

1 Low-cost and clinically applicable copy number profiling 2 using repeat DNA

3 Abujudeh S^{1,†}, Zeki SS^{2,3,†}, van Lanschot MCV², Pusung M², Weaver JMJ^{2,4}, Li X², Noorani
4 A², Metz AJ², Bornschein J², Bower L¹, Miremadi A², Fitzgerald RC^{2,‡}, Morrissey ER^{1,5,‡},
5 Lynch AG^{1,6,‡}, and the Oesophageal Cancer Clinical and Molecular Stratification (OCCAMS)
6 Consortium⁷

7 ¹Cancer Research UK Cambridge Institute, University of Cambridge, Li Ka Shing Centre,
8 Robinson Way, Cambridge CB2 0RE, UK

9 ²Medical Research Council (MRC) Cancer Unit, University of Cambridge, Cambridge, UK

10 ³Gastroenterology, Guy's and St Thomas' Hospital, London, UK

11 ⁴Department of Medical Oncology, The Christie NHS Foundation Trust, Manchester, M20
12 4TX

13 ⁵Weatherall Institute of Molecular Medicine, University of Oxford, Oxford, UK

14 ⁶School of Mathematics and Statistics/School of Medicine, University of St Andrews, St
15 Andrews, UK

16 ⁷A full list of contributors from the OCCAMS Consortium is available at the end of the
17 manuscript

18 [†]These authors contributed equally to this work

19 [‡]These authors contributed equally to this work

20 August 18, 2018

21 Abstract

22 Large-scale cancer genome studies suggest that tumors are driven by somatic copy number alterations
23 (SCNAs) or single-nucleotide variants (SNVs). Due to the low-cost, the clinical use of genomics assays is
24 biased towards targeted gene panels, which identify SNVs. There is a need for a comparably low-cost and
25 simple assay for high-resolution SCNA profiling. Here we present our method, conliga, which infers SCNA
26 profiles from a low-cost and simple assay.

27

28 Somatic copy number alterations (SCNAs) are common in cancer. On average, cancer samples see SCNAs in
29 34% of the genome, with 17% of the genome amplified and 16% deleted [1, 2, 3]. Certain SCNAs, particularly
30 amplifications of oncogenes and deletions of tumor suppressor genes, have been found to be major drivers in
31 tumor development, associated with prognosis and response to therapy [1]. SCNA burden varies considerably
32 between cancer types [3]. For example, oesophageal adenocarcinoma (OAC) has relatively high levels of SCNAs
33 [4, 5, 6, 7], and generally develops from Barrett's oesophagus. Patients with OAC tend to be diagnosed at a late
34 stage, when spread has occurred to lymph nodes and distant organs. This makes treatment more difficult and
35 leads to poor prognosis [8]. Although most patients with Barrett's do not progress, early stage disease (high
36 grade dysplasia or intramucosal adenocarcinoma) can be successfully treated, usually obviating the need for
37 surgery. There is a critical need to develop technologies that can detect early disease and distinguish between
38 patients at low versus high risk for progression. Since most mutations in OAC driver genes are already present in
39 pre-malignant disease [9], but an increased SCNA load distinguishes OAC [10, 11, 12], low-cost SCNA profiling
40 would be a valuable research and clinical tool.

41 SCNAs have been identified using a number of methods, including comparative genomic hybridisation (CGH)
42 [13], array-based CGH [14], single nucleotide polymorphism (SNP) arrays [15], and whole-genome sequencing
43 (WGS) [16]. Recently, low-coverage (LC) WGS has gained popularity due to its reduced cost and strong
44 performance [17]. However, while LC WGS reduces the cost of sequencing, standard WGS library preparation
45 is required with its associated fixed expense and time needed to produce each sample. A technically simple,
46 fast, easily automated, high-resolution and inexpensive alternative method for SCNA detection, with clinical
47 potential, would be extremely valuable.

48 Recent studies have shown the genome can be amplified at multiple (>10,000) genomic loci with the use of a
49 single non-specific primer pair, using the FAST-SeqS method [18, 19]. With this approach, two polymerase chain
50 reaction (PCR) rounds replace the complicated and expensive library preparation steps associated with WGS.
51 The amplified regions are sufficiently short such that the assay can be performed on cell-free DNA as well as
52 DNA extracted from tissue biopsies. The resulting amplicons can be sequenced, with samples multiplexed on the
53 same sequencing lane. With this method, we maintain a similar sequencing depth to 30-50X high-coverage (HC)
54 WGS while sequencing only specific loci. This is in contrast to LC WGS which samples the whole genome but
55 at reduced sequencing depth (Supplementary Fig. 1). The cost involved in sample preparation and sequencing
56 combined is approximately £14 per sample compared with approximately £52-72 for LC WGS, depending on
57 the library preparation kit used (Supplementary Note 1). The sample preparation can be performed in less
58 than an hour with minimal hands-on time, compared to approximately 3 hours or greater for LC WGS.

59 Until now, the use of FAST-SeqS data has been limited to the detection of whole chromosome gains [18] and
60 entire chromosome arm gains and losses [19, 20]. This means that chromosome segment (focal) alterations are
61 not detected, or perhaps falsely considered as whole chromosome or chromosome arm alterations. Moreover, in

62 these methods SCNAs are not quantified and regions are simply classified as amplified, deleted or normal.

63 Here we present a method (and associated tool: ‘conliga’) that uses a fully probabilistic approach to infer
64 relative copy number (RCN) alterations at each locus from FAST-SeqS data. conliga provides a RCN profile
65 per sample and therefore enables this low-cost sequencing approach to be used as a SCNA assay.

66 Based on observations of raw data (Supplementary Note 2, Supplementary Fig. 1), we created a probabilistic
67 model (Methods, Supplementary Note 3). The model takes account of the observed bias in loci counts, which
68 predominantly results from unequal PCR efficiencies between loci. Since neighboring loci are likely to share the
69 same copy number, we use a hidden Markov model (HMM) to model the spatial dependence between loci. This
70 allows loci with high counts to share statistical strength with neighboring loci, enabling us to infer contiguous
71 regions of copy number more accurately. Moreover, we use a Bayesian nonparametric approach (sticky HDP-
72 HMM) [21] to address the issue of the unknown number of copy number levels present in a given sample a
73 priori (Methods). We use Markov Chain Monte Carlo (MCMC) methods to infer the RCN of each locus, plus
74 all other latent variables in the model (Methods, Supplementary Table 1, Supplementary Notes 4, 5 and 6).
75 This enables us to provide the uncertainty of the RCN estimates, summarized by credible intervals, in conliga’s
76 standard output.

77 To test our method, we analysed 11 oesophageal adenocarcinoma tumors (Methods, Supplementary Tables
78 2 and 3), which had been sequenced using HC WGS (>50X) and FAST-SeqS. In addition, we downsampled
79 the WGS data of each sample to nine million reads to simulate typical LC WGS (~ 0.1X coverage) samples
80 (Methods). We compared the copy number calls derived from ASCAT [22] (applied to HC WGS data) with the
81 RCN calls from QDNAseq [17] (LC WGS data) and conliga (FAST-SeqS data). conliga and QDNAseq achieved
82 a median Pearson correlation coefficient with ASCAT of 0.95 and 0.98 respectively (Methods, Supplementary
83 Table 4).

84 In figure 1a-d we demonstrate that similar RCN profiles are obtained with the three methods for an example
85 sample (OAC2) and that high-resolution SCNA information is maintained by sampling genomic loci using
86 FAST-SeqS. Figure 1e and 1f show the performance of conliga and QDNAseq, both obtaining similar Pearson
87 correlation coefficients with ASCAT’s RCN calls across all 11 OAC samples (conliga: 0.953, QDNAseq: 0.987)
88 and residual distributions when compared to ASCAT (Methods). It should be noted that by downsampling
89 reads from the same WGS sample, this analysis is potentially biased in favor of QDNAseq’s results.

90 From the literature [23, 12] we selected a set of 36 genes that have been observed to be recurrently amplified
91 or deleted in OAC (Supplementary Table 5, Methods). We determined the weighted mean of the RCN calls
92 for these genes for each sample via each method (Methods, Supplementary Tables 6 and 7). While FAST-
93 SeqS/conliga would not be the assay of choice if only interested in a small gene panel, in Figure 1g we see
94 that there are only two instances from 396 comparisons (36 genes x 11 samples) where a substantially different
95 result would be achieved. Naturally if an SCNA is so narrow as to fall between two FAST-SeqS loci then it will
96 not be detected in this way, but the detection of many highly-localized events demonstrates how informative
97 FAST-SeqS/conliga can be. Even within this panel of 36, it is notable that some genes harbour FAST-SeqS

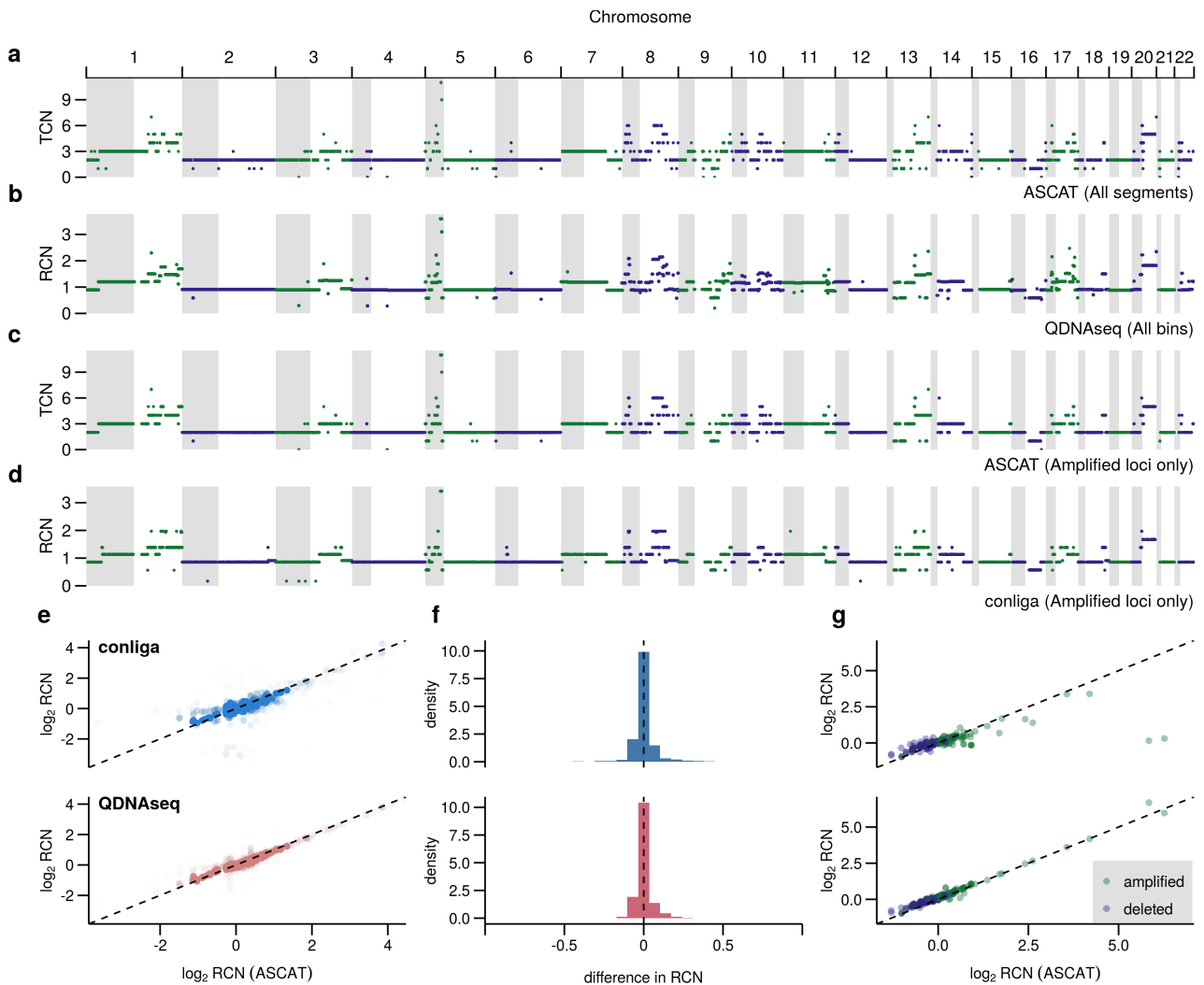


Figure 1: Comparison of conliga method with ASCAT and QDNaseq. (a) Total copy number profile determined by ASCAT from HC WGS data for sample OAC2, showing all copy number segments. (b) Relative copy number profile determined by QDNaseq from LC WGS data for sample OAC2, showing all 15 Kbp bins. (c) Total copy number profile determined by ASCAT from HC WGS data for sample OAC2, showing ASCAT's copy number calls at the intersection of ASCAT's called regions and FAST-SeqS loci. (d) Relative copy number profile determined by conliga from FAST-SeqS data for sample OAC2, at the intersection of ASCAT's called regions and FAST-SeqS loci. (e) Comparison of \log_2 relative copy number calls from 11 samples between conliga and ASCAT (top) and QDNaseq and ASCAT (bottom). All RCN calls at the intersection of ASCAT's called regions, QDNaseq 15Kb bins and FAST-SeqS loci in all 11 OAC samples are shown as points. (f) Distribution of differences between ASCAT RCN calls and conliga RCN estimates for 11 OAC samples (top) and ASCAT RCN calls and QDNaseq RCN estimates for 11 OAC samples (bottom). (g) Comparison of performance at gene level resolution between ASCAT and conliga (top) and ASCAT and QDNaseq (bottom). The values represent the weighted mean of RCN calls at each gene for each of the 11 OAC samples (Methods).

98 loci (Supplementary Tables 8 and 9), providing evidence of intra gene SCNAs in some cases, such as the focal
 99 deletions observed in FHIT, PARK2, and MACROD2 (Supplementary Fig. 2). Focal deletions such as these
 100 may be functionally relevant, potentially rendering tumor suppressor genes inactive.

101 The purity of tumor samples obtained by dissection can vary widely [24], as can samples obtained non-
 102 invasively, e.g ctDNA from plasma [25]. As tumor purity reduces, the copy number signal to noise ratio decreases.
 103 To determine the performance of conliga and QDNaseq under different purity conditions, we generated samples

104 with varying purity by mixing sequencing reads from normal and OAC samples (Methods). FAST-SeqS samples
 105 were generated with two million reads and LC WGS samples were generated with nine million reads.

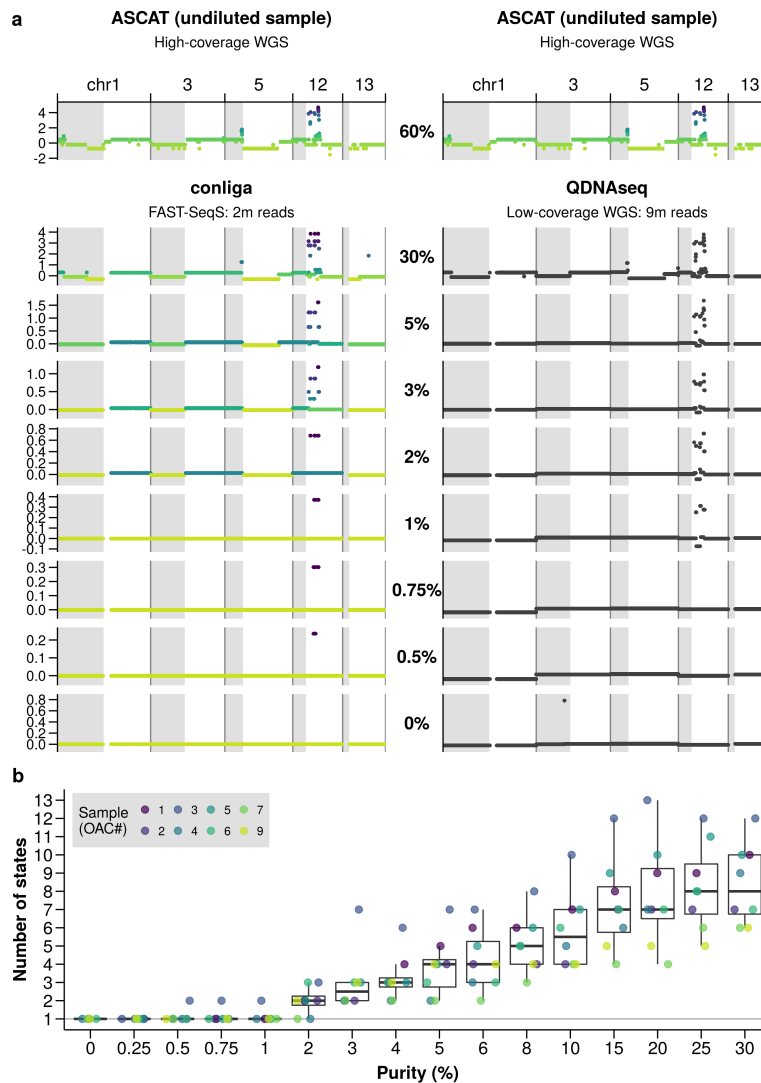


Figure 2: Comparing the performance of SCNA detection in low tumor purity samples and determining the limit of detection. (a) left column: relative copy number calls by conliga at different dilutions of sample OAC3, compared to ASCAT relative copy number profile (top left), discrete copy number states are colored with a gradient (light green to purple), highlighting regions with differing SCNAs. right column: relative copy number calls by QDNaseq at different dilutions of sample OAC3, compared to ASCAT relative copy number profile (top right). (b) The number of copy number states detected by conliga in each of eight OAC samples at differing purity levels. The limit of detection is determined by the lowest purity level in which more than one copy number state is detected.

106 Figure 2a shows the performance of both methods for sample OAC3. At 30% purity, both conliga and
 107 QDNaseq recapitulate the copy number profile as determined by ASCAT. At 5%, other than the focal amplifi-
 108 cation on chromosome 12, QDNaseq fails to detect sub chromosomal SCNAs, whereas conliga shows evidence
 109 of chromosome arm and sub-chromosomal arm changes. At 2% purity, conliga is able to distinguish some of the
 110 more prominent chromosomal arm SCNAs. The focal amplification on chromosome 12 is identified by conliga
 111 at 0.75% and 0.5% purity and not detected by QDNaseq below 1%. At 0.75%, 0.5% and 0% purity, it is hard
 112 to distinguish whole chromosome SCNAs from noise generated by segmentation in the QDNaseq profiles. This

113 highlights the advantage of conliga's ability to assign loci to discrete states, meaning we can easily distinguish
114 when SCNAs are and are not different between loci. Despite using 4.5 fold fewer reads, conliga appears to be
115 more sensitive than QDNAseq.

116 In Figure 2b, we show that conliga is able to detect SCNAs at 3% purity in all samples (eight), five at 2%
117 and one at 0.5%. The limit of detection is dependent on the amplitude and lengths of SCNAs present in the
118 sample. Long chromosomal arm amplifications can be detected at 2-3% purity, while some focal amplifications
119 (particularly those occurring at loci with a bias towards obtaining a high number of counts) can be detected at
120 <1% purity (e.g. chr12 in OAC3, Figure 2a). The limit of detection also depends on the technical variability
121 of the protocol and the total number of reads per sample. Increasing the total number of reads beyond two
122 million and reducing technical variability would further improve the limit of detection.

123 These data demonstrate the potential clinical utility of FAST-SeqS coupled with conliga. Ciriello et al.
124 identified that either somatic single nucleotide variants (SNVs) or SCNAs [3] can drive oncogenesis. Currently,
125 there is a bias towards screening for SNVs using targeted gene panels [26] meaning SCNA-driven cancers may
126 not be detected. To this end, we analyzed samples with pre-malignant disease (Barrett's oesophagus) and were
127 able to detect clinically relevant copy number alterations, such as evidence for focal gains of PRKCI, ERBB2
128 and GATA6 and deletions of regions containing CDKN2A, PTPRD, SMAD4 and TP53 (Supplementary Fig.
129 2). This suggests that there is potential for FAST-SeqS to be used alongside existing low-cost gene panels to
130 detect SCNAs, in addition to SNVs, to screen and surveil patients for the development of cancer.

131 In addition to use as a detection tool, inexpensive production of FAST-SeqS data allows for large cohorts of
132 patients to be studied to find relationships between SCNA profiles and response to therapies, for example. With
133 this in mind, we looked at the average SCNA profiles across small cohorts of patients with OAC, Gastric cancer
134 and Barrett's oesophagus (Supplementary Fig. 2, Methods) which highlighted amplifications of known oncogenes
135 such as EGFR, MYC, GATA4, and MDM2, some with known drug targets, and deletions of tumor suppressor
136 genes, e.g. FHIT, TP53, SMAD4 and RUNX1. Other potential uses include low-cost screening of samples
137 in large-scale cancer genomes studies, such as ICGC or TCGA projects, prior to further genomic analyses.
138 Furthermore, due to the low-cost and low-input DNA required, several spatially or temporally related samples
139 can be analyzed for the purposes of determining how SCNAs accumulate in normal tissues and contribute to
140 tumor evolution, in a similar fashion to previous studies on somatic mutations in the eyelid epidermis [27].

141 Areas for future study could include determining an acceptable number of reads which balances the cost and
142 limit of detection, finding ways to minimise the technical variability, and altering the number of reads obtained
143 at specific loci to increase statistical power in regions of interest.

144 We have shown that FAST-SeqS data can be used as a viable, inexpensive, and simple alternative to LC WGS
145 for the purpose of SCNA detection and quantification. conliga provides accurate and high-resolution SCNA
146 profiles across the genome and at regions of interest such as oncogenes and tumor suppressors. conliga (applied
147 to FAST-SeqS data with two million reads per sample) is particularly useful in detecting and discriminating
148 SCNAs in low purity samples and our results suggest it to be more sensitive than QDNAseq (using LC WGS,

149 nine million reads) for this purpose. We believe that conliga makes FAST-SeqS data a clinically valuable
150 diagnostic assay to detect and monitor patients for the development of cancer, as well as a useful research tool,
151 enabling inexpensive and fast SCNA profiling of cancer samples.

152 Methods

153 conliga: statistical model

154 Statistical model for sample counts

155 We model the sample counts, in L selected loci, by assuming that the count at locus l in chromosome arm r in
156 sample j is distributed:

$$y_{r,l,j} \sim \text{Binomial}(n_j, \theta_{r,l,j}) \quad (1)$$

157 Here, n_j is the total number of sequencing reads aligned to the L loci in sample j , $\theta_{r,l,j}$ represents the probability
158 of observing an aligned read at locus l in chromosome arm r in sample j . We model $\theta_{r,l,j}$ as follows:

$$\theta_{r,l,j} \sim \text{Beta}(s_j \hat{c}_{r,l,j} m_{r,l}, s_j (1 - \hat{c}_{r,l,j} m_{r,l})) \quad (2)$$

159 Here, s_j is the inverse dispersion variable for sample j where $s_j > 0$, $m_{r,l}$ represents the probability of an aligned
160 sequencing read originating from locus l in chromosome arm r in a control sample, where $\sum_r \sum_{l=1}^{L_r} m_{r,l} = 1$
161 and $\hat{c}_{r,l,j}$ is the relative copy number at locus l in chromosome arm r in sample j . The number of loci in each
162 chromosome arm is denoted as L_r and so the total number of loci, $L = \sum_r L_r$.

163 We can interpret m as defining the bias in observing aligned read counts from the FAST-SeqS protocol. This
164 bias can be explained by unequal PCR efficiencies between loci in addition to biases in aligning reads uniquely
165 to FAST-SeqS loci, among other factors. Note that:

$$\mathbb{E}[\theta_{r,l,j}] = \hat{c}_{r,l,j} m_{r,l} \quad (3)$$

166 We can interpret this equation intuitively; the relative copy number scales the probability of reads to
167 align to a locus. For example, if the relative copy number of a locus is 2 we expect the proportion of reads at
168 the locus to double. This fits with our observations shown in Supplementary Fig. 1.

169 The inverse dispersion variable, s_j , is sample specific and reflects our observations that the level of dispersion
170 varies between samples. This variation in dispersion between samples might be due to varying levels of DNA
171 degradation and/or varying quantities of starting material between samples, among other factors. s_j relates to
172 the variance and the mean of $\theta_{r,l,j}$ in the following way:

$$\text{Var}(\theta_{r,l,j}) = \frac{1}{s_j + 1} \left(\mathbb{E}[\theta_{r,l,j}] - \mathbb{E}[\theta_{r,l,j}]^2 \right) \quad (4)$$

173 The expected count, $y_{r,l,j}$, in chromosome arm r at locus l in sample j is:

$$\mathbb{E}[y_{r,l,j} | \theta_{r,l,j}] = \mu = n_j \hat{c}_{r,l,j} m_{r,l} \quad (5)$$

174 The variance of $y_{r,l,j}$ can be written as a quadratic function of μ with the coefficients being a function of n_j
175 and s_j :

$$\text{Var}(y_{r,l,j} | \theta_{r,l,j}) = \left(1 + \frac{n_j - 1}{s_j + 1}\right) \mu - \left(\frac{1}{n_j} + \frac{n_j - 1}{s_j + 1}\right) \mu^2 \quad (6)$$

176 Note that in the limit $s_j \rightarrow \infty$, a Binomial noise model is recovered.

177 Probabilistic generative model of loci counts for control samples

178 We assume that the loci within a control sample, k , have equal copy numbers (diploid). This means that the
179 RCN for each locus is 1. By setting $\hat{c}_{r,l,k} = 1$, we model the generative process of counts from a control sample
180 as follows:

$$\begin{aligned} s_k | \psi &\sim \text{Gamma}(\psi_{\text{shape}}, \psi_{\text{scale}}) \\ m_{r,l} | \phi &\sim \text{Beta}(\phi_{c,r,l}, \phi_{d,r,l}) \\ \theta_{r,l,k} | s_k, m_{r,l} &\sim \text{Beta}(s_k m_{r,l}, s_k(1 - m_{r,l})) \\ x_{r,l,k} | \theta_{r,l,k}, n_k &\sim \text{Binomial}(n_k, \theta_{r,l,k}) \end{aligned} \quad (7)$$

181 Here, $\text{Gamma}(\psi_{\text{shape}}, \psi_{\text{scale}})$ represents the prior distribution over the sample specific inverse dispersion pa-
182 rameter, s_k , and $\text{Beta}(\phi_{c,r,l}, \phi_{d,r,l})$ defines the prior distribution over $m_{r,l}$.

183 Linking FAST-SeqS loci using a hidden Markov model

184 We assume that chromosome arms are independent. By that we mean, the RCN of the first locus in arm q
185 is independent of the RCN of the last locus in arm p from the same chromosome (and all other chromosome
186 arms). As such, we model each chromosome arm as an independent Markov chain for each sample j . We denote
187 (note that for simplicity we have dropped the sample index j):

- 188 • $z_{r,l}$ as the *hidden state* (or *copy number state*) of the Markov chain at locus l in chromosome arm r
- 189 • π^0 as the *initial distribution* of the first locus ($l = 1$), in chromosome r
- 190 • π_u as the *transition distribution* for hidden state, u
- 191 • \hat{c}_u as the *relative copy number* associated with hidden state, u .

192 The first locus of a chromosome arm ($l = 1$) is distributed:

$$z_{r,1} \sim \pi^0 \quad (8)$$

193 For all other loci ($l > 1$):

$$z_{r,l} \mid z_{r,l-1} \sim \pi_{(z_{r,l-1})} \quad (9)$$

194 The count, $y_{r,l}$, at locus l in chromosome arm r is conditionally independent of the hidden states and
 195 observations of other loci:

$$\begin{aligned} \theta_{r,l} \mid \hat{c}, z_{r,l}, m_{r,l}, s &\sim \text{Beta}(s\hat{c}_{z_{r,l}}m_{r,l}, s(1 - \hat{c}_{z_{r,l}}m_{r,l})) \\ y_{r,l} \mid \theta_{r,l}, n &\sim \text{Binomial}(n, \theta_{r,l}) \end{aligned} \quad (10)$$

196 The joint density for L_r loci in chromosome arm r is:

$$\begin{aligned} p(z_{r,1:L_r}, y_{r,1:L_r}, \theta_{r,1:L_r}) &= p(y_{r,1} \mid z_{r,1}, \theta_{r,1})p(\theta_{r,1} \mid z_{r,1})p(z_{r,1}) \\ &\quad \prod_{l=2}^{L_r} p(y_{r,l} \mid z_{r,l}, \theta_{r,l})p(\theta_{r,l} \mid z_{r,l})p(z_{r,l} \mid z_{r,l-1}) \\ &= \pi_{z_{r,1}}^0 p(y_{r,1} \mid z_{r,1}, \theta_{r,1})p(\theta_{r,1} \mid z_{r,1}) \\ &\quad \prod_{l=2}^{L_r} \pi_{z_{r,l-1}, z_{r,l}} p(y_{r,l} \mid z_{r,l}, \theta_{r,l})p(\theta_{r,l} \mid z_{r,l}) \end{aligned} \quad (11)$$

197 where, $z_{r,1:L_r}$ denotes the sequence $\{z_{r,1}, \dots, z_{r,L_r}\}$, $y_{r,1:L_r}$ denotes $\{y_{r,1}, \dots, y_{r,L_r}\}$, and $\theta_{r,1:L_r}$ denotes $\{\theta_{r,1}, \dots, \theta_{r,L_r}\}$.

198 The joint density for all L loci in the genome is given by:

$$p(\mathbf{z}, \mathbf{y}, \boldsymbol{\theta}) = \prod_r p(z_{r,1:L_r}, y_{r,1:L_r}, \theta_{r,1:L_r}) \quad (12)$$

199 Probabilistic generative model of a sample's relative copy number profile

200 The number of copy number states present in a sample is unknown a priori. In samples that have equal copies
 201 of each locus, only one copy number state is present. Conversely, it is possible (although unlikely) that each
 202 locus has its own unique copy number, meaning that there could be up to L copy number states in a sample.
 203 Additionally, we expect neighboring loci to share the same copy number given their genomic distance from
 204 each other (Supplementary Fig. 1). To address these two features of the data, we used the sticky hierarchical
 205 Dirichlet process hidden Markov model (sticky HDP-HMM) [21] as a framework to model the generative process
 206 of a sample's relative copy number profile. By doing so, we adequately model the spatial persistence of copy
 207 number states and allow for countably infinite numbers of states within a sample. The generative model is as
 208 follows:

$$\begin{aligned}
 \beta &| \gamma \sim \text{GEM}(\gamma) \\
 \pi^0 &| \alpha, \beta \sim \text{DP}(\alpha, \beta) \\
 \pi_u &| \alpha, \kappa, \beta \sim \text{DP}\left(\alpha + \kappa, \frac{\alpha\beta + \kappa\delta_u}{\alpha + \kappa}\right) \\
 \hat{c}_u &| H, \lambda \sim H(\lambda) \\
 z_{r,1} &| \pi^0 \sim \pi^0
 \end{aligned} \tag{13}$$

$$z_{r,l} | \{\pi_u\}_{u=1}^\infty, z_{r,l-1} \sim \pi_{z_{r,l-1}}, \text{ for } l > 1$$

$$\tilde{s} | \omega \sim \text{Gamma}(\omega_{shape}, \omega_{scale})$$

$$\tilde{\theta}_{r,l} | \{\hat{c}_u\}_{u=1}^\infty, z_{r,l}, \hat{m}_{r,l}, \tilde{s} \sim \text{Beta}(\tilde{s}\hat{c}_{z_{r,l}}\hat{m}_{r,l}, \tilde{s}(1 - \hat{c}_{z_{r,l}}\hat{m}_{r,l}))$$

$$y_{r,l} | \tilde{\theta}_{r,l}, \tilde{n}, \sim \text{Binomial}(\tilde{n}, \tilde{\theta}_{r,l})$$

209 Note that we use \tilde{n} , \tilde{s} , $\tilde{\theta}_{r,l}$ to distinguish these variables from those in the probabilistic model of control counts
 210 (equation 7) and denote them as specific to the sample with copy number profile. Here, GEM denotes the
 211 stick-breaking construction of the Dirichlet Process as described in Fox *et al.* [21]. γ is a hyperparameter of
 212 the sticky HDP-HMM and represents our prior on the number of copy number states in the sample; the greater
 213 the value of γ , the greater number of copy number states we expect in the sample. Each row of the transition
 214 matrix, π_u , is drawn from a Dirichlet Process and depends on β , α and κ . It can be shown that:

$$\mathbb{E}[\pi_{u,v} | \alpha, \beta, \kappa] = \frac{\alpha\beta_v + \kappa\delta_{u,v}}{\alpha + \kappa} \tag{14}$$

215 where $\delta_{u,v}$ represents the discrete Kronecker delta function. If we define $\rho = \frac{\kappa}{\alpha + \kappa}$ (as in Fox *et al.* [21]) and by
 216 noting that $\alpha = (1 - \rho)(\alpha + \kappa)$, we obtain:

$$\mathbb{E}[\pi_{u,v} | \beta, \rho] = (1 - \rho)\beta_v + \rho\delta_{u,v} \tag{15}$$

217 As such, we see that ρ defines how much weight is placed on self-transition within a copy number state.
 218 The vector, β , itself drawn from a Dirichlet Process, represents the global transition distribution and holds
 219 information about the proportion of loci expected in each copy number state.

220 The variance of the transition probability from copy number state u to v is given by:

$$\text{Var}(\pi_{u,v} | \alpha, \beta, \kappa) = \frac{\mathbb{E}[\pi_{u,v} | \alpha, \beta, \kappa] (1 - \mathbb{E}[\pi_{u,v} | \alpha, \beta, \kappa])}{\alpha + \kappa + 1} \tag{16}$$

221 We see that $\alpha + \kappa$ is inversely proportional to the variance of the state transition probabilities.

222 H is the prior base distribution of the Dirichlet Process and represents a parametric distribution, which in
 223 this case is a Gamma distribution, with parameters λ . It can be viewed as our prior probability distribution on
 224 the relative copy number values of the hidden states.

225 Note that $\hat{m}_{r,l}$ refers to the maximum a posteriori (MAP) value of $m_{r,l}$ and is such assumed to be a known
226 quantity in equation 13. For simplicity, the hyperparameters (α , κ , γ , λ , ω and n) are shown as fixed quantities
227 in the model. In practice, γ , λ , ω and n are treated as fixed, while the model is parameterized in terms of ρ and
228 $(\alpha + \kappa)$, with a Beta prior placed on ρ and a Gamma prior placed on $(\alpha + \kappa)$ as in Fox *et al.* [21]. See the section
229 on inference for further details of prior distributions used and Supplementary Note 3 for further discussion on
230 the model.

231 Inference

232 Inference of loci count proportion bias (m)

233 Given a set of K control samples, and their loci counts, \mathbf{x}_k , we used our model defined in equation 7 and
234 Markov Chain Monte Carlo (MCMC) methods to infer the latent variables \mathbf{m} and \mathbf{s} (the vector of sample
235 specific inverse dispersion parameters). A Metropolis-Hastings MCMC algorithm was used to obtain a sample
236 of the posterior probability of $m_{r,l}$ for all r and l , and s_k , for each sample k . Full details of the algorithms are
237 provided in Supplementary Notes 4 and 5. Count data for samples analyzed in this study, processed by the
238 pipeline described, are provided in Supplementary Table 10.

239 For each sequencing experiment, a suitable set of controls samples were used (see Supplementary Table 11
240 for the list of samples used in each experiment). As described in equation 7, control samples were assumed
241 to have a relative copy number of one at each locus. In all experiments described in this paper, we used the
242 following values for the hyperparameters:

- 243 • $\psi_{shape} = 1.5$, $\psi_{scale} = 10^6$; where ψ_{shape} and ψ_{scale} define the shape and scale of the Gamma prior
244 distribution on s_k , respectively.
- 245 • $\phi_{c,r,l} = 1$ and $\phi_{d,r,l} = 1$ for all r and l ; i.e. we used a flat Beta(1, 1) prior for all $m_{r,l}$

246 In each sequencing experiment, 20,000 iterations of the MCMC were run and the first 5,000 iterations were
247 discarded (burn-in). Maximum a posteriori (MAP) estimates of \mathbf{m} (denoted as $\hat{\mathbf{m}}$) were obtained by determining
248 the mode of the sampled posterior densities for each locus using the KernSmooth R package [28]. Note that the
249 MAP estimates are unlikely to sum to 1 exactly, and as such we enforced this by setting $\hat{m}_{r,l} = \frac{\hat{m}_{r,l}}{\sum_r \sum_l \hat{m}_{r,l}}$.

250 Inference of relative copy number profile

251 Given $\hat{\mathbf{m}}$ and the loci counts (\mathbf{y}) for a sample with unknown copy number profile, we used the generative model
252 defined in equation 13 and MCMC methods (based on algorithm 3 in Fox *et al.* [21]) to infer the latent variables
253 in our model. MCMC methods were used to obtain a sample of the posterior probability of the hidden state of
254 each locus ($z_{r,l}$ for all r and l), the relative copy number of each hidden state (\hat{c}_u), the sample specific inverse
255 dispersion (\tilde{s}), along with other latent variables in our generative model. Full details of the MCMC algorithms
256 can be found in Supplementary Notes 4 and 6. In all experiments described in this paper, we used the following
257 values for the hyperparameters:

- 258 • $\gamma = 1$
- 259 • Gamma(2000, 10) prior distribution (defined by shape and scale) was placed on $(\alpha + \kappa)$
- 260 • Beta(100000, 100) prior was placed on ρ
- 261 • Gamma(3, 1) prior distribution (defined by shape and scale) was placed on the relative copy number value
- 262 of the hidden states; the shape and scale parameters are defined by λ in equation 13
- 263 • $\omega_{shape} = 1.5$, $\omega_{scale} = 10^6$; where ω_{shape} and ω_{scale} define the shape and scale of the Gamma prior
- 264 distribution on \tilde{s} , respectively

265 The output of the MCMC was summarized in two main ways, 1) by marginalizing out the copy number
266 state information and computing the MAP estimate (using KernSmooth R package [28]) and credible interval of
267 the relative copy number of each locus, 2) by making use of the copy number state assignments in the following
268 way:

- 269 1. we determined the MAP number of states observed in the MCMC chain (after burn-in). This was achieved
270 by calculating the number of populated states in each iteration of the MCMC, and then choosing the most
271 frequently observed number of populated states. Note that a state was considered populated in an iteration
272 of the MCMC if at least one locus was assigned to it.
- 273 2. we filtered the iterations of the MCMC (after burn-in), choosing only those iterations that had the number
274 of populated states equal to the MAP number of states.
- 275 3. we used the Stephens algorithm (algorithm 2 in the paper) [29] along with the Hungarian (Munkres)
276 algorithm [30] to relabel the states, to resolve the label switching problem inherent in MCMC methods.
- 277 4. we calculated the MAP estimate and credible intervals for the relative copy number values of each relabeled
278 state.
- 279 5. we assigned each locus to a relabeled state, choosing the relabeled state it was most frequently assigned
280 to in the filtered iterations of the MCMC chain.

281 For the results presented in Figure 2, summarization method 2 was used. For all other results presented
282 in the paper, summarization method 1 was used. For the oesophageal cancer, gastric cancer and Barrett's
283 oesophagus samples, 50,000 iterations of the MCMC were run and the chain was thinned such that every 5th
284 iteration of the MCMC was output to file. Additionally, the first 20,000 iterations of the MCMC were discarded
285 (burn-in), to ensure the Markov chain had reached its equilibrium distribution. For the in silico diluted samples,
286 presented in Figure 2, 30,000 iterations were run, with the chain thinned so that every 5th sample was output
287 to file and the first 5,000 iterations of the MCMC were discarded.

288 **Sample preparation and sequencing of samples**

289 **Sample preparation and generation of FAST-SeqS data**

290 Sequencing libraries were prepared using two rounds of PCR, using a similar protocol to previously published
291 methods [18, 19]. Each extracted DNA sample underwent a 50 μ l first round PCR reaction with 10 μ l 5x Phusion
292 HF Buffer (ThermoFisher Scientific), 1 μ l 10 mM dNTP (ThermoFisher Scientific), 5 μ l of both the forward and
293 reverse primers (0.5 μ M) each (Sigma-Aldrich), 0.5 μ l Phusion Hot Start II DNA Polymerase 2U/ μ l, 5-10 μ l DNA
294 template depending on the extracted concentration, and RNase free water to make the total reaction volume.
295 The cycling conditions for the L1PA7 primers were 98 °C for 120 s followed 2 cycles of 98 °C for 10 s, 57 °C for
296 120 s, and 72 °C for 120 s. The second round was also carried out as a 50 μ l sample reaction using 20 μ l taken from
297 the first round. The rest of the reaction constituents were the same as the first round reaction with the exception
298 of primers (Supplementary Table 12), which contained a unique index for each sample. The cycling conditions
299 for the second round reaction were 98 °C for 120 s followed by 13 cycles of 98 °C for 10 s, 65 °C for 15 s, and 72 °C
300 for 15 s for all the primers. After the second round, samples underwent quantification using the 2200 TapeStation
301 (Agilent), Agilent 2100 Bioanalyser (Agilent) and Kapa quantification (KapaBiosystems) prior to submission
302 for sequencing. The samples were then pooled in equimolar concentrations and gel extracted according to
303 manufacturer's instructions (Qiaquick gel extraction kit, Qiagen). Finally the samples were submitted for
304 sequencing on a MiSeq (Illumina) platform. All samples were run with 20% PhiX to increase complexity for
305 sequencing. Sequencing was performed as 150bp single end. Samples were run with at least three normal
306 controls prepared at the same time and sequenced on the same platform.

307 **Sample preparation and generation of high-coverage WGS data**

308 WGS library preparation and sequencing was performed as previously described by Secrier *et al.* [6].

309 **In silico generation of low-coverage WGS data**

310 For our purposes, LC WGS data was defined as nine million single-end 50 base pair reads per sample because this
311 was the type of data analyzed in Scheinin *et al.* [17]. Samples are typically multiplexed together and sequenced
312 on a single Illumina sequencing lane. After processing and alignment of the reads, we expect approximately
313 0.1X coverage of the genome (as per analysis described in Scheinin *et al.*). We obtained LC WGS data by
314 down-sampling reads from HC WGS BAM files in the following way:

- 315 1. we selected a subset of the alignments, containing only reads sequenced on a single lane (chosen to be the
316 lane from the first read in the BAM file), and trimmed the reads and Phred scores to the first 50 base
317 pairs using a custom Bash script.
- 318 2. The resulting alignments were filtered (using samtools [31] version 0.1.18), excluding those that were
319 secondary alignments (-F 256) and including only those that were first in a pair (-f 64) and output to a
320 new BAM file.

- 321 3. This BAM file was down-sampled to 9 million reads/alignments using the DownsampleSam command from
322 Picard tools (<http://broadinstitute.github.io/picard>, version 2.9.1) using the "Chained" strategy.
- 323 4. The resulting BAM file was converted to FASTQ by SamToFastq (Picard tools).
- 324 5. The FASTQ file was aligned to GRCh38 (GenBank accession: GCA_000001405.15, no alt analysis set)
325 using BWA-backtrack (bwa samse and bwa aln, version 0.7.15-r1140) [32], which is more suitable for reads
326 below 70 base pairs in length.
- 327 6. In the resulting BAM file, we removed PCR duplicates and removed alignments with mapping quality
328 below 37 as per the analysis undertaken by Scheinin *et al.* [17] using samtools (version 0.1.18).

329 We performed these steps for 11 oesophageal samples and their matched normal samples along with an
330 additional four normal samples obtained from other patients (Supplementary Table 1). This resulted in greater
331 than seven million primary alignments per sample.

332 **In silico generation of FAST-SeqS dilution data**

333 We performed an in silico dilution of FAST-SeqS data by mixing sequencing reads from control samples with
334 reads from OAC samples. Since the number of reads in the matched controls were insufficient to create samples
335 with two million reads, we created a pool of control reads (in silico) which were used to dilute the OAC samples.
336 This was done by sub-sampling two million reads from 12 control samples (which were prepared and sequenced in
337 the same batch as the OAC samples). The total number of reads from these 12 control samples was 14,405,596.
338 To obtain a pool of 2 million reads, we used the 'sample' command from seqtk ([urlhttps://github.com/lh3/seqtk](https://github.com/lh3/seqtk),
339 version: 1.2-r101) to sample a proportion (2/14.405596) of each control sample's reads and merged these together
340 into a single FASTQ file. The reads that were sub-sampled were removed from the control samples (using a
341 custom python script) to avoid using the same reads to fit *m*.

342 We mixed the pool of control reads with the OAC samples in varying proportions to achieve a desired diluted
343 tumor purity. The OAC samples did not have a tumor purity of 100%, instead they were themselves a mixture
344 of tumor and normal DNA. The purity of these samples were determined by ASCAT-NGS (version 2.1) [22].
345 Based on ASCAT's purity value, we calculated the number of reads required from the OAC sample to achieve
346 a desired dilution and total number of reads. This was calculated as follows:

$$\text{required tumor reads} = \text{round} \left(\frac{\text{desired purity proportion} \cdot \text{required total reads}}{\text{ASCAT inferred purity proportion}} \right) \quad (17)$$

347 Hence, the number of control reads required were:

$$\text{required control reads} = \text{required total reads} - \text{required tumor reads} \quad (18)$$

348 We produced in silico dilution FASTQ files in the following way:

- 349 1. we used the ‘sample’ command from seqtk to sample the required number of tumor reads from the OAC
350 FAST-SeqS FASTQ file
- 351 2. we used the ‘sample’ command from seqtk to sample the required number of control reads from the pooled
352 control reads FASTQ file
- 353 3. We merged the sampled tumor and control reads into a single FASTQ file

354 We performed these steps for each OAC sample to create diluted samples with two million total reads and
355 the following purity values: 0.3, 0.25, 0.2, 0.15, 0.1, 0.08, 0.06, 0.05, 0.04, 0.03, 0.02, 0.01, 0.0075, 0.005, 0.0025
356 and 0. Here purity is defined as the proportion of tumor reads in the sample. Of the 11 OAC samples, 8
357 (OAC1-7 and 9, Supplementary Table 1) were of sufficient initial tumor purity to feasibly create all the desired
358 dilution levels.

359 **In silico generation of LC WGS dilution data**

360 We produced in silico diluted LC WGS tumor samples by mixing reads from tumor and matched normal LC
361 WGS BAM files (previously downsampled and filtered as described above). We first calculated the number
362 of reads in the tumor BAM and normal BAM files using samtools (`samtools view -F 256 -c [BAM file]`).
363 Next, we calculated the number of reads required using equations 17 and 18. Using the DownsampleSAM
364 command (Picard tools) and the ‘HighAccuracy’ strategy, we sampled the corresponding desired proportion of
365 reads from the tumor BAM file and normal BAM file. We used samtools to merge the resulting sampled tumor
366 BAM file with the normal BAM file into a single file representing the diluted sample. We aimed to obtain seven
367 million filtered primary alignments per diluted sample (as this is what we expect from nine million reads after
368 alignment and filtering) and dilution levels which matched the diluted FAST-SeqS samples. This was performed
369 for 8 OAC samples and their matched normals (OAC1-7 and 9).

370 **Processing of FAST-SeqS sequencing data to counts**

371 Each sequencing run of the Illumina MiSeq platform produced a BCL file which was converted to FASTQ format
372 (using Illumina’s `bcl2fastq` tool). Sequencing reads that failed the Illumina chastity filter were removed. The
373 FASTQ file was demultiplexed into separate FASTQ files corresponding to each sample using the `demuxFQ` tool
374 (<https://genomicsequencing.cruk.cam.ac.uk/glsstatic/lablink/downloads/DemultiplexingGuide.html>)
375 with the default settings. The sample barcodes are provided in Supplementary Table 12. Each sample’s FASTQ
376 file was then processed through a custom pipeline which we describe below.

377 **Identifying forward primer position**

378 For each read in the FASTQ file, the position of the forward primer sequence was detected by searching for the
379 sequence with the minimum hamming distance to the forward primer sequence using a sliding window. Reads
380 with a minimum hamming distance greater than 5 were discarded.

381 **Read trimming**

382 The portion of the reads before and including the forward primer sequence were trimmed. The ends of the
383 reads were also trimmed such that the length of the reads used for downstream analyses were 100 base pairs
384 minus the forward primer length. Any reads shorter than 100 base pairs minus the forward primer length after
385 trimming were discarded.

386 **Quality control**

387 After trimming, reads were discarded if they contained at least one base with a Phred quality score less than
388 20 and/or contained one or more ambiguous base calls (N).

389 **Obtaining unique sequences and counts per unique sequence**

390 To avoid aligning the same sequence multiple times, only unique read sequences were kept. For each unique
391 read, the number of identical fragments were recorded.

392 **Alignment of unique sequences**

393 Unique raw read sequences were aligned with Bowtie 1.0.0 [33] (using the option: -r). Three mismatches were
394 permitted (option: -v3) and reads aligning to multiple locations were discarded (option: -m1). The reads were
395 aligned to GRCh38 (GenBank accession: GCA_000001405.15, no alt analysis set).

396 **Counts and alignments combined**

397 Each sample's unique read alignments and their corresponding unique read counts were combined into a single
398 file consisting of a matrix of counts. The rows corresponded to genomic positions (the union of genomic positions
399 from the alignments in all samples) and columns corresponded to samples. The first three columns of the matrix
400 corresponded to the chromosome, position and strand for the locus, respectively. The matrix of counts used in
401 this analysis can be found in the conliga R package and in Supplementary Table 10.

402 **Selecting loci**

403 Rows of the count matrix corresponding to genomic loci within chromosomes X, Y and within unplaced or
404 unresolved contigs were discarded. For each batch of samples, genomic loci obtaining a zero count in any one of
405 a set of control samples were also discarded. Depending on the sequencing batch we analyzed and the controls
406 chosen to filter loci (Supplementary Table 11), this resulted in approximately 10,000 - 12,000 genomic loci across
407 chromosomes 1 to 22.

408 **Analysis of copy number from FAST-SeqS data**

409 conliga (version 0.1.0) was used to analyze all FAST-SeqS samples in this study (Supplementary Table 1) using
410 R (version 3.2.3) [34] and RcppAramdillo (version 0.6.500.4.0) [35]. Of the 15 OAC samples sequenced, four
411 were excluded due to their obtaining fewer than 350,000 reads. Two control samples were excluded due to their
412 inferred RCN profiles having two main hidden states incompatible with their supposed ‘normal’ status. The
413 values for the priors used and MCMC settings are stated in the inference sections above. The samples used as
414 a basis to filter loci and fit \hat{m} for each experiment are listed in Supplementary Table 9.

415 **Analysis of copy number from high coverage WGS data**

416 High coverage WGS samples were processed and aligned using BWA-MEM [36] (version 0.5.9) and total copy
417 number (TCN) profiles and normal contamination estimates were provided by ASCAT-NGS (version 2.1) using
418 a pipeline previously described by Secrier *et al.* [6]. The only exception to this was that the reads were aligned
419 to GRCh38 (GenBank accession: GCA_000001405.15, no alt analysis set) rather than GRCh37.

420 **Analysis of copy number from low-coverage WGS data**

421 QDNAseq (version 1.6.1) was used to obtain relative copy number calls for all LC WGS data. The bin size
422 used was 15Kb as per the analysis performed in Scheinin *et al.* [17] for 0.1X LC WGS. The bins were created
423 using GRCh38 (BSgenome.Hsapiens.NCBI.GRCh38) and a mappability file (bigWig format) for 50-mers was
424 created for GRCh38 using the GEM library (GEM-binaries-Linux-x86_64-core_i3-20130406-045632) <https://sourceforge.net/projects/gemlibrary/>. 15 normal LC WGS samples (Supplementary Table 1), were
425 used to run the applyFilters and iterateResiduals functions. 11 of these 15 samples correspond to the matched
426 normals of the oesophageal samples (Supplementary Table 1). We did not run the functions normalizeBins and
427 normalizeSegmentedBins which scale the read counts by the median value. This was not necessary and would
428 make the comparison between ASCAT, QDNAseq and conliga results more difficult to interpret.
429

430 **Comparison of copy number between methods**

ASCAT outputs total copy number (TCN) in contiguous genomic regions, QDNAseq outputs relative copy
number (RCN) in 15 Kb bins across the genome and conliga outputs RCN values at specific FAST-SeqS loci.
To make a fair comparison between the tools, it was necessary to convert ASCAT’s TCN calls to RCN as follows:

$$\text{RCN}_i = \frac{(1 - \text{normal}) \cdot \text{TCN}_i + \text{normal} \cdot 2}{\text{mean TCN}} \quad (19)$$

Here, normal represents the estimated normal contamination value provided by ASCAT and i represents a
contiguous genomic region or a discrete locus or fragment. In the case of a contiguous region, the mean TCN

(or ploidy) was calculated as follows:

$$\text{mean TCN} = \frac{\sum_i (\text{TCN}_i \cdot \text{length}_i)}{\sum_i \text{length}_i} \quad (20)$$

and in the case of discrete loci or fragments:

$$\text{mean TCN} = \frac{\sum_i \text{TCN}_i}{L} \quad (21)$$

431 where L represents the total number of loci or fragments considered.

432 In Figure 1e and f, we compared the RCN values at the intersection of genomic loci across ASCAT, QDNAseq
433 and conliga. Since this intersection represented a subset of each method's genomic loci, the RCN values were
434 rescaled considering only this subset. QDNAseq and conliga RCN values were rescaled by the sample's mean
435 RCN of the considered loci. ASCAT's RCN was calculated using equations 19 and 21.

In figure 1g, we compared RCN values in genes of interest. Recurrently amplified and deleted genes were
obtained from Dulak *et al.* [23] and Ross-innes *et al.* [12]. Here, ASCAT's RCN values were calculated using
equations 19 and 20 using all called regions for each sample. For each gene in each sample, the weighted mean
of the relative copy number (weighted by the length of the overlapping called region) was computed for ASCAT
and QDNAseq. This was calculated as follows:

$$\text{RCN}_{\text{gene}} = \frac{\sum_i \text{RCN}_i \cdot l_i}{\sum_i l_i} \quad (22)$$

436 where l_i represents the length of the overlapping portion of the called region with the gene.

437 For conliga, if loci occurred within the gene, the mean of the RCN values within the gene was used, otherwise
438 the loci directly upstream and downstream, i.e. either side, of the gene were used and a mean value was taken.
439 See Supplementary Table 4 for the full list of genes used in the analysis.

440 **Computing Pearson correlation**

441 For each sample, the Pearson correlation coefficient between ASCAT and conliga was calculated. We used
442 ASCAT's TCN and conliga RCN values at the intersection of genomic loci between ASCAT and conliga. The
443 median value of the sample's correlation coefficients was reported (all sample correlation coefficients can be
444 found in Supplementary Table 3).

445 For each sample, the Pearson correlation coefficient between ASCAT and QDNAseq was calculated. We
446 used the intersection of QDNAseq bins with ASCAT copy number regions, using the length-weighted mean of
447 ASCAT's overlapping TCN values.

448 When calculating the Pearson correlation for all calls across all samples, we used the re-scaled RCN value at
449 the intersecting genomic loci between ASCAT, QDNAseq and conliga, using the rescaled RCN values described
450 above for Figures 1e and f.

451 Code availability

452 conliga source code is freely available under an open-source GPLv2 license at [https://github.com/samabs/](https://github.com/samabs/conliga)
453 [conliga](https://github.com/samabs/conliga) and as Supplementary Software.

454 Data availability

455 The WGS and FAST-SeqS data can be found at the European Genome-phenome Archive (EGA) under accession
456 EGAD00001004289. The copy number results obtained from ASCAT, QDNAseq and conliga can be found
457 https://osf.io/bhx6f/?view_only=ed25e2fb521d46239e5274c032350f0b

458 Oesophageal Cancer Clinical and Molecular Stratification (OCCAMS) 459 Consortium

460 Rebecca C. Fitzgerald¹, Ayesha Noorani¹, Paul A.W. Edwards^{1,2}, Nicola Grehan¹, Barbara Nutzinger¹, Caitri-
461 ona Hughes¹, Elwira Fidziukiewicz¹, Jan Bornschein¹, Shona MacRae¹, Jason Crawte¹, Alex Northrop¹, Gian-
462 marco Contino¹, Xiaodun Li¹, Rachel de la Rue¹, Maria O'Donovan^{1,3}, Ahmad Miremadi^{1,3}, Shalini Malhotra^{1,3},
463 Monika Tripathi^{1,3}, Simon Tavaré², Andy G. Lynch², Matthew Eldridge², Maria Secrier², Lawrence Bower²,
464 Ginny Devonshire², Juliane Perner², Sriganesh Jammula², Jim Davies⁵, Charles Crichton⁵, Nick Carroll⁶, Pe-
465 ter Safranek⁶, Andrew Hindmarsh⁶, Vijayendran Sujendran⁶, Stephen J. Hayes^{7,14}, Yeng Ang^{7,8,29}, Shaun R.
466 Preston⁹, Sarah Oakes⁹, Izhar Bagwan⁹, Vicki Save¹⁰, Richard J.E. Skipworth¹⁰, Ted R. Hupp¹⁰, J. Robert
467 O'Neill^{10,23}, Olga Tucker^{11,33}, Andrew Beggs^{11,28}, Philippe Taniere¹¹, Sonia Puig¹¹, Timothy J. Underwood^{12,13},
468 Fergus Noble¹², Jack Owsley¹², Hugh Barr¹⁵, Neil Shepherd¹⁵, Oliver Old¹⁵, Jesper Lagergren^{16,25}, James
469 Gossage^{16,24}, Andrew Davies^{16,24}, Fujun Chang^{16,24}, Janine Zylstra^{16,24}, Ula Mahadeva¹⁶, Vicky Goh²⁴, Francesca
470 D. Ciccarelli²⁴, Grant Sanders¹⁷, Richard Berrisford¹⁷, Catherine Harden¹⁷, Mike Lewis¹⁸, Ed Cheong¹⁸,
471 Bhaskar Kumar¹⁸, Simon L Parsons¹⁹, Irshad Soomro¹⁹, Philip Kaye¹⁹, John Saunders¹⁹, Laurence Lovat²⁰,
472 Rehan Haidry²⁰, Laszlo Igali²¹, Michael Scott²², Sharmila Sothi²⁶, Sari Suortamo²⁶, Suzy Lishman²⁷, George
473 B. Hanna³¹, Christopher J. Peters³¹, Anna Grabowska³², Richard Turkington³⁴.

474 ¹ Medical Research Council Cancer Unit, Hutchison/Medical Research Council Research Centre, University of
475 Cambridge, Cambridge, UK

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

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

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

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

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

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

- 482 ⁸ Wigan and Leigh NHS Foundation Trust, Wigan, Manchester, UK, WN1 2NN
- 483 ⁹ Royal Surrey County Hospital NHS Foundation Trust, Guildford, UK, GU2 7XX
- 484 ¹⁰ Edinburgh Royal Infirmary, Edinburgh, UK, EH16 4SA
- 485 ¹¹ University Hospitals Birmingham NHS Foundation Trust, Birmingham, UK, B15 2GW
- 486 ¹² University Hospital Southampton NHS Foundation Trust, Southampton, UK, SO16 6YD
- 487 ¹³ Cancer Sciences Division, University of Southampton, Southampton, UK, SO17 1BJ
- 488 ¹⁴ Faculty of Medical and Human Sciences, University of Manchester, UK, M13 9PL
- 489 ¹⁵ Gloucester Royal Hospital, Gloucester, UK, GL1 3NN
- 490 ¹⁶ Guy's and St Thomas's NHS Foundation Trust, London, UK, SE1 7EH
- 491 ¹⁷ Plymouth Hospitals NHS Trust, Plymouth, UK, PL6 8DH
- 492 ¹⁸ Norfolk and Norwich University Hospital NHS Foundation Trust, Norwich, UK, NR4 7UY
- 493 ¹⁹ Nottingham University Hospitals NHS Trust, Nottingham, UK, NG7 2UH
- 494 ²⁰ University College London, London, UK, WC1E 6BT
- 495 ²¹ Norfolk and Waveney Cellular Pathology Network, Norwich, UK, NR4 7UY
- 496 ²² Wythenshawe Hospital, Manchester, UK, M23 9LT
- 497 ²³ Edinburgh University, Edinburgh, UK, EH8 9YL
- 498 ²⁴ King's College London, London, UK, WC2R 2LS
- 499 ²⁵ Karolinska Institutet, Stockholm, Sweden, SE 77
- 500 ²⁶ University Hospitals Coventry and Warwickshire NHS, Trust, Coventry, UK, CV2 2DX
- 501 ²⁷ Peterborough Hospitals NHS Trust, Peterborough City Hospital, Peterborough, UK, PE3 9GZ
- 502 ²⁸ Institute of Cancer and Genomic sciences, University of Birmingham, B15 2TT
- 503 ²⁹ GI science centre, University of Manchester, UK, M13 9PL.
- 504 ³⁰ Queen's Medical Centre, University of Nottingham, Nottingham, UK, NG7 2UH
- 505 ³¹ Imperial College NHS Trust, Imperial College London, UK, W2 1NY
- 506 ³² Queen's Medical Centre, University of Nottingham, Nottingham, UK
- 507 ³³ Heart of England NHS Foundation Trust, Birmingham, UK, B9 5SS
- 508 ³⁴ Centre for Cancer Research and Cell Biology, Queen's University Belfast, Northern Ireland, UK, BT7 1NN.

509 **Acknowledgements**

510 Funding for sample sequencing was through the Oesophageal Cancer Clinical and Molecular Stratification (OC-
511 CAMS) Consortium as part of the International Cancer Genome Consortium and was funded by a programme
512 grant from Cancer Research UK (RG66287). SA was funded by Wellcome Trust award [102272]. ERM was
513 funded by MRC Computational Biology Fellowship (MC_UU_12025, MRC Strategic Alliance Funding: MRC
514 Weatherall Institute of Molecular Medicine). We acknowledge the support of The University of Cambridge,
515 Cancer Research UK and Hutchison Whampoa Limited. In particular we acknowledge the support of the Can-

516 cer Research UK Cambridge Institute’s Genomics Core Facility. We thank the Human Research Tissue Bank,
517 which is supported by the National Institute for Health Research (NIHR) Cambridge Biomedical Research
518 Centre, from Addenbrooke’s Hospital. Additional infrastructure support was provided from the CRUK funded
519 Experimental Cancer Medicine Centre. The authors also wish to thank Sarah Field, Wing-Kit Leung, Charlie
520 Massie, Paul Coupland and Astrid Wendler for helpful discussions and Ginny Devonshire for her help with
521 uploading the sequencing data to EGA.

522 Author contributions statement

523 RCF, JMJW, SSZ conceived the clinical utility of the experimental approach for upper GI cancer studies.
524 SA, ERM, AGL conceived the broad computational approach for analysis of these experiments. MCVvL, MP,
525 XL, AnN, AJM, JB, AhM, led by SSZ, took the biological samples from patients to data. SA developed and
526 implemented the conliga method, analyzed and preprocessed the data under supervision from ERM and AGL.
527 LB analyzed the high-coverage WGS data. SA wrote the manuscript in consultation with AGL, ERM, RCF
528 and SSZ. All authors read and approved the final version of the manuscript.

529 References

- 530 [1] Beroukhi, R. *et al.* The landscape of somatic copy-number alteration across human cancers. *Nat.* **463**,
531 899–905 (2010). DOI [10.1038/nature08822](https://doi.org/10.1038/nature08822).
- 532 [2] Zack, T. I. *et al.* Pan-cancer patterns of somatic copy number alteration. *Nat. Genet.* **45**, 1134–1140
533 (2013). URL <http://dx.doi.org/10.1038/ng.2760>. DOI [10.1038/ng.2760](https://doi.org/10.1038/ng.2760).
- 534 [3] Ciriello, G. *et al.* Emerging landscape of oncogenic signatures across human cancers. *Nat. Genet.* **45**,
535 1127–1133 (2013). URL <http://dx.doi.org/10.1038/ng.2762>. DOI [10.1038/ng.2762](https://doi.org/10.1038/ng.2762).
- 536 [4] Cancer Genome Atlas Research Network *et al.* The Cancer Genome Atlas Pan-Cancer analysis project.
537 *Nat. Genet.* **45**, 1113–1120 (2013). URL <http://dx.doi.org/10.1038/ng.2764>. DOI [10.1038/ng.2764](https://doi.org/10.1038/ng.2764).
- 538 [5] Nones, K. *et al.* Genomic catastrophes frequently arise in esophageal adenocarcinoma and drive tu-
539 morigenesis. *Nat. Commun.* **5**, 1–9 (2014). URL <http://dx.doi.org/10.1038/ncomms6224>. DOI
540 [10.1038/ncomms6224](https://doi.org/10.1038/ncomms6224).
- 541 [6] Secrier, M. *et al.* Mutational signatures in esophageal adenocarcinoma define etiologically distinct subgroups
542 with therapeutic relevance. *Nat. Genet.* **48**, 1131–1141 (2016). DOI [10.1038/ng.3659](https://doi.org/10.1038/ng.3659).
- 543 [7] Frankell, A. M. *et al.* The landscape of selection in 551 Esophageal Adenocarcinomas defines genomic
544 biomarkers for the clinic. *bioRxiv* (2018). URL [https://www.biorxiv.org/content/early/2018/04/28/](https://www.biorxiv.org/content/early/2018/04/28/310029)
545 [310029](https://doi.org/10.1101/310029). [/dx.doi.org/10.1101/310029](https://doi.org/10.1101/310029).

- 546 [8] Bird-Lieberman, E. L. & Fitzgerald, R. C. Early diagnosis of oesophageal cancer. *Br. J. Can-*
547 *cer* **101**, 1–6 (2009). URL <http://www.nature.com/doi/10.1038/sj.bjc.6605126>. DOI
548 10.1038/sj.bjc.6605126.
- 549 [9] Weaver, J. M. *et al.* Ordering of mutations in preinvasive disease stages of esophageal carcinogenesis. *Nat.*
550 *Genet.* **46**, 837–843 (2014). URL <http://dx.doi.org/10.1038/ng.3013>. DOI 10.1038/ng.3013.
- 551 [10] Paulson, T. G. *et al.* Chromosomal Instability and Copy Number Alterations in Barrett’s Esophagus and
552 Esophageal Adenocarcinoma. *Clin. Cancer Res.* **15**, 3305–3315 (2009). DOI 10.1158/1078-0432.CCR-08-
553 2494.
- 554 [11] Li, X. *et al.* Temporal and Spatial Evolution of Somatic Chromosomal Alterations : A Case-Cohort Study
555 of Barrett’s Esophagus. *Cancer Prev. Res.* **7**, 114–128 (2014). DOI 10.1158/1940-6207.CAPR-13-0289.
- 556 [12] Ross-innes, C. S. *et al.* Whole-genome sequencing provides new insights into the clonal architecture of Bar-
557 rett’s esophagus and esophageal adenocarcinoma. *Nat. Genet.* **47**, 1038–1046 (2015). DOI 10.1038/ng.3357.
- 558 [13] Kallioniemi, A. *et al.* Comparative genomic hybridization for molecular cytogenetic analysis of solid tumors.
559 *Sci.* **258**, 818–821 (1992).
- 560 [14] Solinas-toldo, S., Lampel, S., Stilgenbauer, S., Nickolenko, J. & Benner, A. Matrix-Based Comparative
561 Genomic Hybridization: Biochips to Screen for Genomic Imbalances. *Genes, Chromosom. Cancer* **20**,
562 399–407 (1997).
- 563 [15] Kennedy, G. C. *et al.* Large-scale genotyping of complex DNA. *Nat. Biotechnol.* **21**, 1233–1237 (2003).
564 DOI 10.1038/nbt869.
- 565 [16] Ley, T. J. *et al.* DNA sequencing of a cytogenetically normal acute myeloid leukaemia genome. *Nat.* **456**,
566 66–72 (2008). DOI 10.1038/nature07485.
- 567 [17] Scheinin, I. *et al.* DNA copy number analysis of fresh and formalin-fixed specimens by shallow whole-genome
568 sequencing with identification and exclusion of problematic regions in the genome assembly. *Genome Res.*
569 **24**, 2022–2032 (2014). DOI 10.1101/gr.175141.114.Freely.
- 570 [18] Kinde, I., Papadopoulos, N., Kinzler, K. W. & Vogelstein, B. FAST-SeqS: A simple and efficient method
571 for the detection of aneuploidy by massively parallel sequencing. *PLoS ONE* **7** (2012). DOI 10.1371/jour-
572 nal.pone.0041162.
- 573 [19] Belic, J. *et al.* Rapid identification of plasma DNA samples with increased ctDNA levels by a modified
574 FAST-SeqS approach. *Clin. Chem.* **61**, 838–849 (2015). DOI 10.1373/clinchem.2014.234286.
- 575 [20] Douville, C., Springer, S., Kinde, I., Cohen, J. D. & Hruban, R. H. Detection of aneuploidy in patients
576 with cancer through amplification of long interspersed nucleotide elements (LINEs). *PNAS* (2018). DOI
577 10.1073/pnas.1717846115.

- 578 [21] Fox, E. B., Sudderth, E. B., Jordan, M. I. & Willsky, A. S. A Sticky HDP-HMM with application to
579 speaker diarization. *The Annals Appl. Stat.* **5**, 1020–1056 (2011). DOI 10.1214/10-AOAS395.
- 580 [22] Loo, P. V. *et al.* Allele-specific copy number analysis of tumors. *PNAS* **107**, 16910–16915 (2010). DOI
581 10.1073/pnas.1009843107.
- 582 [23] Dulak, A. M. *et al.* Gastrointestinal Adenocarcinomas of the Esophagus, Stomach, and Colon Exhibit
583 Distinct Patterns of Genome Instability and Oncogenesis. *Cancer Res.* **72**, 4383–4394 (2012). DOI
584 10.1158/0008-5472.CAN-11-3893.
- 585 [24] Aran, D., Sirota, M. & Butte, A. J. Systematic pan-cancer analysis of tumour purity. *Nat. Commun.* **6**,
586 1–11 (2015). URL <http://dx.doi.org/10.1038/ncomms9971>. DOI 10.1038/ncomms9971.
- 587 [25] Diehl, F. *et al.* Detection and quantification of mutations in the plasma of patients with colorectal tumors.
588 *PNAS* **102**, 16368–16373 (2005). DOI 10.1073/pnas.0507904102.
- 589 [26] Macintyre, G., Ylstra, B. & Brenton, J. D. Sequencing Structural Variants in Cancer for Precision Ther-
590 apeutics. *Trends Genet.* **32**, 530–542 (2016). URL <http://dx.doi.org/10.1016/j.tig.2016.07.002>.
591 DOI 10.1016/j.tig.2016.07.002.
- 592 [27] Martincorena, I. *et al.* High burden and pervasive positive selection of somatic mutations in normal human
593 skin. *Sci.* **348**, 880–886 (2015). URL <http://www.sciencemag.org/cgi/doi/10.1126/science.aaa6806>.
594 DOI 10.1126/science.aaa6806.
- 595 [28] Wand, M. *KernSmooth: Functions for Kernel Smoothing Supporting Wand & Jones (1995)* (2015). URL
596 <https://CRAN.R-project.org/package=KernSmooth>. R package version 2.23-15.
- 597 [29] Stephens, M. Dealing with label switching in mixture models. *J. Royal Stat. Soc. B* **62**, 795–809 (2000).
- 598 [30] Munkres, J. Algorithms for the Assignment and Transportation Problems. *J. Soc. for Ind. Appl. Math.* **5**,
599 32–38 (1957). URL <http://epubs.siam.org/doi/10.1137/0105003>. DOI 10.1137/0105003.
- 600 [31] Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinforma.* **25**, 2078–2079 (2009). DOI
601 10.1093/bioinformatics/btp352. [1006.1266v2](https://doi.org/10.1093/bioinformatics/btp352).
- 602 [32] Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinforma.*
603 **25**, 1754–1760 (2009). DOI 10.1093/bioinformatics/btp324. [1303.3997](https://doi.org/10.1093/bioinformatics/btp324).
- 604 [33] Langmead, B., Trapnell, C., Pop, M. & Salzberg, S. L. Ultrafast and memory-efficient alignment of short
605 DNA sequences to the human genome. *Genome Biol.* **10** (2009). DOI 10.1186/gb-2009-10-3-r25.
- 606 [34] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical
607 Computing, Vienna, Austria (2013). URL <http://www.R-project.org/>.

- 608 [35] Eddelbuettel, D. & Sanderson, C. RcppArmadillo: Accelerating R with high-performance C++ linear
609 algebra. *Comput. Stat. Data Analysis* **71**, 1054–1063 (2014). DOI 10.1016/j.csda.2013.02.005.
- 610 [36] Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows-Wheeler transform. *Bioinforma.*
611 **26**, 589–595 (2010). DOI 10.1093/bioinformatics/btp698. [1303.3997](https://doi.org/10.1093/bioinformatics/btp698).

612 Figure legends

613 Figure 1

614 Comparison of conliga method with ASCAT and QDNAseq. (a) Total copy number profile determined by
615 ASCAT from HC WGS data for sample OAC2, showing all copy number segments. (b) Relative copy number
616 profile determined by QDNAseq from LC WGS data for sample OAC2, showing all 15 Kbp bins. (c) Total
617 copy number profile determined by ASCAT from HC WGS data for sample OAC2, showing ASCAT's copy
618 number calls at the intersection of ASCAT's called regions and FAST-SeqS loci. (d) Relative copy number
619 profile determined by conliga from FAST-SeqS data for sample OAC2, at the intersection of ASCAT's called
620 regions and FAST-SeqS loci. (e) Comparison of \log_2 relative copy number calls from 11 samples between conliga
621 and ASCAT (top) and QDNAseq and ASCAT (bottom). All RCN calls at the intersection of ASCAT's called
622 regions, QDNAseq 15Kb bins and FAST-SeqS loci in all 11 OAC samples are shown as points. (f) Distribution
623 of differences between ASCAT RCN calls and conliga RCN estimates for 11 OAC samples (top) and ASCAT
624 RCN calls and QDNAseq RCN estimates for 11 OAC samples (bottom). (g) Comparison of performance at gene
625 level resolution between ASCAT and conliga (top) and ASCAT and QDNAseq (bottom). The values represent
626 the weighted mean of RCN calls at each gene for each of the 11 OAC samples (Methods).

627 Figure 2

628 Comparing the performance of SCNA detection in low tumor purity samples and determining the limit of
629 detection. (a) left column: relative copy number calls by conliga at different dilutions of sample OAC3, compared
630 to ASCAT relative copy number profile (top left), discrete copy number states are colored with a gradient (light
631 green to purple), highlighting regions with differing SCNAs. right column: relative copy number calls by
632 QDNAseq at different dilutions of sample OAC3, compared to ASCAT relative copy number profile (top right).
633 (b) The number of copy number states detected by conliga in each of eight OAC samples at differing purity
634 levels. The limit of detection is determined by the lowest purity level in which more than one copy number
635 state is detected.

636 Supplementary Figure 1

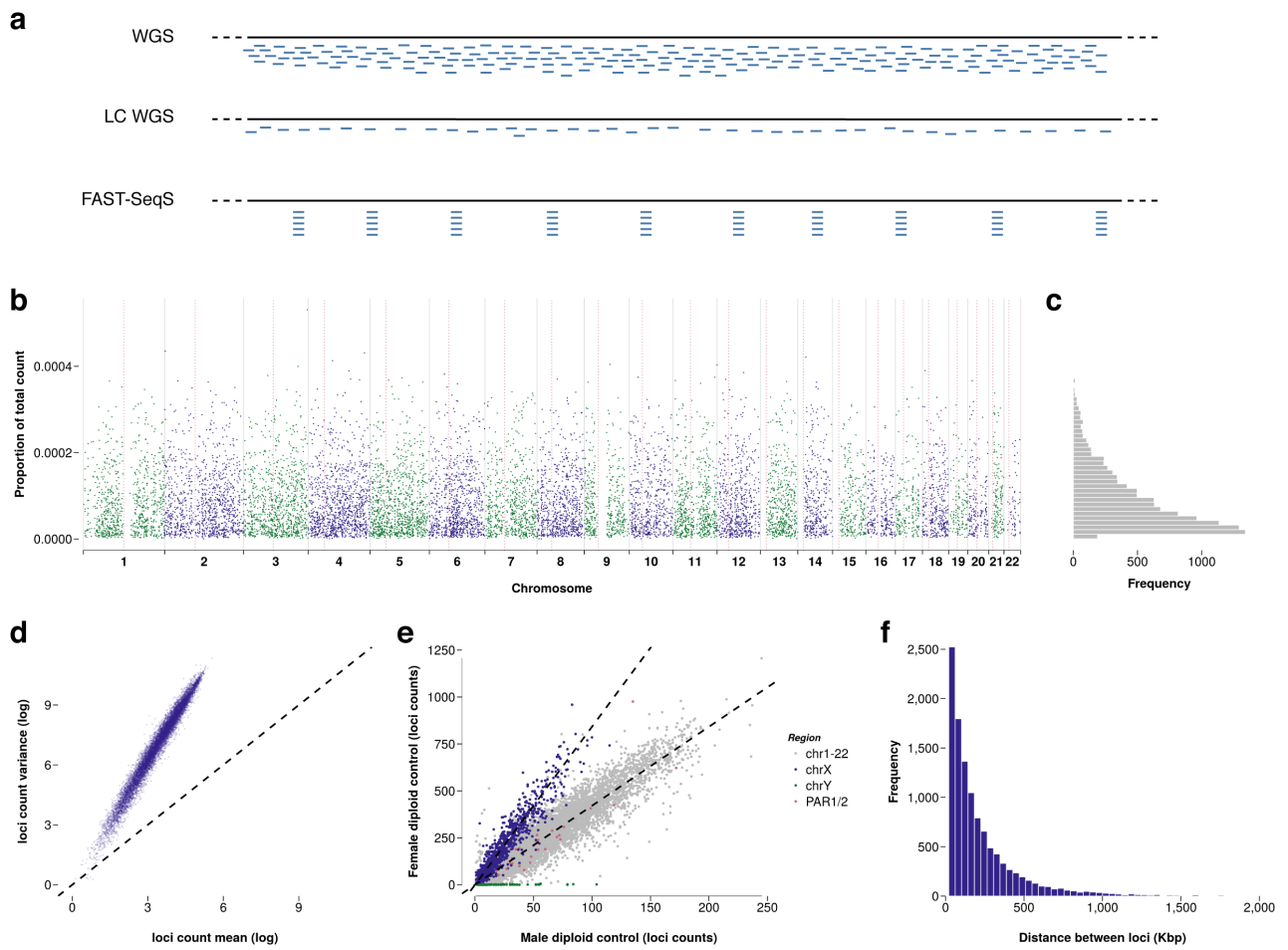
637 Aspects of FAST-SeqS data. (a) a graphical representation of the different approaches to sequencing for the
638 purposes of SCNA profiling; high-coverage WGS (top), low-coverage WGS (middle), FAST-SeqS (bottom). (b)

639 The proportion of reads obtained at each locus in chr1-22 for control sample (NORM1). (c) Histogram of the
640 proportion of reads obtained at each locus across in chr1-22 for control sample NORM1. (d) log mean vs log
641 variance for each locus in control samples. (e) A male control sample (NORM2) counts plotted against a female
642 control sample (NORM1) counts, showing a relative doubling of count proportions in chrX for the female control
643 sample vs male and absence of counts from chrY in the female sample. (f) Histogram of distances between loci
644 with a mean distance of approximately 200Kbp between loci.

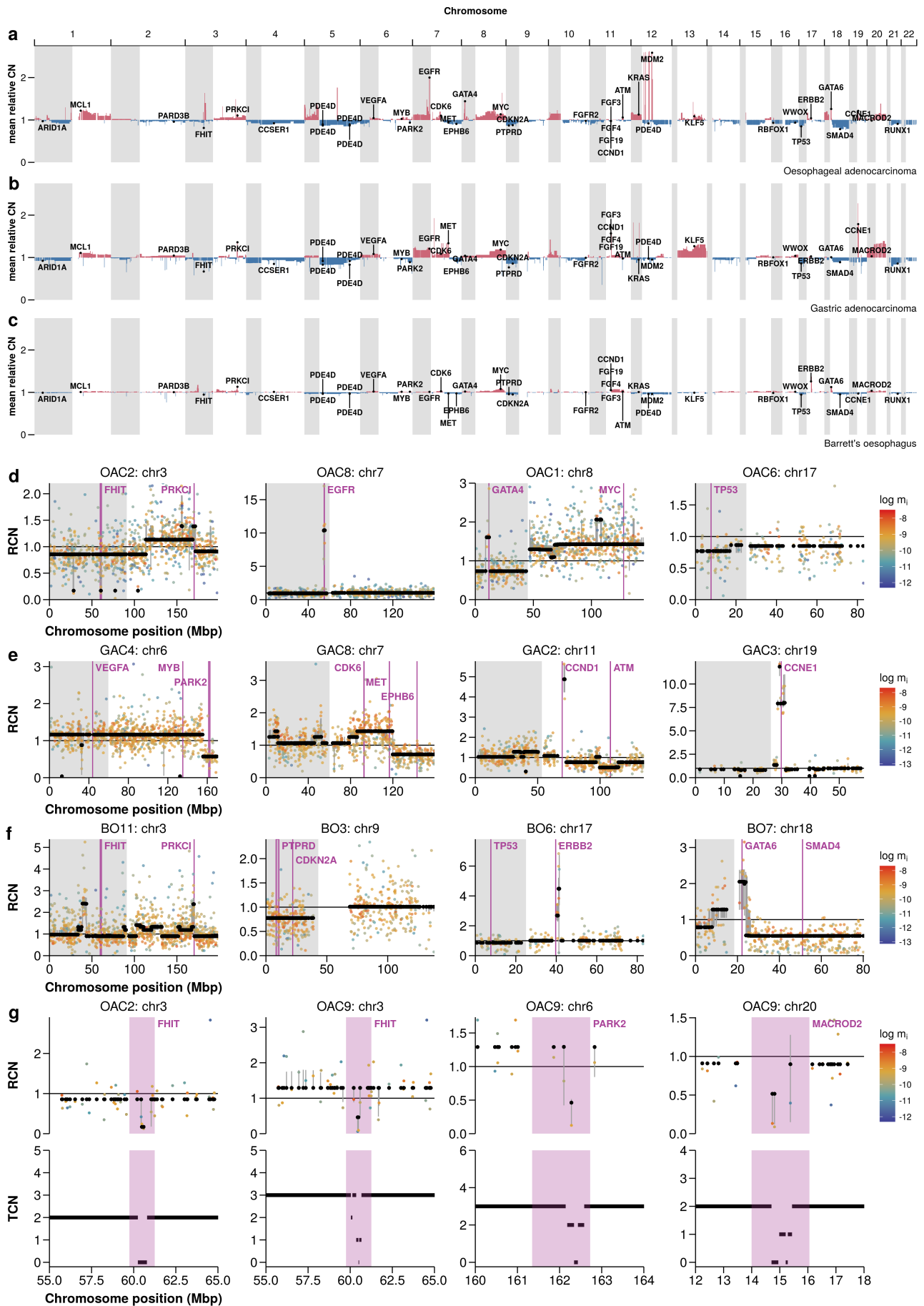
645 **Supplementary Figure 2**

646 Copy number profile summary of patient cohorts used in this study. (a) Mean relative copy number profile for
647 11 oesophageal adenocarcinoma samples. (b) Mean relative copy number profile for 8 gastric adenocarcinoma
648 samples. (c) Mean relative copy number profile for 16 Barrett's oesophagus samples, with varying levels of
649 dysplasia. (d)-(f) Examples of relative copy number profiles for various chromosomes from different samples for
650 OAC, GAC and BO respectively. Black points represent the maximum a posteriori (MAP) relative copy number
651 for each locus, the colored points represent the proportion of reads expected in a control sample (log), with
652 red representing a high proportion and blue representing a low proportion, grey lines represent 90% credible
653 intervals. (g) Zoomed-in regions of chromosomes 3, 6 and 20 showing intra-gene deletion of FHIT, PARK2 and
654 MACROD2. conliga results (top) with comparison to ASCAT (bottom).

655 **Supplementary Figures**



Supplementary Figure 1



Supplementary Figure 2