429 ²⁰ University College London, London, UK, WC1E 6BT

430 ²¹ Norfolk and Waveney Cellular Pathology Network, Norwich, UK, NR4 7UY

431 ²² Wythenshawe Hospital, Manchester, UK, M23 9LT

432 ²³ Edinburgh University, Edinburgh, UK, EH8 9YL

433 ²⁴ King's College London, London, UK, WC2R 2LS

434 ²⁵ Karolinska Institutet, Stockholm, Sweden, SE-171 77

435 ²⁶ University Hospitals Coventry and Warwickshire NHS, Trust, Coventry, UK, CV2 2DX

436 ²⁷ Peterborough Hospitals NHS Trust, Peterborough City Hospital, Peterborough, UK, PE3 9GZ

437 ²⁸ Institute of Cancer and Genomic sciences, University of Birmingham, B15 2TT

438 ²⁹ GI science centre, University of Manchester, UK, M13 9PL.

439 ³⁰ Queen's Medical Centre, University of Nottingham, Nottingham, UK, NG7 2UH

440 ³¹ Department of Surgery and Cancer, Imperial College London, UK, W2 1NY

441 ³² Queen's Medical Centre, University of Nottingham, Nottingham, UK

442 ³³ Heart of England NHS Foundation Trust, Birmingham, UK, B9 5SS.

443 ³⁴ Centre for Cancer Research and Cell Biology, Queen's University Belfast, Northern Ireland, UK, BT7
444 1NN.

445

446

447 **Author contributions:**

448 RCF and AMF conceived the overall study. AMF and SJ analysed the genomic data and performed

449 statistical analyses. RCF, AMF and XL designed the experiments. AMF, XL and JM performed the

450 experiments. GC contributed to the Structural variant analysis and data visualisation. SK helped

451 compile the clinical data and aided statistical analyses. JP and SA produced and QC'ed the RNA-seq
452 data. EO aided the whole genome sequencing of EAC cell lines. SM and NG coordinated the clinical
453 centres and were responsible for sample collections. ME benchmarked our mutation calling
454 pipelines. MO led the pathological sample QC for sequencing. LB and GD ran variant calling
455 pipelines. RCF and ST supervised the research. RCF and ST obtained funding. AMF and RCF wrote the
456 manuscript. All authors approved the manuscript.

457

458 **The authors declare no competing interests.**

459 **Sequencing data will be deposited in a publicly assessable database before publication**

460 **Code associated with the analysis is available upon request.**

461 **The study was registered (UKCRNID 8880), approved by the Institutional Ethics**
462 **Committees (REC 07/H0305/52 and 10/H0305/1), and all subjects gave individual**
463 **informed consent.**

464

465 **OCCAMS was funded by a programme grant from Cancer Research UK (RG66287). We thank**
466 **the Human Research Tissue Bank, which is supported by the National Institute for Health**
467 **Research (NIHR) Cambridge Biomedical Research Centre, from Addenbrooke's**
468 **Hospital. Additional infrastructure support was provided from the CRUK funded**
469 **Experimental Cancer Medicine Centre.**

470

471 **Acknowledgements**

472 We would like to thank Dr. Adam Bass and Dr. Nic Waddel for providing data in Dulak et al 2013 and
473 Nones et al 2014 respectively, also included in our previous publication Secrier et al 2016. Inclusion

474 of this data allowed augmentation of our ICGC cohort and a greater sensitivity for the detection of
475 EAC driver genes.

476

477 **Figure Legends:**

478 **Figure 1 Detection of EAC driver Genes.** a. Types of driver-associated features used to detect
479 positive selection in mutations and copy number events with examples of genes containing such
480 features b. Coding driver genes identified and their driver-associated features. c. Non-coding driver
481 elements detected and their element types.

482

483 **Figure 2. Copy number variation under positive selection.** a. Recurrent copy number changes across
484 the genome identified by GISTIC. Frequency of different CNV types are indicated as well as the position
485 of CNV high confidence driver genes and candidate driver genes. The q value for expression correlation
486 with amplification and homozygous deletion is shown for each gene within each amplification and
487 deletion peaks respectively and occasions of significant association between LOH and mutation are
488 indicated in green. Purple deletion peaks indicate fragile sites. b. Examples of Extrachromosomal-like
489 amplifications suggested by very high read support SVs at the boundaries of highly amplified regions
490 produced from a single copy number step. In the first example (bi) two populations of
491 extrachromosomal DNA are apparent (biii), one amplifying only MYC and the second also
492 incorporating ERBB2 from a different chromosome. In the second example (bii) an inversion has
493 occurred before circularization and amplification around KRAS (biv). c. Relationship between copy
494 number and expression in CN driver genes.

495

496 **Figure 3. The driver gene landscape of Esophageal Adenocarcinoma.** a. Driver mutations or CNVs are
497 shown for each patient. Amplification is defined as >2 Copy number adjusted ploidy (2 x ploidy of that

498 case) and extrachromosomal amplification as >10 Copy number adjusted ploidy (10 x ploidy for that
499 case). Driver associated features for each driver gene are displayed to the left. On the right the
500 percentages of different mutation and copy number changes are displayed, differentiating between
501 driver and passenger mutations using dNdScv, and the % of predicted drivers by mutation type is
502 shown. Above the plot are the number of driver mutations per sample with an indication of the mean
503 (red line = 5). **b.** Assessment of driver event types per case and comparison to exome-wide excess of
504 mutations generated by dNdScv. **c.** Expression changes in EAC driver genes in comparison to normal
505 intestinal tissues. Genes with expression changes of note are shown.

506

507 **Figure 4. Biological pathways undergoing selective dysregulation in EAC.** **a.** Biological Pathways
508 dysregulated by driver gene mutation and/or CNVs. WT cases for a pathway are not shown. Inter
509 and intra-pathway interactions are described and mutual exclusivities and/or associations between
510 genes in a pathway are annotated. GATA4/6 amplifications have a mutually exclusive relationship
511 although this does not reach statistical significance (fisher's exact test $p=0.07$ OR =0.52). **b.** Pairwise
512 assessment of mutual exclusivity and association in EAC driver genes and pathways.

513

514 **Figure 5. Clinical significance of Driver events in EAC.** **a.** Hazard ratios and 95% confidence
515 intervals for Cox regression analysis across all drivers genes with at least a 5% frequency of driver
516 alterations * = $q < 0.05$ after BH adjustment. **b.** Kaplan-Meier curves for EACs with different status of
517 significant prognostic indicators (GATA4 and SMAD4). **c.** Kaplan-Meier curves for different
518 alterations in the TGFbeta pathway. **d.** Kaplan-Meier curves showing verification GATA4 prognostic
519 value in GI cancers using a pancreatic TCGA cohort. **e.** Kaplan-Meier curves showing verification
520 SMAD4 prognostic value in Gastroesophageal cancers using a gastroesophageal TCGA cohort. **f.**
521 Differentiation bias in tumours containing events in Wnt pathway driver genes. **g.** Relative frequency
522 of KRAS mutations and TP53 mutations driver gene events in females vs males (fishers exact test).

523

524 **Figure 6. CDK4/6 inhibitors in EAC. a.** Drug classes for which sensitivity is indicated by EAC driver
525 genes with data from the Cancer Biomarkers database³⁶. **b.** Area under the curve (AUC) of sensitivity
526 is shown in a panel of 13 EAC and Be high grade dysplasia cell lines with associated WGS and their
527 corresponding driver events, based on primary tumour analysis. Also AUC is shown for two control
528 lines T47D, an ER +ve breast cancer line (+ve control) and MDA-MB-468 a Rb negative breast cancer
529 (-ve control). *CCNE1 is a known marker of resistance to CDK4/6 inhibitors due to its regulation of
530 Rb downstream of CDK4/6 hence bypassing the need for CDK4/6 activity (see figure 4). **c.** Response
531 of organoid cultures to three FDA approved CDK4/6 inhibitors and corresponding driver events.

532

533 **Supplementary figure legends**

534 **Supplementary figure 1.** Distribution of small scale mutations (SNVs and Indels) across the 551 EAC
535 cohort. Red line indicates the median mutations per case (6.4)

536

537 **Supplementary Figure 2. Concordance between driver gene detection methods. A.** Hierarchical
538 clustering between tools based on gene identified. **B** Genes identified by each tool.

539

540 **Supplementary Figure 3.** Frequency and significance of EAC non-coding drivers from
541 ActiveDriverWGS. **a.** The observed and expected mutation counts found on each element in
542 ActiveDriverWGS. **b.** The *fdr* for each element in ActiveDriverWGS.

543

544 **Supplementary Figure 4.** Frequency of Extrachromosomal like events (CN adjusted Ploidy >10)
545 in GISTIC amplification peaks and presence of high confidence drivers in those peaks indicated.

546

547 **Supplementary figure 5. Examples of Normal amplification (PLiody-adjusted CN >2 & <10) and**
548 **Extrachromosomal-like amplification (ploidy-adjusted CN >10) events.** 1-10 = Extrachromosomal-
549 like amplification and 11-20 = Normal amplification events. Events were picked at random using
550 runif() function in R. SV and CNAs surrounding events are shown. Features indicative of
551 extrachromosomal double minute (DM) formations include sharp, large CN steps, SVs with high read
552 support at the edges of these steps and when not derived from a continuous region of the genome
553 CN regions in the DM may have the same CN status (taking into account other additional events
554 which may have occurred in that region). These features are enriched in the extrachromosomal-like
555 events, although example 20 may be a low-copy number extrachromosomal event. It should be
556 noted that SV calling using short read sequencing techniques such as in this study has a relatively
557 low sensitivity and accuracy for the precise localisation of many SV break points. Examples continue
558 over four pages.

559

560 **Supplementary Figure 6.** A scheme demonstrating how to use mutational clustering along with dnds
561 ratios to estimate the probability of a particular mutation being a driver. In this case the dnds ratio
562 suggests 2/3 of missense mutations are drivers hence 10/15. 8 missense mutation lie in a mutational
563 cluster, in this case of known significance in the N-terminal of B-Catenin, making it likely that these
564 are drivers and hence most (2/7) other mutations are passengers. Similarly, mutations on amino
565 acids known to be hyper mutated in other cancer types (see Supplementary table 5, for instance if
566 we found a single KRAS G12 mutation) can be considered likely drivers.

567

568 **Supplementary Figure 7. Hierarchical Clustering of samples based on presence of driver variants**
569 **with genes ordered by pathway membership.**

570

571 **Supplementary Figure 8.** A detailed breakdown of mutation and copy number types per case and a
572 breakdown of exome wide dnds excess for different mutation types (note that exome wide indel
573 cannot be calculated excess as they have no synonymous mutation equivalent, although a null
574 model is used in the per gene dnds method to use them to detect driver genes). Error bars indicate
575 95% confidence intervals for exome-wide dnds mutation excess assessment.

576

577 **Supplementary Figure 9.** TP53 expression in different TP53 mutation types in comparison to TP53
578 WT tumours and normal duodenum and gastric cardia tissues.

579

580 **Supplementary Figure 10.** Expression of all EAC driver genes across different genomic states for the
581 gene in question in 116 EAC tumours, and in comparison to duodenum and gastric cardia tissues.

582

583 **Supplementary Figure 11.** Growth inhibition responses of EAC cell lines and control lines to CDK4/6
584 inhibitors Palbociclib and Ribociclib. A subset of cell lines also received treatment with Abemaciclib
585 which shows efficacy in such cell lines as well as in organoids (Fig 6C).

586

587 **Methods**

588 **Cohort, sequencing and calling of genomic events**

589 380 cases (69%) of our EAC cohort were derived from the esophageal adenocarcinoma WGS ICGC
590 study, for which samples are collected through the UK wide OCCAMS (Oesophageal Cancer
591 Classification and Molecular Stratification) consortium. The procedures for obtaining the samples,
592 quality control processes, extractions and whole genome sequencing are as previously described¹⁸.

593 Strict pathology consensus review was observed for these samples with a 70% cellularity
594 requirement before inclusion. Comprehensive clinical information was available for the ICGC-
595 OCCAMS cases. In addition, previously published samples were included in the analysis from Dulak
596 et al 2013²⁰ – 139 WES and 10 WGS (total 27%) and Nones et al 2014²¹ with 22 WGS samples (4%) to
597 total 551 genome characterised EACs. RNA-seq data was available from our ICGC WGS samples
598 (116/380). BAM files for all samples (include those from Dulak et al 2013 and Nones et al 2014) were
599 run through our alignment (BWA-MEM), mutation (Strelka), copy number (ASCAT) and structural
600 variant (Manta) calling pipelines, as previously described¹⁸. Our methods were benchmarked against
601 various other available methods and have among the best sensitivity and specificity for variant
602 calling (ICGC benchmarking exercise⁵³). Mutation and copy number calling on cell lines was
603 performed as previously described⁴⁶.

604 Total RNA was extracted using All Prep DNA/RNA kit from Qiagen and the quality was checked on
605 Agilent 2100 Bioanalyzer using RNA 6000 nano kit (Agilent). Qubit High sensitivity RNA assay kit from
606 thermo fisher was used for quantification. Libraries were prepared from 250ng RNA, using TruSeq
607 Stranded Total RNA Library Prep Gold (Ribo-zero) kit and ribosomal RNA (nuclear, cytoplasmic and
608 mitochondrial rRNA) was depleted, whereby biotinylated probes selectively bind to ribosomal RNA
609 molecules forming probe-rRNA hybrids. These hybrids were pulled down using magnetic beads and
610 rRNA depleted total RNA was reverse transcribed. The libraries were prepared according to Illumina
611 protocol⁵⁴. Paired end 75bp sequencing on HiSeq4000 generated the paired end reads. For normal
612 expression controls we chose gastric cardia tissue, from which some hypothesise Barrett's may arise,
613 and duodenum which contains intestinal histology, including goblet cells, which mimics that of
614 Barrett's. We did not use Barrett's tissue itself as a normal control given the heterogeneous and
615 plentiful phenotypic and genomic changes which it undergoes early in its pathogenesis.

616

617

618

619 **Analysing EAC mutations for selection**

620 To detect positively selected mutations in our EAC cohort, a multi-tool approach across various
621 selection related 'Features' (Recurrence, Functional impact, Clustering) was implemented in order to
622 provide a comprehensive analysis. This is broadly similar to several previous approaches^{8,12}.
623 dNdScv²², MutsigCV⁹, e-Driver²⁶, ActivedriverWGS and e-Driver3D²⁷ were run using the default
624 parameters. To run OncodriverFM²⁴, Polyphen⁵⁵ and SIFT⁵⁶ were used to score the functional impact
625 of each missense non-synonomous mutation (from 0, non-impactful to 1 highly impactful),
626 synonymous mutation were given a score of 0 impact and truncating mutations (Non-sense and
627 frameshift mutations) were given a score of 1. Any gene with less than 7 mutations, unlikely to
628 contain detectable drivers using this method, was not considered to decrease the false discovery
629 rate. OncodriveClust was run using a minimum cluster distance of 3, minimum number of mutations
630 for a gene to be considered of 7 and with a stringent probability cut off to find cluster seeds of $p =$
631 $Ex10^{-13}$ to prevent infiltration of large numbers of, likely, false positive genes. For all tool outputs we
632 undertook quality control including Q-Q plots to ensure no tool produces inflated q-values and each
633 tool produced at least 30% known cancer genes. Two tools were removed from the analysis due to
634 failure for both of these parameters at quality control (Activedriver⁵⁷ and Hotspot³²). For three of the
635 QC-approved tools (dNdScv, OncodriveFM, MutsigCV) where this was possible we also undertook an
636 additional *fdr* reducing analysis by re-calculating q values based on analysis of known cancer genes
637 only^{22,28,29} as has been previously implemented^{22,58}. Significance cut offs were set at $q < 0.1$ for coding
638 genes. Tool outputs were then put through various filters to remove any further possible false
639 positive genes. Specifically, genes where <50% of EAC cases had no expression (TPM<0.1) in our
640 matched RNA-seq cohort were removed and, using dNdScv, genes with no significant mutation
641 excess (observed: expected ratio > 1.5:1) of any single mutation type were also removed. We also
642 removed two (MT-MD2, MT-MD4) mitochondrial genes which were highly enriched for truncating

643 mutations and were frequently called in OncodriveFM as well as other tools. This is may be due to
644 the different mutational dynamics, caused by ROS from the mitochondrial electron transport chain,
645 and the high number of mitochondrial genomes per cell which enables significantly more
646 heterogeneity. These factors prevent the tools used from calculating an accurate null model for
647 these genes however they may be worthy of functional investigation. For non-coding elements
648 called by ActivedriverWGS filtering for expression or dnds was not possible and despite recent
649 benchmarking³⁰ are not so well established. Hence we took a more cautious approach with general
650 significance cut offs of $q < 0.001$ and $q < 0.1$ for previously identified elements in PCAWG¹². Q values
651 were not recalculated for Driver elements only but $q < 0.1$ for known elements was based on all
652 elements. To calculate exome-wide mutational excess hypermutated cases (>500 exonic mutations)
653 were removed and the global non-synonymous dnds ratios were applied to all dndscv annotated
654 mutations excluding “synonymous” and “no SNV” annotations as described in Martincorena et al²².

655

656 **Detecting selection in CNVs**

657 ASCAT raw CN values were used to detected frequently deleted or amplified regions of the genome
658 using GISTIC2.0¹⁵. To determine which genes in these regions confer a selective advantage, CNVs
659 from each gene within a GISTIC identified loci were correlated with TPM from matched RNA-seq in a
660 sub-cohort of 116 samples and with mutations across all 551 samples. To call copy number in genes
661 which spanned multiple copy number segments in ASCAT we considered the total number of full
662 copies of the gene (ie the lowest total copy number). Occasionally ASCAT is unable to confidently call
663 the copy number in a highly aberrant genomic regions. We found that the expression of genes in
664 such regions matched well what we would expect given the surrounding copy number and hence we
665 used the mean of the two adjacent copy number fragments to call copy number in the gene in
666 question. We found amplification peak regions identified by GISTIC2.0 varied significantly in precise
667 location both in analysis of different sub-cohorts and when comparing to published GISTIC data from

668 EACs^{10,16,17}. A peak would often sit next to but not overlapping a well characterised oncogene or
669 tumour suppressor. To account for this, we widened the amplification peak sizes upstream and
670 downstream by twice the size of each peak to ensure we captured all possible drivers. Our
671 expression analysis allows us to then remove false positives from this wider region and called drivers
672 were still highly enriched for genes closer to the centre of GISTIC peak regions.

673 To detect genes in which amplification correlated with increased expression we compared
674 expression of samples with a high CN for that gene (above 10th percentile CN/Ploidy) with those
675 which have a normal CN (median +/- 1) using the Wilcox rank-sum test and using the specific
676 alternative hypothesis that high CN would lead to increased expression. Q-values were then
677 generated based on Benjamini & Hochberg method, not considering genes without significant
678 expression in amplified samples (at least 75% amplified samples with TPM > 0.1) and considering
679 $q < 0.001$ as significant. We also included an additional known driver gene only FDR reduction analysis
680 as previously described for mutational drivers with $q < 0.1$ considered as significant given the
681 additional evidence for these genes in other cancer types. We also included MYC despite its $q = 0.11$
682 for expression correlation. This is due to frequent non-amplification associated overexpression of
683 MYC when compared to normal controls and otherwise MYC is well evidence by a very close
684 proximity to the peak centre (top 4 genes) and its high rate of amplification (19%). We took the
685 same approach to detect genes in which homozygous deletion correlated with expression loss.
686 Expression modulation was a highly specific marker for known CN driver genes and was not a
687 widespread feature in most recurrently copy number variant genes. However, while expression
688 modulation is a requirement for selection of CNV only drivers, it is not sufficient evidence alone and
689 hence we grouped such genes into those which have been characterised as drivers previously in
690 other cancer types (high confidence EAC CN drivers) and other genes (Candidate EAC CN drivers)
691 which await functional validation. We used fragile site regions detected in Wala et al 2017⁵⁹. We also
692 defined regions which may be recurrently heterozygous deleted, without any significant expression
693 modulations, to allow LOH of tumour suppressor gene mutations. To do this we analysed genes with

694 at least 5 mutations in the matched RNA cohort for association between LOH (ASCAT minor allele =
695 0) and mutation using fisher's exact test and generated q values using the Benjamini & Hochberg
696 method. The analysis was repeated on known cancer genes only for reduced FDR and $q < 0.05$
697 considered significant for both analyses. For those high confidence drivers we chose to define
698 amplification as CN/ploidy (referred to as Ploidy adjusted copy number) this produces superior
699 correlation with expression. We chose a cut off for amplification at CN/ploidy = 2 as has been
700 previously used, and as causes a highly significant increase in expression in our CN-driver genes.

701

702 **Pathways and relative distributions of genomic events**

703 The relative distribution of driver events in each pathway was analysed using a fisher's exact test in
704 the case of pair-wise comparisons including WT cases. In the case of multi-gene comparisons such as
705 the Cyclins we calculate the p value and odds ratio for each pair in the group by fisher's exact test
706 and combine p values using the Fisher method, Genes without comparable Odds ratios to the rest of
707 the genes in question were removed. For this analysis we also remove highly mutated cases (>500
708 exonic mutations, 41/551) as they bias distribution of genes towards co-occurrence. We repeated
709 this analyses across all pairs of driver genes using BH multiple hypothesis correction. We validated
710 these relationships in independent TCGA cohorts of other GI cancers where we could find cohorts
711 with reasonable numbers of the genomic events in question (not possible for GATA4/6 for instance)
712 using the cBioportal web interface tool⁶⁰.

713

714 **Correlating genomics with the clinical phenotype**

715 To find genomic markers for prognosis we undertook univariate Cox regression for those driver
716 genes present in >5% of cases (16) along with Benjamini & Hochberg false discovery correction. We
717 considered only these genes to reduce our false discover rate and because other genes were unlikely

718 to impact on clinical practise given their low frequency in EAC. We validated SMAD4, in the TCGA
719 gastroesophageal cohort which had a comparable frequency of these events, but notably is
720 composed mainly of gastric cancers, and GATA4 in the TCGA pancreatic cohort using the cBioportal
721 web interface tool. We also validated these markers as independent predictors of survival both in
722 respect of each other and stage using a multivariate Cox regression in our 551 case cohort. When
723 assessing for genomic correlates with differentiation phenotypes we found only very few cases with
724 well differentiated phenotypes (<5% cases) and hence for statistical analyses we collapse these cases
725 with moderate differentiation to allow a binary fisher's exact test to compare poorly differentiated
726 with well-moderate differentiated phenotypes.

727

728 **Therapeutics**

729 The cancer biomarker database was filtered for drugs linked to biomarkers found in EAC drivers and
730 supplementary table 6 constructed using the cohort frequencies of EAC biomarkers. 10 EAC cell lines
731 (SKGT4, OACP4C, OACM5.1, ESO26, ESO51, OE33, MFD, OE19, Flo-1 and JHesoAD) and 3 BE high
732 grade dysplasia cell lines (CP-B, CP-C and CP-D) with WGS data⁴⁶ were used in proliferation assays to
733 determine drug sensitivity to CDK4/6 inhibitors, Palbociclib (Biovision) and Ribociclib (Selleckchem).
734 Cell lines were grown in their normal growth media (methods table 1). Proliferation was measured
735 using the Incucyte live cell analysis system (Incucyte ZOOM Essen biosciences). Each cell line was
736 plated at a starting confluency of 10% and growth rate measured across 4-7 days depending on basal
737 proliferation rate. For each cell-line drug combination concentrations of 16, 64, 250, 1000 and 4000
738 nanomolar were used each in 0.3% DMSO and compared to 0.3% DMSO only. Each condition was
739 performed in at least triplicate. The time period of the exponential growth phase in the untreated
740 (0.3% DMSO) condition was used to calculate GI50 and AUC. Accurate GI50s could not be calculated
741 in cases where a cell line had >50% proliferation inhibition even with the highest drug concentration
742 and hence AUC was used to compare cell line sensitivity. T47D had a highly similar GI50 for

743 Palbociclib to that previously calculated in other studies (112 nM vs 127 nM)⁶¹. Primary organoid
744 cultures were derived from EAC cases included in the OCCAMS/ICGC sequencing study. Detailed
745 organoid culture and derivation method have been previously described (cite nat comms Li et al).
746 Regarding the drug treatment, the seeding density for each line was optimised to ensure cell growth
747 in the logarithmic growth phase. Cells were seeded in complete medium for 24 hours then treated
748 with compounds at a 5-point 4-fold serial dilutions for 6 days or 12 days. Cell viability was assessed
749 using CellTiter-Glo (Promega) after drug incubation.

750 References

- 751 1. Ferlay J, Soerjomataram I, Dikshit R, et al. Cancer incidence and mortality worldwide:
752 sources, methods and major patterns in GLOBOCAN 2012. *Int J Cancer*. Mar 1
753 2015;136(5):E359-386.
- 754 2. Coleman HG, Xie SH, Lagergren J. The Epidemiology of Esophageal Adenocarcinoma.
755 *Gastroenterology*. Jan 2018;154(2):390-405.
- 756 3. Smyth EC, Lagergren J, Fitzgerald RC, et al. Oesophageal cancer. *Nat Rev Dis Primers*. Jul 27
757 2017;3:17048.
- 758 4. Campbell PJ, Getz G, Stuart JM, Korbel JO, Stein LD. Pan-cancer analysis of whole genomes.
759 *bioRxiv*. 2017.
- 760 5. Ciriello G, Miller ML, Aksoy BA, Senbabaoglu Y, Schultz N, Sander C. Emerging landscape of
761 oncogenic signatures across human cancers. *Nat Genet*. Oct 2013;45(10):1127-1133.
- 762 6. Secrier M, Li X, de Silva N, et al. Mutational signatures in esophageal adenocarcinoma define
763 etiologically distinct subgroups with therapeutic relevance. *Nat Genet*. Oct
764 2016;48(10):1131-1141.
- 765 7. Stratton MR, Futreal PA. Cancer: understanding the target. *Nature*. Jul 1 2004;430(6995):30.
- 766 8. Tamborero D, Gonzalez-Perez A, Perez-Llamas C, et al. Comprehensive identification of
767 mutational cancer driver genes across 12 tumor types. *Sci Rep*. Oct 2 2013;3:2650.
- 768 9. Lawrence MS, Stojanov P, Polak P, et al. Mutational heterogeneity in cancer and the search
769 for new cancer-associated genes. *Nature*. Jul 11 2013;499(7457):214-218.
- 770 10. Integrated genomic characterization of oesophageal carcinoma. *Nature*. Jan 12
771 2017;541(7636):169-175.
- 772 11. Lin DC, Dinh HQ, Xie JJ, et al. Identification of distinct mutational patterns and new driver
773 genes in oesophageal squamous cell carcinomas and adenocarcinomas. *Gut*. Aug 31 2017.
- 774 12. Rheinbay E, Nielsen MM, Abascal F, et al. Discovery and characterization of coding and non-
775 coding driver mutations in more than 2,500 whole cancer genomes. *bioRxiv*. 2017.
- 776 13. Comprehensive molecular characterization of urothelial bladder carcinoma. *Nature*. Mar 20
777 2014;507(7492):315-322.
- 778 14. Comprehensive molecular characterization of gastric adenocarcinoma. *Nature*. Sep 11
779 2014;513(7517):202-209.
- 780 15. Mermel CH, Schumacher SE, Hill B, Meyerson ML, Beroukhi R, Getz G. GISTIC2.0 facilitates
781 sensitive and confident localization of the targets of focal somatic copy-number alteration in
782 human cancers. *Genome Biol*. 2011;12(4):R41.

- 783 16. Dulak AM, Schumacher SE, van Lieshout J, et al. Gastrointestinal adenocarcinomas of the
784 esophagus, stomach, and colon exhibit distinct patterns of genome instability and
785 oncogenesis. *Cancer Res.* Sep 1 2012;72(17):4383-4393.
- 786 17. Frankel A, Armour N, Nancarrow D, et al. Genome-wide analysis of esophageal
787 adenocarcinoma yields specific copy number aberrations that correlate with prognosis.
788 *Genes Chromosomes Cancer.* Apr 2014;53(4):324-338.
- 789 18. Secrier M, Fitzgerald RC. Signatures of Mutational Processes and Associated Risk Factors in
790 Esophageal Squamous Cell Carcinoma: A Geographically Independent Stratification Strategy?
791 *Gastroenterology.* May 2016;150(5):1080-1083.
- 792 19. Zack TI, Schumacher SE, Carter SL, et al. Pan-cancer patterns of somatic copy number
793 alteration. *Nat Genet.* Oct 2013;45(10):1134-1140.
- 794 20. Dulak AM, Stojanov P, Peng S, et al. Exome and whole-genome sequencing of esophageal
795 adenocarcinoma identifies recurrent driver events and mutational complexity. *Nat Genet.*
796 May 2013;45(5):478-486.
- 797 21. Nones K, Waddell N, Wayte N, et al. Genomic catastrophes frequently arise in esophageal
798 adenocarcinoma and drive tumorigenesis. *Nat Commun.* Oct 29 2014;5:5224.
- 799 22. Martincorena I, Raine KM, Gerstung M, et al. Universal Patterns of Selection in Cancer and
800 Somatic Tissues. *Cell.* Nov 16 2017;171(5):1029-1041 e1021.
- 801 23. Wadi L, Uuskula-Reimand L, Isaev K, et al. Candidate cancer driver mutations in super-
802 enhancers and long-range chromatin interaction networks. *bioRxiv.* 2017.
- 803 24. Gonzalez-Perez A, Lopez-Bigas N. Functional impact bias reveals cancer drivers. *Nucleic Acids*
804 *Res.* Nov 2012;40(21):e169.
- 805 25. Tamborero D, Gonzalez-Perez A, Lopez-Bigas N. OncodriveCLUST: exploiting the positional
806 clustering of somatic mutations to identify cancer genes. *Bioinformatics.* Sep 15
807 2013;29(18):2238-2244.
- 808 26. Porta-Pardo E, Godzik A. e-Driver: a novel method to identify protein regions driving cancer.
809 *Bioinformatics.* Nov 1 2014;30(21):3109-3114.
- 810 27. Porta-Pardo E, Hrabe T, Godzik A. Cancer3D: understanding cancer mutations through
811 protein structures. *Nucleic Acids Res.* Jan 2015;43(Database issue):D968-973.
- 812 28. Futreal PA, Coin L, Marshall M, et al. A census of human cancer genes. *Nat Rev Cancer.* Mar
813 2004;4(3):177-183.
- 814 29. Kandoth C, McLellan MD, Vandin F, et al. Mutational landscape and significance across 12
815 major cancer types. *Nature.* Oct 17 2013;502(7471):333-339.
- 816 30. Shuai S, Gallinger S, Stein LD. DriverPower: Combined burden and functional impact tests for
817 cancer driver discovery. *bioRxiv.* 2017.
- 818 31. Taylor AM, Shih J, Ha G, et al. Genomic and Functional Approaches to Understanding Cancer
819 Aneuploidy. *Cancer cell.* Apr 9 2018;33(4):676-689 e673.
- 820 32. Chang MT, Asthana S, Gao SP, et al. Identifying recurrent mutations in cancer reveals
821 widespread lineage diversity and mutational specificity. *Nat Biotechnol.* Feb 2016;34(2):155-
822 163.
- 823 33. Zaretsky JM, Garcia-Diaz A, Shin DS, et al. Mutations Associated with Acquired Resistance to
824 PD-1 Blockade in Melanoma. *N Engl J Med.* Sep 1 2016;375(9):819-829.
- 825 34. Chen Z, Shi T, Zhang L, et al. Mammalian drug efflux transporters of the ATP binding cassette
826 (ABC) family in multidrug resistance: A review of the past decade. *Cancer Lett.* Jan 1
827 2016;370(1):153-164.
- 828 35. Giannakis M, Mu XJ, Shukla SA, et al. Genomic Correlates of Immune-Cell Infiltrates in
829 Colorectal Carcinoma. *Cell reports.* Oct 18 2016;17(4):1206.
- 830 36. Pei XH, Xiong Y. Biochemical and cellular mechanisms of mammalian CDK inhibitors: a few
831 unresolved issues. *Oncogene.* Apr 18 2005;24(17):2787-2795.
- 832 37. Comprehensive molecular characterization of human colon and rectal cancer. *Nature.* Jul 18
833 2012;487(7407):330-337.

- 834 38. Waddell N, Pajic M, Patch AM, et al. Whole genomes redefine the mutational landscape of
835 pancreatic cancer. *Nature*. Feb 26 2015;518(7540):495-501.
- 836 39. Leiserson MD, Vandin F, Wu HT, et al. Pan-cancer network analysis identifies combinations
837 of rare somatic mutations across pathways and protein complexes. *Nat Genet*. Feb
838 2015;47(2):106-114.
- 839 40. Grasso CS, Giannakis M, Wells DK, et al. Genetic Mechanisms of Immune Evasion in
840 Colorectal Cancer. *Cancer discovery*. Jun 2018;8(6):730-749.
- 841 41. Singhi AD, Foxwell TJ, Nason K, et al. Smad4 loss in esophageal adenocarcinoma is associated
842 with an increased propensity for disease recurrence and poor survival. *Am J Surg Pathol*. Apr
843 2015;39(4):487-495.
- 844 42. Levy L, Hill CS. Alterations in components of the TGF-beta superfamily signaling pathways in
845 human cancer. *Cytokine Growth Factor Rev*. Feb-Apr 2006;17(1-2):41-58.
- 846 43. Tamborero D, Rubio-Perez C, Deu-Pons J, et al. Cancer Genome Interpreter Annotates The
847 Biological And Clinical Relevance Of Tumor Alterations. *bioRxiv*. 2017.
- 848 44. Weaver JMJ, Ross-Innes CS, Shannon N, et al. Ordering of mutations in preinvasive disease
849 stages of esophageal carcinogenesis. *Nature genetics*. Aug 2014;46(8):837-843.
- 850 45. Ross-Innes CS, Becq J, Warren A, et al. Whole-genome sequencing provides new insights into
851 the clonal architecture of Barrett's esophagus and esophageal adenocarcinoma. *Nature*
852 *genetics*. Sep 2015;47(9):1038-1046.
- 853 46. Contino G, Eldridge MD, Secrier M, et al. Whole-genome sequencing of nine esophageal
854 adenocarcinoma cell lines. *F1000Res*. 2016;5:1336.
- 855 47. Liston DR, Davis M. Clinically Relevant Concentrations of Anticancer Drugs: A Guide for
856 Nonclinical Studies. *Clin Cancer Res*. Jul 15 2017;23(14):3489-3498.
- 857 48. Herrera-Abreu MT, Palafox M, Asghar U, et al. Early Adaptation and Acquired Resistance to
858 CDK4/6 Inhibition in Estrogen Receptor-Positive Breast Cancer. *Cancer Res*. Apr 15
859 2016;76(8):2301-2313.
- 860 49. Li X, Francies HE, Secrier M, et al. Organoid cultures recapitulate esophageal
861 adenocarcinoma heterogeneity providing a model for clonality studies and precision
862 therapeutics. *Nature communications*. Jul 30 2018;9(1):2983.
- 863 50. Patnaik A, Rosen LS, Tolaney SM, et al. Efficacy and Safety of Abemaciclib, an Inhibitor of
864 CDK4 and CDK6, for Patients with Breast Cancer, Non-Small Cell Lung Cancer, and Other
865 Solid Tumors. *Cancer discovery*. Jul 2016;6(7):740-753.
- 866 51. D'Antonio M, Ciccarelli FD. Integrated analysis of recurrent properties of cancer genes to
867 identify novel drivers. *Genome Biol*. May 29 2013;14(5):R52.
- 868 52. Zehir A, Benayed R, Shah RH, et al. Mutational landscape of metastatic cancer revealed from
869 prospective clinical sequencing of 10,000 patients. *Nat Med*. Jun 2017;23(6):703-713.
- 870 53. Ding J, McConechy MK, Horlings HM, et al. Systematic analysis of somatic mutations
871 impacting gene expression in 12 tumour types. *Nat Commun*. Oct 5 2015;6:8554.
- 872 54. Nagai K, Kohno K, Chiba M, et al. Differential expression profiles of sense and antisense
873 transcripts between HCV-associated hepatocellular carcinoma and corresponding non-
874 cancerous liver tissue. *Int J Oncol*. Jun 2012;40(6):1813-1820.
- 875 55. Adzhubei I, Jordan DM, Sunyaev SR. Predicting functional effect of human missense
876 mutations using PolyPhen-2. *Curr Protoc Hum Genet*. Jan 2013;Chapter 7:Unit7 20.
- 877 56. Ng PC, Henikoff S. Predicting the effects of amino acid substitutions on protein function.
878 *Annu Rev Genomics Hum Genet*. 2006;7:61-80.
- 879 57. Reimand J, Wagih O, Bader GD. The mutational landscape of phosphorylation signaling in
880 cancer. *Sci Rep*. Oct 2 2013;3:2651.
- 881 58. Northcott PA, Buchhalter I, Morrissy AS, et al. The whole-genome landscape of
882 medulloblastoma subtypes. *Nature*. Jul 19 2017;547(7663):311-317.
- 883 59. Wala JA, Shapira O, Li Y, et al. Selective and mechanistic sources of recurrent
884 rearrangements across the cancer genome. *bioRxiv*. 2017.

- 885 60. Gao J, Aksoy BA, Dogrusoz U, et al. Integrative analysis of complex cancer genomics and
886 clinical profiles using the cBioPortal. *Sci Signal*. Apr 2 2013;6(269):p11.
- 887 61. Finn RS, Dering J, Conklin D, et al. PD 0332991, a selective cyclin D kinase 4/6 inhibitor,
888 preferentially inhibits proliferation of luminal estrogen receptor-positive human breast
889 cancer cell lines in vitro. *Breast Cancer Res*. 2009;11(5):R77.
- 890

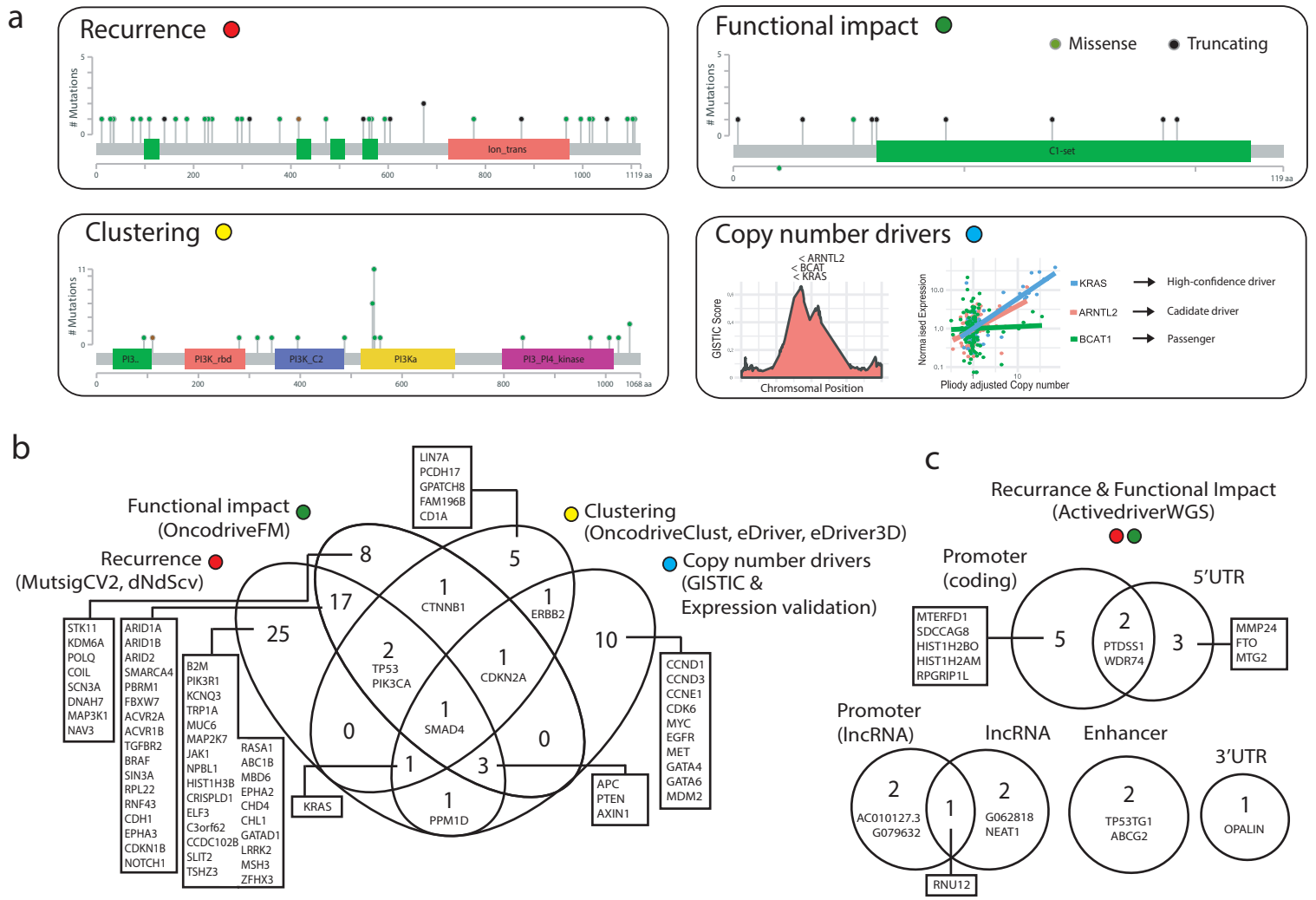


Figure 1 Detection of EAC Driver Genes. **a.** Types of driver-associated features used to detect positive selection in mutations and copy number events with examples of genes containing such features **b.** Coding driver genes identified and their driver-associated features. **c.** Non-coding driver elements detected and their element types.

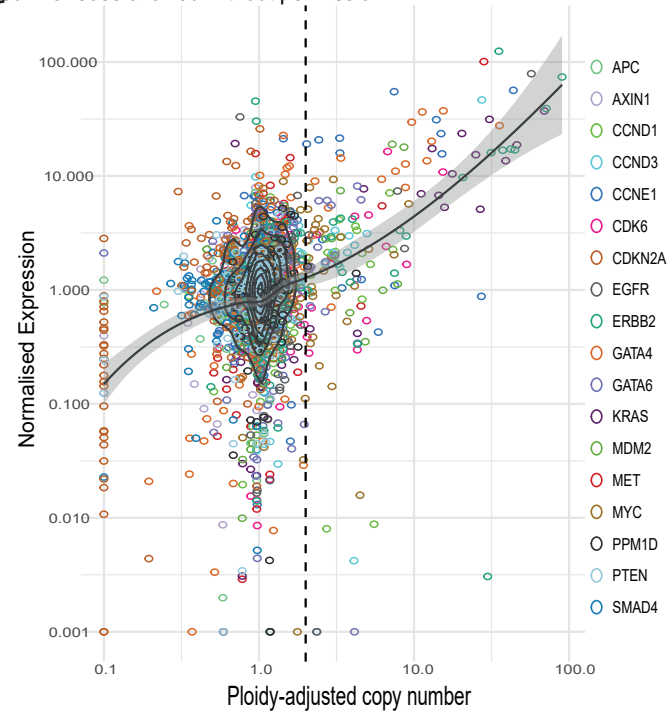
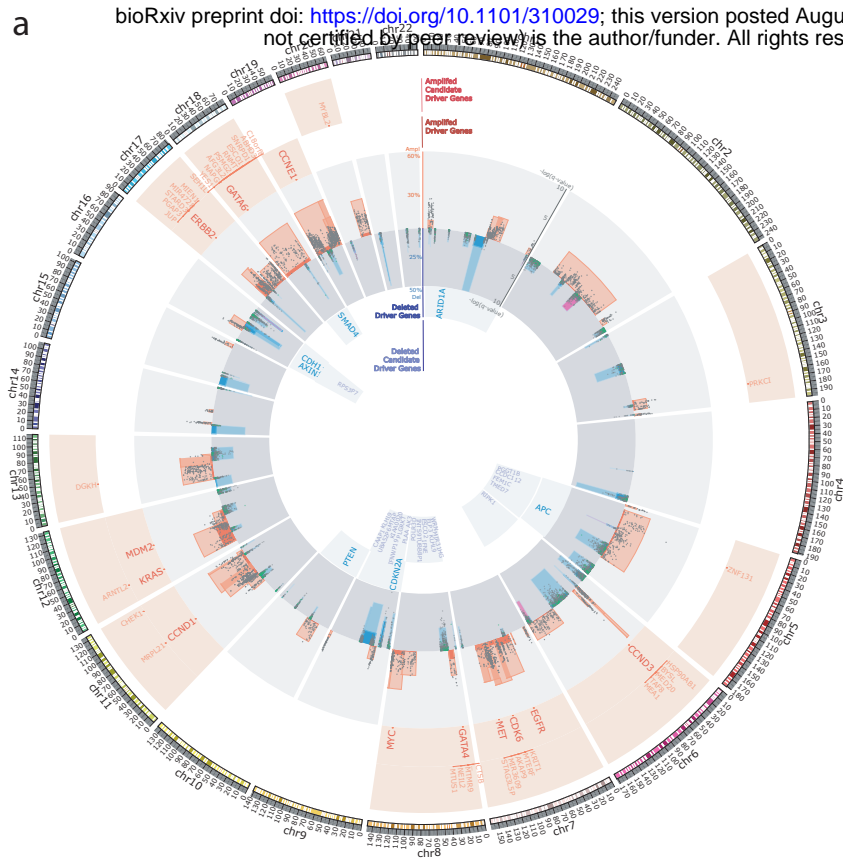
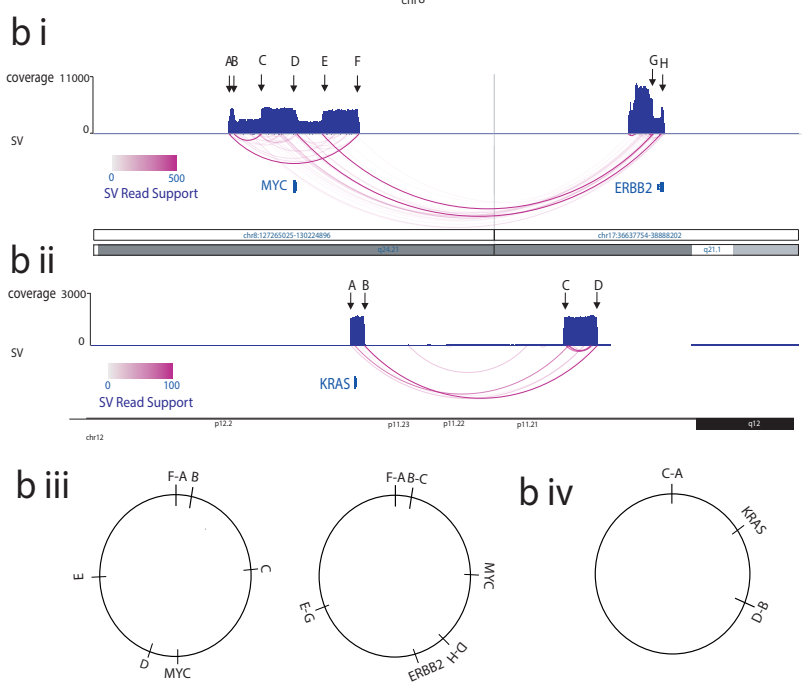


Figure 2. Copy number variation under positive selection.



a. Recurrent copy number changes across the genome identified by GISTIC. Frequency of different CNV types are indicated as well as the position of CNV high confidence driver genes and candidate driver genes. The q value for expression correlation with amplification and homozygous deletion is shown for each gene within each amplification and deletion peaks respectively and occasions of significant association between LOH and mutation are indicated in green. Purple deletion peaks indicate fragile sites. **b.** Examples of extrachromosomal-like amplifications suggested by very high read support SVs at the boundaries of highly amplified regions produced from a single copy number step. In the first example (bi) two populations of extrachromosomal DNA are apparent (biii), one amplifying only MYC and the second also incorporating ERBB2 from a different chromosome. In the second example (bii) an inversion has occurred before circularization and amplification around KRAS (biv). **c.** Relationship between copy number and expression in CN driver genes.

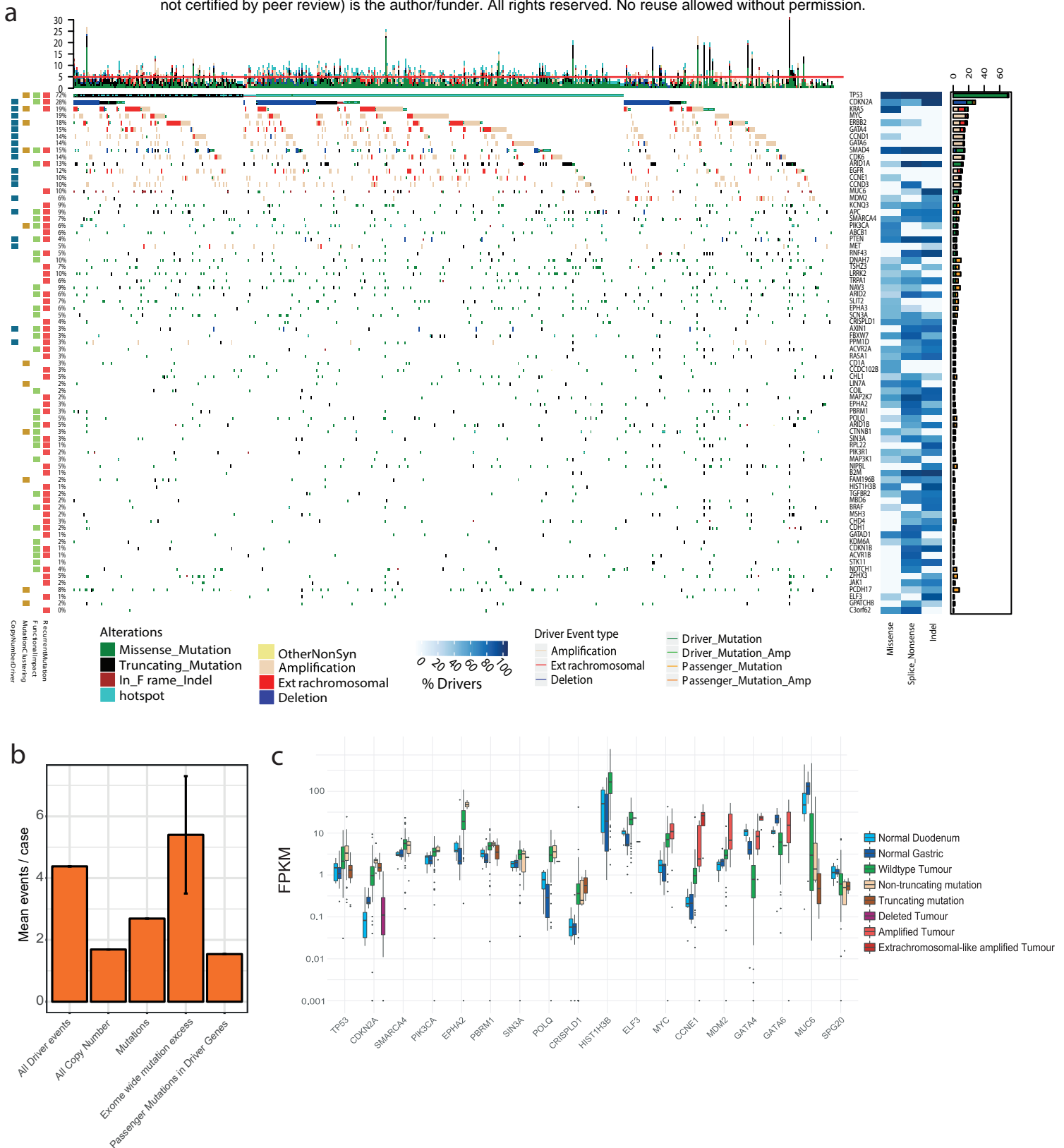


Figure 3. The driver gene landscape of Esophageal Adenocarcinoma. a. Driver mutations or CNVs are shown for each patient. Amplification is defined as >2 Copy number adjusted ploidy ($2 \times$ ploidy of that case) and extrachromosomal amplification as >10 Copy number adjusted ploidy ($10 \times$ ploidy for that case). Driver associated features for each driver gene are displayed to the left. On the right the percentages of different mutation and copy number changes are displayed, differentiating between driver and passenger mutations using dNdScv, and the % of predicted drivers by mutation type is shown. Above the plot are the number of driver mutations per sample with an indication of the median (red line = 5). **b.** Assessment of driver event types per case and comparison to exome-wide excess of mutations generated by dNdScv. **c.** Expression changes in EAC driver genes in comparison to normal intestinal tissues. Only genes with significant expression changes of note are shown.

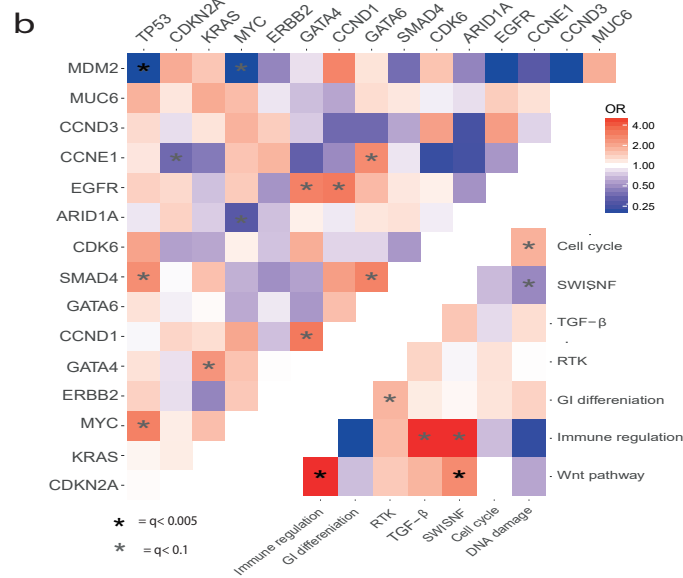
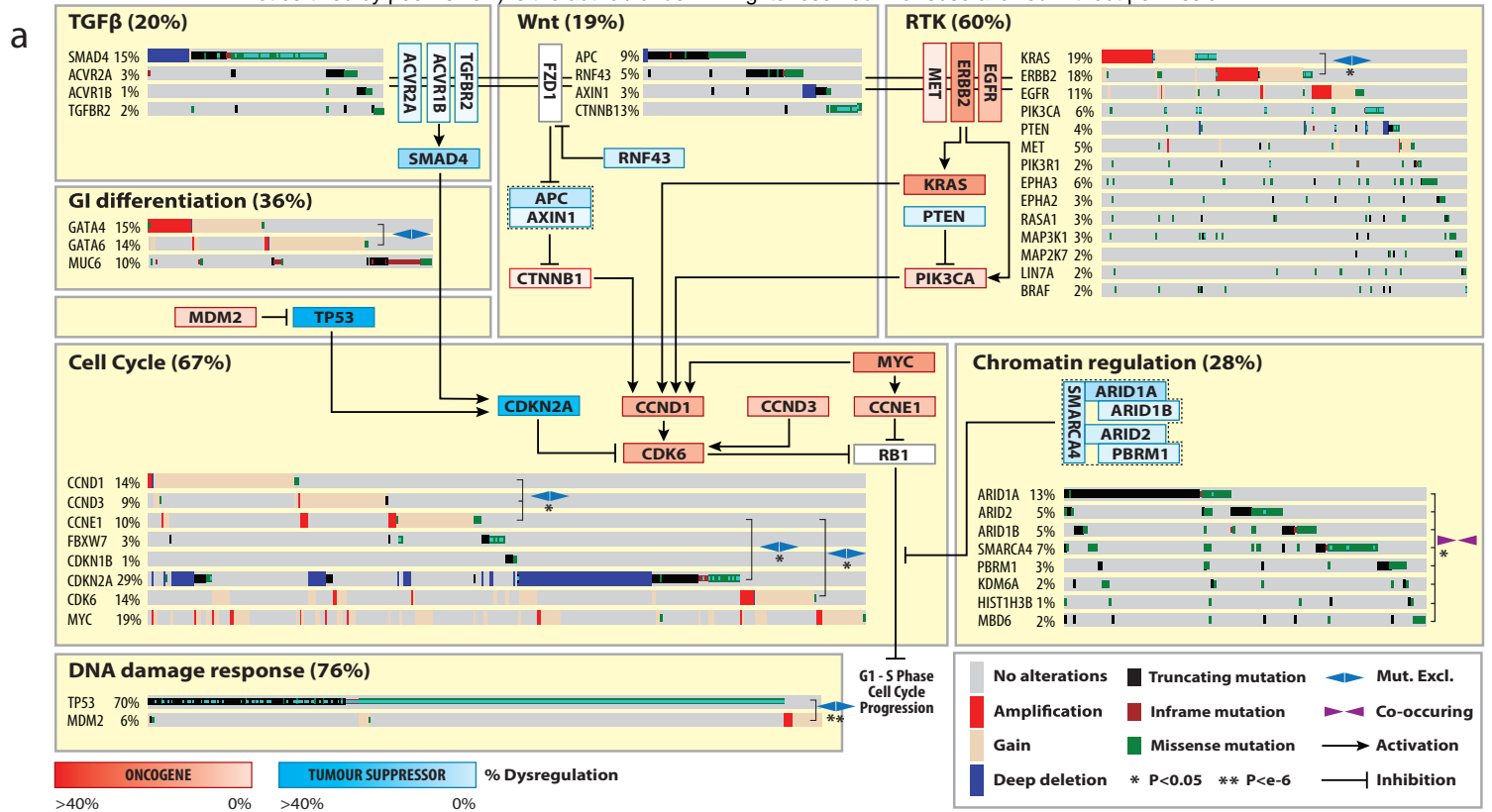


Figure 4. Biological pathways undergoing selective dysregulation in EAC. a. Biological Pathways dysregulated by driver gene mutation and/or CNVs. WT cases for a pathway are not shown. Mutual exclusivities and/or associations between genes in a pathway are annotated. GATA4/6 amplifications have a mutually exclusive relationship (ie GATA4 amplification is more common in GATA6 WT cases) although this does not reach statistical significance (fisher's exact test $p=0.07$ OR =0.52). **b.** Pairwise assessment of mutual exclusivity and association in EAC driver genes and pathways.

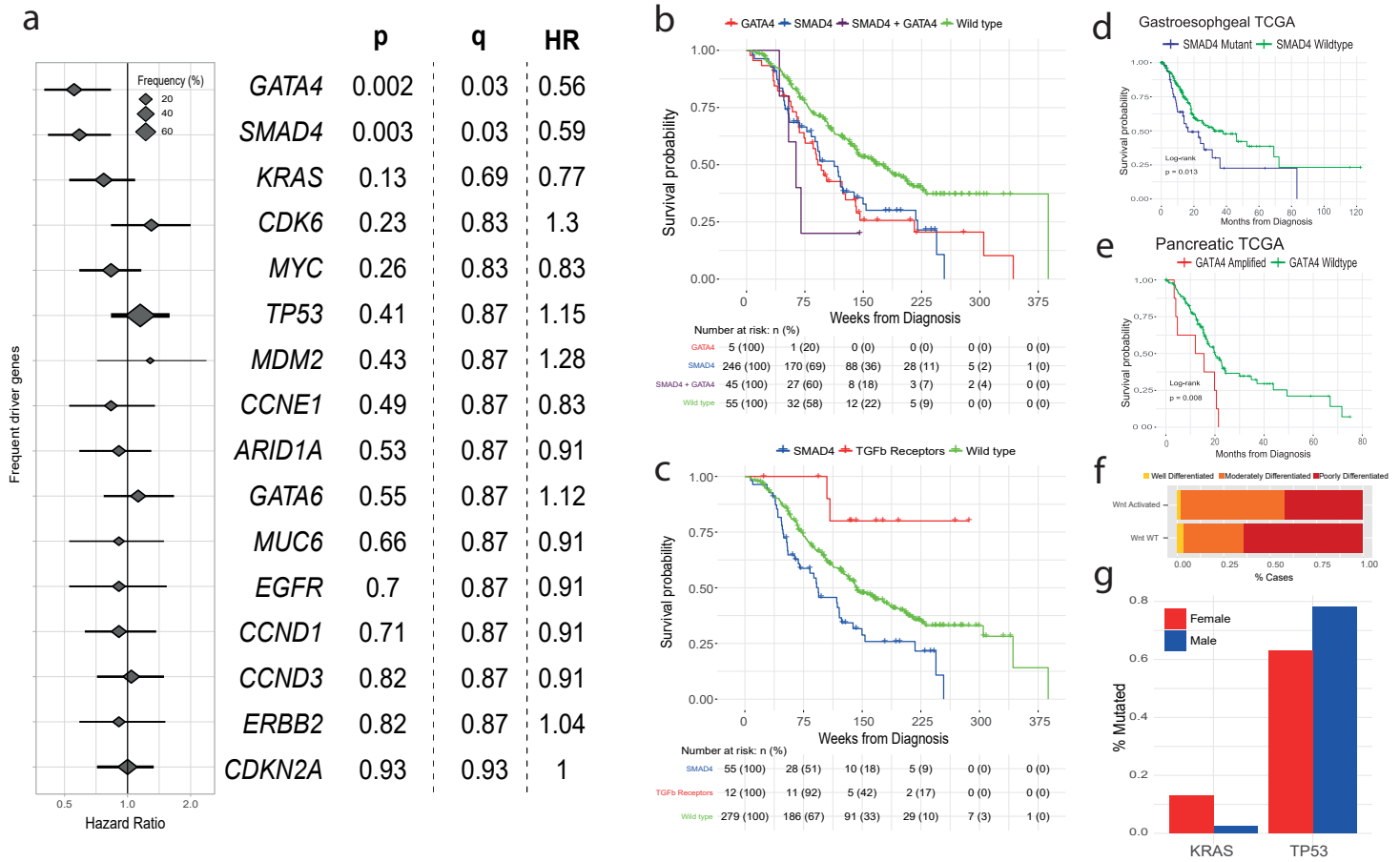


Figure 5. Clinical significance of Driver events in EAC. **a.** Hazard ratios and 95% confidence intervals for Cox regression analysis across all driver genes with at least a 5% frequency of driver alterations. P values are generated from the Wald test and q values generated using BH correction. **b.** Kaplan-Meier curves for EACs with different status of significant prognostic indicators (GATA4 and SMAD4). **c.** Kaplan-Meier curves for different alterations in the TGFbeta pathway. **d.** Kaplan-Meier curves showing verification GATA4 prognostic value in GI cancers using a pancreatic TCGA cohort. **e.** Kaplan-Meier curves showing verification SMAD4 prognostic value in gastroesophageal cancers using a gastroesophageal TCGA cohort. **f.** Differentiation bias in tumours containing events in Wnt pathway driver genes. **g.** Relative frequency of KRAS mutations and TP53 mutations driver gene events in females vs males (fisher's exact test).

