

Genomic evidence supports a clonal diaspora model for metastases of esophageal adenocarcinoma

DOI:

[10.1038/s41588-019-0551-3](https://doi.org/10.1038/s41588-019-0551-3)

Document Version

Accepted author manuscript

[Link to publication record in Manchester Research Explorer](#)

Citation for published version (APA):

Noorani, A., Li, X., Goddard, M., Crawte, J., Alexandrov, L. B., Secrier, M., Eldridge, M. D., Bower, L., Weaver, J., Lao-Sirieix, P., Martincorena, I., Debiram-Beecham, I., Grehan, N., MacRae, S., Malhotra, S., Miremadi, A., Thomas, T., Galbraith, S., Petersen, L., ... Fitzgerald, R. C. (2020). Genomic evidence supports a clonal diaspora model for metastases of esophageal adenocarcinoma. *Nature Genetics*. <https://doi.org/10.1038/s41588-019-0551-3>

Published in:

Nature Genetics

Citing this paper

Please note that where the full-text provided on Manchester Research Explorer is the Author Accepted Manuscript or Proof version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version.

General rights

Copyright and moral rights for the publications made accessible in the Research Explorer are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Takedown policy

If you believe that this document breaches copyright please refer to the University of Manchester's Takedown Procedures [<http://man.ac.uk/04Y6Bo>] or contact uml.scholarlycommunications@manchester.ac.uk providing relevant details, so we can investigate your claim.



1 **1. Extended Data**

2

Figure #	Figure title One sentence only	Filename This should be the name the file is saved as when it is uploaded to our system. Please include the file extension. i.e.: <i>Smith_ED Fig1.jpg</i>	Figure Legend If you are citing a reference for the first time in these legends, please include all new references in the Online Methods References section, and carry on the numbering from the main References section of the paper.
Extended Data Fig. 1	Flowchart describing key steps taken to construct phylogenetic trees	Noorani et al_Extended Data1_201911 18.tif	A. The details of phylogenetic tree reconstruction is further elaborated in Supplementary methods, Mutation clustering and phylogenetic tree construction (p.25).
Extended Data Fig. 2	Phylogenetic tree construction for example case S3	Noorani et al_Extended Data2_201911 18.tif	1) Battenberg algorithm to determine total copy number (purple line) and minor allele (blue line). Y-axis =number of chromosome copies, X-axis= chromosome and position. The average ploidy, aberrant cell fraction (cellularity) and goodness of fit to the model are shown for each sample, Primary E1, E2, Lymph node L1 and Distant metastasis D1. The goodness of fit is a measure of the amount of the genome with clonal, rather than subclonal copy number states. D1 has a subclonal mix of different copy number states resulting in noninteger total copy number, for example on chromosome 2, resulting in a goodness of fit below 100%. 2) Bayesian Dirichlet Process to cluster SNVs based on CCF in each sample. The density plots show the posterior probability of a mutational cluster, these are produced for every pair of samples and selected plots are shown High density at CCF of (0,0) indicates subclones that are not present in the pair of samples shown in a particular plot. 3) Clustering of results – Clusters are identified as local maxima in the posterior density. The table shows the number of SNVs assigned to each cluster, and their associated CCFs. 4) Unscaled Tree construction using the sum rule and crossing rule as detailed in Supplementary Methods p25. 5) Final Tree -The tree is drawn as seen in Figures 2 and Extended Data2, branch lengths are proportional to

			the number of SNVs assigned to each subclone. Scales vary on a per case basis depending on the total number of SNVs, in order to fit cases on one figure. Trees are annotated with the gene names of known drivers, and the colour of each branch represents a trunk (pink), branch (purple) or leaf (yellow). The grey circles represent clones and subclones and their CCFs are shown in Supplementary Table5 and 6.
Extended Data Fig. 3	Phylogenetic trees of cases in cohort with only nodal or distant organ disease, as derived from H-WGS	Noorani et al_Extended Data3_201911 18.tif	E=esophagus, D=distant organ, L=lymph node, B= Barrett's. For precise anatomical locations, refer to Supplementary Table3 and 4. MRCA=most recent common ancestor. Pink=trunk (shared events), Purple=branch (shared by more than one sample), Yellow=leaf (unique to one sample). Grey dots at the end of the lines (truncal, branches or leaves) represent subclones or clones, whose CCFs are shown in Supplementary Tables5 and 6. Trees are annotated with key driver events as identified from the literature ^{14,16,19} . Black=point mutations, Red=copy number alterations, purple= structural variants. The adjacent scales are relative to the number of SNVs in that particular case and hence constructed on a case by case basis.
Extended Data Fig. 4	Structural variation of 18 metastatic esophageal adenocarcinoma cases	Noorani et al_Extended Data4_201911 18.tif	a. Similarity matrix based clustering for all SVs 122 genomes across 19 cases. SVs were deemed to refer to the same rearrangement event across cases if their corresponding breakpoint locations fell within a window of maximum 50 bp. The individual sample types are shown as a separate row on the x axis with the color key depicting the sample type. The purple scale indicates the number of shared SVs. (L=lymph nodes; M=metastasis; T=tumor). b. Histogram showing the percentage of rearranged genes that are concordant, unique to tumors and unique to metastases. Two-tailed Welch test P=0.2674 demonstrating no overall difference between total number of SVs in primary, local lymph nodes and distant metastases c. Stacked bar charts showing the composition of various SVs in each sample on a per patient basis INV= inversion, ME= mobile element, BND= translocation DEL=deletion, DUP=duplication, INS= insertion.

Extended Data Fig. 5	Random simulation model for S-WGS cluster detection	Noorani et al_Extended Data5_201911 18.tif	The number of mutations detected correlates strongly with the CCF of the cluster (Pearson $r=0.992$, $n=100$). Number of mutations in each cluster =1000.
Extended Data Fig. 6	Correlation of fraction of mutations detected with CCF as a function of cluster size using simulated S-WGS data	Noorani et al_Extended Data6_201911 18.tif	Pearson correlation coefficient is above 0.97 for clusters with 200 or more mutations.
Extended Data Fig. 7	Bar chart demonstrating the Pearson correlation coefficient of VAF at 1xWGS and High Depth Resequencing (n=33)	Noorani et al_Extended Data7_201911 18.tif	
Extended Data Fig. 8	Detection of Selection in subsets of mutations	Noorani et al_Extended Data8_201911 18.tif	SNVs and indels from all cases (n=18) were aggregated into 4 different subsets: clonal = variants found in the MRCA (n=378453); subclonal = variants not found in the MRCA (n=516136); pre-diaspora = variants found above the diaspora founder clone in the phylogenetic tree (n=313545); post-diaspora = variants found in the diaspora founder or in clones below the founder in the phylogenetic tree (n=295316). Within each subset, dN/dS analysis was performed separately on: missense variants; truncating variants. Bars show maximum likelihood estimates of dN/dS values, with values greater than 1 (dashed line) indicating positive selection. Vertical lines = 95% confidence intervals, estimated using Wald test.
Extended Data Fig. 9	Percentage of truncal and branch clusters in tissue from earlier time-points	Noorani et al_Extended Data9_201911 18.tif	Stacked horizontal bar chart representing the percentage of truncal and branch clusters present in tissue from earlier time-points on the x-axis and the Case ID on the y-axis. P1 diagnosis* is a frozen sample, while the rest are FFPE. Blue = truncal, maroon = branch, grey = not present. The number of clusters (n) is demonstrated for each case.
Extended Data Fig. 10	ctDNA analysis from historical	Noorani et al_Extended	Digital PCR traces of mutant allele fraction for TP53 on the Y-axis and days from

	plasma samples	Data10_20191118.tif	diagnosis on the X-axis, and grey areas indicate periods of therapy. Where subclones and clones are seen at 1xWGS on plasma, they are highlighted on the 50x phylogenetic tree (coloured blue). The samples in which these subclones and clones are present in are shown in Supplementary Table3. There was no TP53 data for S3 as it was wild type for TP53 mutations. Copy number traces for P1 are shown, with the arrow demonstrating an area of MET amplification.
--	----------------	---------------------	---

3
4

5 2. Supplementary Information:

6 A. Flat Files

7

Item	Present?	Filename This should be the name the file is saved as when it is uploaded to our system, and should include the file extension. The extension must be .pdf	A brief, numerical description of file contents. i.e.: <i>Supplementary Figures 1-4, Supplementary Note, and Supplementary Tables 1-4.</i>
Supplementary Information	Yes	Noorani et al_Supplementar yInformation_20191118.pdf	Supplementary Figures 1-5, Supplementary Tables 1-16, Supplementary Note
Reporting Summary	Yes	Reportingsummary.pdf	

8

9 B. Additional Supplementary Files

Type	Number If there are multiple files of the same type this should be the numerical indicator. i.e. "1" for Video 1, "2" for Video 2, etc.	Filename This should be the name the file is saved as when it is uploaded to our system, and should include the file extension. i.e.: <i>Smith_Supplementary Video 1.mov</i>	Legend or Descriptive Caption Describe the contents of the file
Supplementary Table	1	Noorani et al_SupplementaryTables_Excel_20191118.xlsx	Excel Spreadsheets for Supplementary Table 5, 6, 7, 8, 9, 11, 16

10

11 **Genomic evidence supports a clonal diaspora model for metastases of esophageal**
12 **adenocarcinoma**

13
14 Ayesha Noorani¹, Xiaodun Li¹, Martin Goddard², Jason Crawte¹, Ludmil B. Alexandrov³, Maria
15 Secrier⁴, Matthew D. Eldridge⁴, Lawrence Bower⁴, Jamie Weaver¹, Pierre Lao-Sirieix¹, Inigo
16 Martincorena⁵, Irene Debiram-Beecham¹, Nicola Grehan¹, Shona MacRae¹, Shalini
17 Malhotra⁶, Ahmad Miremadi⁶, Tabitha Thomas⁷, Sarah Galbraith⁸, Lorraine Petersen⁷,
18 Stephen D. Preston², David Gilligan⁹, Andrew Hindmarsh¹⁰, Richard H. Hardwick¹, Michael R.
19 Stratton⁵, David C. Wedge^{11, 12*} and Rebecca C. Fitzgerald^{1*}

20 ¹MRC Cancer Unit, University of Cambridge, Biomedical Campus, Cambridge, CB2 0XZ, UK

21 ²Department of Histopathology, Papworth Hospital NHS Trust, Cambridge, CB23 3RE, UK

22 ³Theoretical Biology and Biophysics (T-6), Los Alamos National Laboratory, New Mexico, 87545, USA

23 ⁴Cancer Research UK Cambridge Research Institute, Cambridge, CB2 0RE, UK

24 ⁵Wellcome Trust Sanger Institute, Cambridge, CB10 1SA, UK

25 ⁶Department of Histopathology, Cambridge University Hospitals NHS Foundation Trust, Cambridge, CB2 0QQ,
26 UK

27 ⁷Arthur Rank Hospice Charity, Cambridge, CB22 3FB, UK

28 ⁸Department of Palliative Care, Cambridge University Hospitals NHS Foundation Trust, Cambridge, CB2 0QQ,
29 UK

30 ⁹Oncology Centre, Cambridge University Hospitals NHS Foundation Trust, Cambridge, CB2 0QQ, UK

31 ¹⁰Oesophago-Gastric Centre, Cambridge University Hospitals NHS Foundation Trust, Cambridge, CB2 0QQ, UK

32 ¹¹Big Data Institute, University of Oxford, Oxford, OX3 7LF, UK

33 ¹²Oxford NIHR Biomedical Research Centre, Oxford, OX4 2PG, UK

34

35

36 *Correspondence to Rebecca Fitzgerald rcf29@mrc-cu.cam.ac.uk or David Wedge david.wedge@bdi.ox.ac.uk

37 **Abstract (95 words)**

38 The poor outcomes in esophageal adenocarcinoma (EAC) prompted us to interrogate the
39 pattern and timing of metastatic spread. Whole genome sequencing and phylogenetic
40 analysis of 388 samples across 18 EAC cases demonstrated in 90% of cases that multiple
41 subclones from the primary tumor spread very rapidly from the primary site to form
42 multiple metastases, including lymph nodes and distant tissues, a mode of dissemination
43 that we term 'clonal diaspora'. Metastatic subclones at autopsy were present in tissue and
44 blood samples from earlier time-points. These findings have implications for our
45 understanding and clinical evaluation of EAC.

46

47 **Introduction**

48 Metastatic spread to distant sites accounts for the majority of cancer deaths¹.
49 Understanding the anatomical extent of disease is essential to determine the optimum
50 treatment strategy. This is challenging since cancer continually evolves at a microscopic
51 scale, often beyond the resolution of clinical imaging techniques. Furthermore, the patterns
52 of metastatic spread are often unpredictable in terms of time-course and anatomical
53 location. Treatments may therefore be unnecessarily toxic (e.g. radical lymphadenectomy
54 and chemotherapy) or insufficiently aggressive, leading to high recurrence rates²⁻⁴.

55 Esophageal cancer is the sixth most common cause of cancer-related death worldwide and
56 the current median survival time is still <1 year⁵. Incidence rates for esophageal
57 adenocarcinoma (EAC) have risen sharply and it is now the predominant subtype in
58 developed countries. Prognosis is highly variable for EAC patients as shown by the wide
59 range of 5-year survival (18-47% with lymph node involvement), making it difficult to advise
60 patients when embarking on a long course of grueling treatment^{2,6}.

61 Theoretical and experimental studies attempt to understand how tumor cell populations
62 respond to selective pressures over time⁷. A number of models of tumor evolution have
63 been proposed, including linear, branching, neutral and punctuated evolution, but the
64 extent to which these are specific to a given cancer type or co-occur is controversial^{8,9}.
65 Genome sequencing studies have attempted to delineate different models of evolution¹⁰.
66 However, many of these studies have focused solely on evolution within the primary site,
67 and knowledge of how genetic diversity emerges during metastasis remains limited. The

68 lack of understanding is in part due to the practical challenge of collecting multiple samples
69 over space and time from advanced stage cancer patients.

70 To better understand the evolution of EAC, we designed a prospective study with extensive
71 sampling over time including samples from diagnosis, surgery and at warm autopsy (Figure
72 1). We used whole genome sequencing (WGS) at high depth (50x), to identify mutations,
73 and at shallow (1x) coverage, to track known variants, to interrogate the clonal architecture
74 across time and space.

75

76 **Results**

77 **Genomic architecture of 18 cases**

78 Eighteen cases were included in the study and the clinical demographics of these cases are
79 shown in Supplementary Table 1 and 2, with details of the individual samples given in
80 Supplementary Table 3 and 4. In the first part of the study (Figure 1a, Extended Data Fig.
81 1,2) we used 50x WGS to construct a phylogenetic tree for each case, to understand the
82 relationship between the primary and metastatic tumors (Figure 2, Extended Data Fig. 3 ,
83 Supplementary Figure 1, Supplementary Table 3, 4). Mutation clustering was performed,
84 and the fractions of tumor cells carrying each set of mutations (Cancer Cell Fraction, CCF)
85 within each sample were used to determine: 1) the clonal and sub-clonal architecture of
86 each tumor (subclonal CCF <95%, clonal CCF \geq 95%); 2) the hierarchy of events; and 3) the
87 distance of these sub-clonal or clonal clusters from the most recent common ancestor
88 (MRCA) (Figure 1a, Extended Data Fig. 1,2). The CCF and number of single nucleotide
89 variants (SNVs) associated with each clone and subclone are shown in Supplementary Table
90 5 and 6, as is the tumor purity of each sample using the Battenberg algorithm¹¹, in
91 Supplementary Table 7 and the confidence intervals of the clonal and subclonal CCFs in
92 Supplementary Table 8. Detailed information on experimental design is provided in the Life
93 Sciences Reporting Summary.

94 These analyses enabled us to construct phylogenetic trees (Methods). In all cases we
95 observed a long trunk compared to the rest of the tree (median 19,034 SNVs, IQR 11,299-
96 63,908), consistent with previous studies in EAC^{12,13}. The median size of clonal or subclonal
97 clusters across all cases was 3,069 SNVs (IQR 1332-63908) and only 2/157 contained fewer
98 than 200 SNVs (S1_3 and P5_11), Extended Data Fig. 3 and Supplementary Table 6.

99 The key driver events^{14,15} are depicted on each phylogenetic tree (Figure 2 and Extended
100 Data Fig. 3). The events identified as most frequent in previous studies occurred in the
101 trunks of the phylogenetic trees, consistent with their previous classification as drivers. *TP53*
102 was mutated in the trunk of 16 out of 18 cases, consistent with our knowledge of the
103 disease^{14,16-19}. Amplifications (gene names in red) were often truncal, but also observed on
104 the branches of the phylogenetic tree, providing evidence of divergence during later
105 evolutionary stages (Figure 2, Extended Data Fig. 3). The majority of events in driver genes
106 were copy number alterations (CNAs) rather than SNVs or InDels (Figure 2, Extended Data
107 Fig. 3)^{14,19,20}. There was no significant difference in the overall number of structural variants
108 between primary and metastatic samples ($p=0.41$, generalized linear model; Extended Data
109 Fig. 4b). However, a larger proportion of structural variants in metastatic samples were
110 retro-transpositions of mobile elements than in the primary samples ($p=0.045$, Extended
111 Data Fig. 4c). This contrasts with pancreatic cancer, where deletions and fold-back
112 inversions are more common in metastases, and breast cancer where tandem duplications
113 dominate²¹. Interestingly, the high rate of L1 transposon activity in EAC has recently been
114 associated with high activity in the germline²². Our results suggest a further increase in L1
115 activity in metastatic EAC. Furthermore, the proportion of structural variants found uniquely
116 in metastases or in primary sites was higher than that of SNVs (Figure 2, Extended Data 4a),
117 suggesting an increase in genomic instability in later stages of the disease. However, it
118 cannot be ruled out that some structural variants have not been identified in every sample
119 as a result of lower sensitivity in the detection of structural variants than SNVs.

120 Across the eighteen cases, 8 mutational signatures were observed, consistent with previous
121 studies²³⁻²⁶ (Figure 3a), with varying prevalence across the cases. None of the signatures that
122 we observe in patients in our cohort who had oncologic therapy have been associated with
123 treatment with alkylating antineoplastic agents²⁷, platinum therapy²⁸ or radiation therapy²⁹.

124

125 **Early seeding of oligometastases**

126 Ten of eighteen patients (S3, S4, P1-4, P6, P8-10) had both nodal and solid organ
127 metastases, allowing a direct comparison of the genomic architecture between different
128 metastatic sites (Figure 2).

129 In four of these ten cases, an isolated clone or subclone confined to 1 or 2 distant
130 metastases, i.e. an oligometastasis, depicted as a dashed black node on the first branch of

131 the phylogenetic tree, shared the highest congruence to the MRCA, (P1, P4, P10, S3 in
132 Figure 2; Subclones P1_2, P4_3, P10_2, S3_2 in Supplementary Table 5). In P1, this clone
133 (P1_2) was observed only in the primary tumor and a pleural metastasis. In S3 and P4, the
134 clone involved in this isolated seeding was identified at a single distant site and not in the
135 primary tumor (S3_2: liver metastasis (D1), P4_3: para-aortic lymph node (L3)). In P10, the
136 early seeding clone (P10_2) was shared between a distant para-aortic node and a sub-clonal
137 metastasis in the right hemi-diaphragm. The subclones associated with these isolated
138 seeding events showed little divergence from the MRCA across these 4 cases (median 1,913
139 SNVs, range 832-8,591), suggesting early seeding to distant metastases. Notably, in P9 a
140 subclone (P9_10, Supplementary Table 5) was found in a premalignant area of Barrett's
141 esophagus and a pleural metastasis but not in any of four areas of the primary tumor
142 subject to 50x WGS. This subclone lineage shares no variants with the main lineage and
143 appears to be an independent second cancer (Figure 2).

144 **A single clone gives rise to multiple metastatic sites**

145 A striking observation was that 9/10 cases had a clone (outlined in red on the phylogenetic
146 tree in Figure 2) that was followed by dispersion of multiple subclones from the primary to
147 discrete metastatic sites, resulting in a model of metastasis that we term 'clonal diaspora'.
148 In most cases, this dispersion was visually stellate in nature, this being defined as a feature
149 of a phylogenetic tree involving 3 or more branches leading from a single founder clone (see
150 details in Discussion). The subclones forming diasporas were located in both primary and
151 metastatic tissue in eight cases (P1, P2, P3, S4, P4, P6, P8, P10) and in P9 were unique to
152 metastases (Figure 2). The only two cases lacking a stellate pattern on the phylogenetic tree
153 were P10 and S3. The latter is a non-autopsy case with limited tissue sampling and the early
154 distant seeding in this case is consistent with a pattern of parallel evolution (Figure 2).

155

156 **Subclonal spread is not constrained by location or tissue**

157 In the second step of the study we tracked the spread of metastases across a wider range of
158 lymph node and distant tissue sites by performing 1x WGS in a further 248 tissue samples
159 from 6 autopsy cases (Figure 1a,c). We did not call new mutations, as this would not be
160 possible at 1x sequencing, but used this method to detect the spread of clones and
161 subclones previously identified using 50x WGS (bioinformatic validation of methods in

162 Extended Data Fig. 5 and 6, Supplementary Note; wet lab validation in Extended Data Fig. 7,
163 Supplementary Table 9). The samples used in this part of the study are outlined in
164 Supplementary Table 10. The median size of subclonal and clonal clusters (identified
165 previously at 50x WGS) that we aimed to detect using 1x WGS was 3,784 (IQR 1,966-49,955).
166 Sample sites were grouped according to their similarity based on the presence of subclones
167 and clones previously detected with 50x WGS (Supplementary Note). The resulting groups
168 of samples are color coded and numbered, and each sample site, colored by group, is shown
169 on the adjacent body map (Figure 4, see also Supplementary Note). Notably, the samples
170 that grouped together based on shared clonal origins were widely dispersed anatomically.
171 Four out of six cases with extensive spatial sampling (Figure 4) had liver metastases
172 evaluated and three of these contained samples that were more similar to local lymph node
173 metastases than neighboring liver metastases (P4, P6, P8 but not P10). The high number of
174 groups within the liver (up to four) suggested seeding by multiple subclones (seen in P4, P6,
175 P8), whereas the single group in the liver of P10 (orange, group 3) indicated seeding by a
176 common progenitor or a set of closely related cells.

177 A comparison of lymph node location and genomic contiguity showed no evidence of
178 tropism, i.e. genomically similar lymph nodes did not occupy nearby anatomical locations.
179 Lymph nodes above and below the diaphragm were frequently seeded from common
180 events (P2: groups 1, 3; P4: groups 5, 6; P6: group 5; P8: groups 2, 3,5, 6; P10: group 4), at
181 odds with a progression from local to distant nodes. Similarly, a comparison of lymph node
182 and solid organ metastases showed scant evidence for tropism, with the exception of P1
183 (Supplementary Note). This patient underwent surgical resection and subsequently had
184 metastatic disease recurrence. In this cancer, separate subclones seeded lymph node and
185 pleural metastases (Figure 2, 4). Notably, the distant metastasis (D1) was an early branching
186 oligometastasis whereas the lymph nodes (L1, L2) constituted the later diaspora event
187 (black and red circles, respectively, in Figure 2).

188 We further traced regions of the primary tumor at autopsy that had similar subclonal
189 compositions to each of the metastases, shown as adjacent tumor maps (Figure 4, bottom
190 left of each case). Subclones occupied spatially distinct areas in the primary tumor.

191 We also looked for driver amplifications post MRCA or post diaspora on a per case basis and
192 identified selection in 6/10 cases. However, this is likely to be an under-estimate, since there
193 may be non-copy number drivers present in additional cases. The ratio of non-synonymous

194 to synonymous SNVs (dN/dS) was analyzed across all cases in order to assess the presence
195 or absence of positive selection³⁰. Results indicated positive selection in both clonal and
196 subclonal genomes, albeit with lower levels of selection within subclones (Extended Data
197 Fig. 8).

198

199 **Metastatic spread is rapid in EAC**

200 To examine the timing and speed of metastatic spread we analyzed base substitution
201 mutational signatures, particularly the aging signature which features a predominance of
202 C>T transition in the NpCpG trinucleotide context (Figure 1a, Figure 3).

203 Signature 1 arises from the spontaneous or enzymatic deamination of methylated cytosines,
204 which is an endogenous process that occurs continuously in both healthy and cancerous
205 cells. This has been shown to act as a molecular clock^{27,31-35}, and was therefore used here as
206 a method to examine the temporal relationship between metastases. Using a previously
207 described method for deconvolving mutational signatures³⁵, we observed that signature 1
208 was present in the trunk but absent in all subclones that constituted diaspora (following the
209 red parental clone in Figure 2) for P2, P4, P6, P9, P10, S4 and it was significantly reduced for
210 P1 (21% to 3%) and P3 (16% to 9%) (Wilcoxon signed rank test $p=0.039$, Figure 3c). To
211 account for the possibility that the number of signature 1 mutations in branch subclones
212 was below the resolution of our deconvolution methods, we also identified the number of
213 mutations with the characteristic feature of signature 1, i.e. C>T mutations in a CpG context.
214 To estimate the time of appearance of diaspora, we compared the number of these
215 characteristic mutations that occurred along the trunk to the parental red clone marking the
216 onset of diaspora with those that occurred on the longest branch leading from this point.
217 The median proportion of such mutations occurring prior to the onset of diaspora was 0.911
218 (Figure 3b). Thus, in the majority of cases one might deduce that little time has elapsed
219 between the appearance of the cell that is ancestral to disseminating cells and the individual
220 cells that seeded each of the metastases. With the exception of P8, the proportion of
221 mutations attributed to signature 1 was significantly lower after the parental (red) clone on
222 the phylogenetic tree ($p < 9.1 \times 10^{-5}$, Chi-squared test across all cases; Figure 3c) suggesting
223 an increase in the activity of other processes in later evolutionary stages (Supplementary
224 Table 11). Of note, there was an increase in the proportion of signature 3 in subclonal SNVs

225 compared to clonal SNVs (Wilcoxon signed rank test $p=0.019$, Figure 3b), suggesting failure
226 of DNA double strand break repair is predominantly a late-stage event in EAC.

227

228 **Early detection from diagnostic samples**

229 Next, we investigated eight cases (P1-4, P6, P8-10) for which the esophageal diagnostic FFPE
230 biopsy or surgical sample (primary tumor at resection for P1 and lymph node from surgery
231 for P9) were available, with a median time prior to autopsy of 12 months (range 5-30
232 months) (Figure 1). The diagnostic sample for P1 was snap frozen and sequenced to 50x
233 (Figure 2; highlighted with * in Extended Data Fig. 9), while 1x WGS was performed on the
234 remainder of the cases. Between 8% and 36% of the subclones and clones observed in
235 samples taken from autopsy were also present in the diagnostic samples (Supplementary
236 Note and Extended Data Fig. 9). In six cases, all subclones identified from the biopsy samples
237 were also found in the primary samples from autopsy. Two diagnostic endoscopic samples
238 from P4 also contained many of the mutations found in the lymph node L2 at autopsy,
239 which had not been previously identified in the primary tumor at autopsy (Figure 2,
240 subclone P4_17, Supplementary Table 5). Similarly, the biopsy sample from P10 contained a
241 substantial number of mutations from both the oligometastasis that seeded D2 and L4
242 (Supplementary Table 5, P10_2), and the lineage that later metastasized to multiple sites
243 (Figure 2). Notably, P4 and P10 had shorter survival times after diagnosis than the remaining
244 patients (5 and 4 months, respectively).

245

246 **Plasma sample analysis at autopsy and earlier time-points**

247 We assessed the clonal composition of circulating tumor DNA (ctDNA) at earlier time-points
248 in seven blood samples from five cases (Figure 1, Figure 5a,c; Extended Data Fig. 10,
249 Supplementary Table 12). Combined 1x WGS subclone/clone detection, copy number
250 aberrations and *TP53* fraction using digital PCR data are displayed for two of these cases (P6
251 and P10) in Figure 5a. Notably, P6 was a patient being treated with curative intent and had
252 no radiological evidence of distant nodal or organ metastases at the time of clinical staging.
253 However, at the time of diagnosis mutations from the truncal cluster and three subclonal
254 clusters later found in the metastases were already present in the plasma (Figure 5a) along
255 with amplifications in *MYC* and *GATA4*. Case S4 is noteworthy as the brain metastases (D1,
256 D2 in Figure 2) appeared to have originated from a subclone shared between the primary

257 and a local lymph node, both of which were removed at the time of surgery (Extended Data
258 Fig. 10c). However, mutations from the truncal cluster and four subclonal clusters were
259 already present in ctDNA prior to radiological recurrence.

260 In eight cases, plasma was available from rapid autopsy. One case (P3) failed wet lab SNV
261 validation and was hence removed from the SNV subclone analysis (Supplementary Note).
262 Analysis of ctDNA demonstrated that in all cases the truncal cluster from autopsy was also
263 represented in plasma (Figure 5c). In addition, mutations from between 0 and 7 subclonal
264 clusters were identified from plasma (Figure 5c). The ratio of mutations detected from each
265 subclone was very consistent between blood from earlier time points and autopsy (Pearson
266 r range [0.851, 0.994], maximum P-value 8.9×10^{-4}) and in 2 of 5 cases the proportion of
267 mutations detected was higher in the earlier sample, suggesting an opportunity for earlier
268 detection of heterogeneous cancer cell populations. Further, subclonal proportions
269 estimated from exome sequencing of plasma samples were highly correlated with those
270 from 1x WGS (Supplementary Table 9).

271 The majority of driver CNAs identified in the MRCA of each tumor from 50x WGS of tissue
272 samples were also identified in plasma both at autopsy and at earlier time-points (Figure
273 5a,b). In addition, MET amplification, which was not present in the MRCA in P1 (Figure 2),
274 was identified in plasma both at autopsy and an earlier time point (Extended Data Fig. 10a),
275 suggesting opportunities for early detection of metastatic subclones. Notably, however,
276 amplifications found only in oligometastases or in post-diaspora subclones from 50x
277 sequencing were not identified in plasma, despite many of them being detected in 1x
278 sequencing of tissue samples (Figure 5b). A plausible explanation for this observation is that
279 each of the many metastasizing subclones contributed insufficient material to the sum of
280 detected ctDNA to enable confident detection of CNAs.

281

282 **Discussion**

283 We have gathered multiple lines of evidence which suggest that, for the majority of EACs, a
284 complex mode of spread is operative. These lines of evidence can be summarized as follows
285 (Figure 6). We observe multiple subclones, each seeding multiple metastatic sites. These
286 subclones are frequently derived from a single parental clone, generally resulting in a
287 stellate pattern on the phylogenetic tree. Metastases in solid organs can bypass nodal
288 involvement and samples within solid organ sites frequently resemble distant metastases

289 more closely than neighboring metastases within the same organ, i.e. no tropism is
290 observed. All metastases appear to have spread directly from the primary site, with little or
291 no evidence of metastasis-to-metastasis seeding.

292 These features differ in some important respects from previously described models of
293 metastasis and we propose that they may constitute a distinct, additional model of
294 evolution. We suggest that this pattern be referred to as a 'diaspora', by extension of the
295 anthropological term to cancer³⁶. Within this context, it is associated with the observation
296 that multiple cell populations in metastatic sites are directly linked to the primary site of
297 origin and that individual subclones seed multiple tissue types, analogous to a diaspora
298 crossing multiple national boundaries.

299 A number of features were frequently associated with this phenomenon (Figure 6), with
300 nine of the cases (all except S3) displaying at least two of the four following features: i)
301 stellate pattern on the phylogenetic tree defined as three or more subclones emerging from
302 the founder clones; ii) lack of signature 1 mutations post MRCA or post-diaspora; iii) spread
303 of subclones to multiple organs of different type; iv) evidence for selection in post diaspora
304 genotypes.

305 Until recently the genomic architectures of metastatic samples have not been defined with
306 enough resolution to discern temporal or spatial patterns of metastatic spread. Several
307 distinct patterns are now emerging which are not necessarily mutually exclusive or cancer-
308 type specific. In pancreatic cancer, Yachida et al. demonstrated that distant organ seeding
309 was a late event consistent with a linear progression model²⁴. In prostate cancer, linear
310 progression is often succeeded by multiple waves of seeding³⁷. The same study further
311 demonstrated widespread subclonal evolution in metastases and metastasis-to-metastasis
312 spread, in keeping with the relatively long longevity of prostate cancer. Strikingly, a stellate
313 pattern was not observed in any of the cases in that study, despite using a similar design to
314 that used here.

315 In Supplementary Table 13 we compare the features of our proposed Diaspora model to the
316 previously posited linear³⁸ and parallel⁸ models. Whereas the linear model predicts that a
317 single subclone seeding lymph node sites is followed by transmission to distant organs, the
318 diaspora model posits simultaneous seeding of multiple sites directly from the primary.
319 Unlike the parallel model, the diaspora model implies that metastasis formation occurs after
320 the majority of evolution has occurred in the primary tumor, resulting in multiple subclones

321 found in common between primary and metastatic tumors. Lymphatic and distant
322 metastases in colon cancer have been shown to arise from independent subclones in the
323 primary tumor with disparate evolutionary trajectories³⁹. In contrast, in EAC we find that
324 individual subclones frequently seed both lymph node and distant organs suggesting that
325 disparate trajectories for nodal and solid organ metastases do not exist for this disease
326 (Figure 2, 3). Of note we acknowledge that, despite the extensive and systematic sampling
327 across all autopsy cases, further sampling may add further branches to our phylogenetic
328 tree, although this is unlikely to affect the diaspora event itself.

329 In common with the Big Bang Model proposed for colorectal cancer⁴⁰, our model predicts
330 the occurrence of highly branching phylogenies. However, the Big Bang Model proposes
331 neutral dynamics, whereas we observe strong evidence for selection in subclonal
332 populations in the form of dN/dS ratios and the occurrence of subclonal driver
333 amplifications (Figure 2, Extended Data Figure 8, Supplementary Figure 2). Moreover, the
334 clonal maps of the primary tumor demonstrate subclones that occupy spatially discrete
335 areas of the primary tumor (Figure 4), in contrast to the intermixed subclones predicted by
336 the Big Bang Model⁴⁰.

337 The sequence of events in metastatic progression may have clinical implications that require
338 further study (Supplementary Table 13). Clonal architecture in EAC defies anatomical
339 location of lymph node stations and distant sites, which is the current basis for the TNM
340 staging and determines whether curative therapy is appropriate. It has been suggested that
341 the high recurrence rate, 52% within one year, results from seeding of distant metastases
342 that are not detected at the time of diagnosis²⁶. This study provides molecular evidence for
343 this observation and highlights the need for different systemic approaches to disease
344 management, including consideration of more aggressive adjuvant therapy which is not
345 currently the mainstay of treatment⁴¹⁻⁴⁴. With advances in the sensitivity of ctDNA assays,
346 metastatic subclones may be detectable in the blood, helping to determine when systemic
347 therapy is required post-surgery and in detecting heterogeneity of acquired resistance⁴⁵.
348 Copy number variation in plasma may also be a future early detection strategy⁴⁶.

349

350 The occurrence of metastasis is a pivotal event in the life history of a cancer. Understanding
351 the drivers behind such an event would have potential relevance to patient stratification
352 and predicting and preventing metastatic spread⁴⁷. While we have identified many drivers

353 on the trunks of the trees, prior to diaspora (Figure 2), we cannot be certain which event, if
354 any, was the immediate trigger of diaspora in individual cases. In a number of cases,
355 diaspora was coincident with an increase in the proportion of signature 3 mutations,
356 associated with failure of DNA double-strand break-repair by homologous recombination
357 (Figure 3b). Our findings are in keeping with the failure of DNA repair driving the
358 appearance of genomic heterogeneity. Whether the heterogeneity observed is itself the
359 driver of diaspora or merely a symptom is an important area for future study. Our
360 investigations of the potential drivers of diaspora were limited to genomic factors, and
361 further multi-platform studies looking at epigenetic and transcriptomic factors are other
362 important avenues of future research. We anticipate that analyses of single cells or small
363 clusters from primary sites, disseminated tumor cells and circulating tumor cells will also
364 yield finer resolution of the processes of dissemination and metastasis.
365 In cancer there are currently very few in-depth studies examining the spatial and temporal
366 evolution of metastases⁴⁸. Further studies are required to ascertain the extent to which our
367 diaspora theory pertains to other cancers.

368

369 **Acknowledgements**

370 Above all, we are indebted to the patients who donated tissue samples to this project and
371 thank them and their families who supported them through it. We would also like to thank
372 the following individuals for their help with study set-up, patient liaison and tissue
373 collection, Ben Smith, Nyrai Chinyama, Vijay Sujendran, Peter Safranek, Athanosios Xanthos,
374 Tara Nuckcheddy-Grant, Rachel de la Rue, Sebastian Zeki, Rachael Fels Elliott, Peter Collins,
375 Kitty Puttock, Sophie Rabey and staff at Arthur Rank Hospice and Luke A Wylie for scientific
376 discussion and contribution. We would like to thank the Oesophageal Cancer Clinical and
377 Molecular Stratification (OCCAMS) Consortium for providing the vehicle through which
378 funding for the International Cancer Genome Consortium (ICGC) was obtained. We are
379 grateful to Professor Simon Tavaré, FRS for his guidance and support for the esophageal
380 whole genome sequencing project as a part of the International Cancer Genome Consortium
381 (ICGC). We would like to thank Jo Westmoreland, LMB visual aids for her graphic art
382 expertise. Thanks also go to the Cancer Research UK Cambridge Institute Genomics Core for
383 their technical expertise. We thank the Human Research Tissue Bank, which is supported by
384 the National Institute for Health Research (NIHR) Cambridge Biomedical Research Centre,

385 from Addenbrooke's Hospital. Additional infrastructure support was provided from the
386 CRUK funded Experimental Cancer Medicine Centre in Cambridge. Computation by DCW
387 used the Oxford Biomedical Research Computing (BMRC) facility, a joint development
388 between the Wellcome Centre for Human Genetics and the Big Data Institute supported by
389 Health Data Research UK and the NIHR Oxford Biomedical Research Centre.

390 Ayesha Noorani was funded through an MRC Clinical Research Fellowship. The work was
391 funded through the above and an MRC core grant (RG84369) and an NIHR Research
392 Professorship (RG67258) to Rebecca Fitzgerald. Funding for sample sequencing (50x WGS)
393 was through the International Cancer Genome Consortium and was funded by a programme
394 grant from Cancer Research UK (RG66287). All OCCAMS samples which were part of the
395 surgical/endoscopy cohort were obtained from Cambridge patients. David Wedge is funded
396 by the Li Ka Shing foundation and the National Institute for Health Research (NIHR) Oxford
397 Biomedical Research Centre.

398

399 **Author Contributions**

400 AN designed and implemented the rapid autopsy study, collected the samples, performed
401 the experiments, analyzed data and wrote the manuscript. MG and S.D.P contributed
402 expertise in pathology and sample collection for the rapid autopsy study. ID-B and NG
403 assisted in study implementation, and along with JC, assisted with sample collection at
404 autopsy. M.S performed the structural variant analysis. M.D.E performed genomic data
405 generation and QC. LB conducted data management. XL, PL-S and JW were involved with
406 autopsy sample collection, advice on experiments and data analysis, and XL contributed to
407 experiments, paper writing, and figure design. LA and IM assisted with data analysis. NG
408 assisted with study Implementation. SMac coordinated the sequencing of samples from the
409 OCCAMS project and contributed to paper writing. SM and AM provided pathology
410 data. TT, SG, LP and DG assisted in implementation and ethical conduct of the autopsy
411 study. R.H.H and AH were involved in surgical sample collection and providing surgical
412 expertise. M.R.S contributed to critical evaluation of the study data and manuscript. D.C.W
413 was responsible for data analysis, paper writing, and assuring integrity of data. The OCCAMS
414 consortium was the vehicle through which the infrastructure and funding was obtained to
415 support the study and the consortium contributed to discussions on the ICGC data and the

416 clinical ramifications. R.C.F provided grant funding and was responsible for study design,
417 supervision of the project, writing the paper and assuring integrity of the data.

418

419 The authors declare no competing interests.

420 **References**

421

- 422 1. Sporn, M.B. The war on cancer. *Lancet* **347**, 1377-81 (1996).
- 423 2. Waterman, T.A. *et al.* The prognostic importance of immunohistochemically detected
424 node metastases in resected esophageal adenocarcinoma. *Ann Thorac Surg* **78**, 1161-
425 9; discussion 1161-9 (2004).
- 426 3. Matsuda, S., Takeuchi, H., Kawakubo, H. & Kitagawa, Y. Three-field lymph node
427 dissection in esophageal cancer surgery. *J Thorac Dis* **9**, S731-S740 (2017).
- 428 4. Lou, F. *et al.* Esophageal cancer recurrence patterns and implications for surveillance.
429 *J Thorac Oncol* **8**, 1558-62 (2013).
- 430 5. Smyth, E.C. *et al.* Oesophageal cancer. *Nat Rev Dis Primers* **3**, 17048 (2017).
- 431 6. Cunningham, D. *et al.* Capecitabine and oxaliplatin for advanced esophagogastric
432 cancer. *N Engl J Med* **358**, 36-46 (2008).
- 433 7. Greaves, M. & Maley, C.C. Clonal evolution in cancer. *Nature* **481**, 306-13 (2012).
- 434 8. Klein, C.A. Parallel progression of primary tumours and metastases. *Nat Rev Cancer*
435 **9**, 302-12 (2009).
- 436 9. Davis, A., Gao, R. & Navin, N. Tumor evolution: Linear, branching, neutral or
437 punctuated? *Biochim Biophys Acta* **1867**, 151-161 (2017).
- 438 10. Yates, L.R. & Campbell, P.J. Evolution of the cancer genome. *Nat Rev Genet* **13**, 795-
439 806 (2012).
- 440 11. Nik-Zainal, S. *et al.* The life history of 21 breast cancers. *Cell* **149**, 994-1007 (2012).
- 441 12. Murugaesu, N. *et al.* Tracking the genomic evolution of esophageal adenocarcinoma
442 through neoadjuvant chemotherapy. *Cancer Discov* **5**, 821-831 (2015).
- 443 13. Gerstung, M. *et al.* The evolutionary history of 2,658 cancers. *bioRxiv* (2017).
- 444 14. Secrier, M. *et al.* Mutational signatures in esophageal adenocarcinoma define
445 etiologically distinct subgroups with therapeutic relevance. *Nat Genet* **48**, 1131-41
446 (2016).
- 447 15. Dulak, A.M. *et al.* Gastrointestinal adenocarcinomas of the esophagus, stomach, and
448 colon exhibit distinct patterns of genome instability and oncogenesis. *Cancer Res* **72**,
449 4383-93 (2012).
- 450 16. Weaver, J.M. *et al.* Ordering of mutations in preinvasive disease stages of esophageal
451 carcinogenesis. *Nat Genet* **46**, 837-43 (2014).
- 452 17. Ross-Innes, C.S. *et al.* Whole-genome sequencing provides new insights into the
453 clonal architecture of Barrett's esophagus and esophageal adenocarcinoma. *Nat Genet*
454 **47**, 1038-46 (2015).
- 455 18. Dulak, A.M. *et al.* Exome and whole-genome sequencing of esophageal
456 adenocarcinoma identifies recurrent driver events and mutational complexity. *Nat*
457 *Genet* **45**, 478-86 (2013).
- 458 19. Nones, K. *et al.* Genomic catastrophes frequently arise in esophageal adenocarcinoma
459 and drive tumorigenesis. *Nat Commun* **5**, 5224 (2014).
- 460 20. Frankell, A.M. *et al.* The landscape of selection in 551 Esophageal Adenocarcinomas
461 defines genomic biomarkers for the clinic. *bioRxiv* (2018).
- 462 21. Yates, L.R. *et al.* Subclonal diversification of primary breast cancer revealed by
463 multiregion sequencing. *Nat Med* **21**, 751-9 (2015).
- 464 22. Rodriguez-Martin, B. *et al.* Pan-cancer analysis of whole genomes reveals driver
465 rearrangements promoted by LINE-1 retrotransposition in human tumours. *bioRxiv*,
466 179705 (2018).
- 467 23. Ajani, J.A. *et al.* Esophageal and esophagogastric junction cancers, version 1.2015. *J*
468 *Natl Compr Canc Netw* **13**, 194-227 (2015).

- 469 24. Yachida, S. *et al.* Distant metastasis occurs late during the genetic evolution of
470 pancreatic cancer. *Nature* **467**, 1114-7 (2010).
- 471 25. Sottoriva, A. *et al.* Intratumor heterogeneity in human glioblastoma reflects cancer
472 evolutionary dynamics. *Proc Natl Acad Sci U S A* **110**, 4009-14 (2013).
- 473 26. Mariette, C. *et al.* Pattern of recurrence following complete resection of esophageal
474 carcinoma and factors predictive of recurrent disease. *Cancer* **97**, 1616-23 (2003).
- 475 27. Alexandrov, L.B. *et al.* Signatures of mutational processes in human cancer. *Nature*
476 **500**, 415-21 (2013).
- 477 28. Liu, D. *et al.* Mutational patterns in chemotherapy resistant muscle-invasive bladder
478 cancer. *Nat Commun* **8**, 2193 (2017).
- 479 29. Behjati, S. *et al.* Mutational signatures of ionizing radiation in second malignancies.
480 *Nat Commun* **7**, 12605 (2016).
- 481 30. Dentre, S.C. *et al.* Portraits of genetic intra-tumour heterogeneity and subclonal
482 selection across cancer types. *bioRxiv* (2018).
- 483 31. Lodato, M.A. *et al.* Aging and neurodegeneration are associated with increased
484 mutations in single human neurons. *Science* **359**, 555-559 (2018).
- 485 32. Gao, Z., Wyman, M.J., Sella, G. & Przeworski, M. Interpreting the Dependence of
486 Mutation Rates on Age and Time. *PLoS Biol* **14**, e1002355 (2016).
- 487 33. Letouze, E. *et al.* Mutational signatures reveal the dynamic interplay of risk factors
488 and cellular processes during liver tumorigenesis. *Nat Commun* **8**, 1315 (2017).
- 489 34. Blokzijl, F. *et al.* Tissue-specific mutation accumulation in human adult stem cells
490 during life. *Nature* **538**, 260-264 (2016).
- 491 35. Alexandrov, L.B. *et al.* Clock-like mutational processes in human somatic cells. *Nat*
492 *Genet* **47**, 1402-7 (2015).
- 493 36. Pienta, K.J., Robertson, B.A., Coffey, D.S. & Taichman, R.S. The cancer diaspora:
494 Metastasis beyond the seed and soil hypothesis. *Clin Cancer Res* **19**, 5849-55 (2013).
- 495 37. Gudem, G. *et al.* The evolutionary history of lethal metastatic prostate cancer.
496 *Nature* **520**, 353-357 (2015).
- 497 38. Foulds, L. The experimental study of tumor progression: a review. *Cancer Res* **14**,
498 327-39 (1954).
- 499 39. Naxerova, K. *et al.* Origins of lymphatic and distant metastases in human colorectal
500 cancer. *Science* **357**, 55-60 (2017).
- 501 40. Sottoriva, A. *et al.* A Big Bang model of human colorectal tumor growth. *Nat Genet*
502 **47**, 209-16 (2015).
- 503 41. Sjoquist, K.M. *et al.* Survival after neoadjuvant chemotherapy or chemoradiotherapy
504 for resectable oesophageal carcinoma: an updated meta-analysis. *Lancet Oncol* **12**,
505 681-92 (2011).
- 506 42. Gabriel, E. *et al.* Novel Calculator to Estimate Overall Survival Benefit from
507 Neoadjuvant Chemoradiation in Patients with Esophageal Adenocarcinoma. *J Am*
508 *Coll Surg* **224**, 884-894 e1 (2017).
- 509 43. Burt, B.M. *et al.* Utility of Adjuvant Chemotherapy After Neoadjuvant
510 Chemoradiation and Esophagectomy for Esophageal Cancer. *Ann Surg* **266**, 297-304
511 (2017).
- 512 44. Pasquali, S. *et al.* Survival After Neoadjuvant and Adjuvant Treatments Compared to
513 Surgery Alone for Resectable Esophageal Carcinoma: A Network Meta-analysis. *Ann*
514 *Surg* **265**, 481-491 (2017).
- 515 45. Parikh, A.R. *et al.* Liquid versus tissue biopsy for detecting acquired resistance and
516 tumor heterogeneity in gastrointestinal cancers. *Nat Med* **25**, 1415-1421 (2019).
- 517 46. Van Roy, N. *et al.* Shallow Whole Genome Sequencing on Circulating Cell-Free
518 DNA Allows Reliable Noninvasive Copy-Number Profiling in Neuroblastoma
519 Patients. *Clin Cancer Res* **23**, 6305-6314 (2017).

- 520 47. Hu, Z. *et al.* Quantitative evidence for early metastatic seeding in colorectal cancer.
521 *Nat Genet* **51**, 1113-1122 (2019).
522 48. Robinson, D.R. *et al.* Integrative clinical genomics of metastatic cancer. *Nature* **548**,
523 297-303 (2017).
524

525 **Figure Legends**

526 **Figure 1 Overall project strategy and study design**

527 a. Overall Strategy to identify clonal evolution in metastatic EAC. There were three main
528 steps in this study which comprised: Clonal discovery at autopsy (see Supplementary Note
529 High Depth Whole Genome Sequencing (50x WGS), Mutation clustering and phylogenetic
530 tree construction, dN/dS analysis and Mutational Signature Analysis); Spatial tracking at
531 autopsy (see Supplementary Note Shallow Whole Genome Sequencing (1x WGS) and
532 Temporal tracking at earlier time-points (see Supplementary Note Shallow Whole Genome
533 Sequencing (1x) for Subclone identification, Supplementary Table 12 for precise samples for
534 plasma and Extended Data Fig. 9 for FFPE diagnostic samples). Colored circles depict clones
535 and subclones respectively. b. Sampling Strategy at Rapid Autopsy. Areas sampled for the
536 50x WGS part of the study are shown in blue and for 1x WGS are shown in orange. c. Study
537 Design and Sequencing Strategy. The flow chart demonstrates the study design and how this
538 relates to sequencing. Clonal Discovery is in blue and Clonal Tracking in orange. The sample
539 distribution for 50x WGS and 1x WGS are shown. 50x WGS = High depth WGS (50x), 1x WGS
540 = Shallow WGS (1x). n = number of cases, s = number of samples. †=248 solid tissue
541 samples, and 8 ctDNA at autopsy. CNA, copy number alteration; SNV, single nucleotide
542 variant; MRCA, most recent common ancestor.

543 **Figure 2 Phylogenetic Analysis of ten cases with nodal and distant metastases**

544 Patient body maps (S=surgical case, P=rapid autopsy) are shown. Green circles denote
545 lymph node metastases and yellow circles distant metastases. The labels within each circle
546 describe the specific location (see Supplementary Table 3, 4). An organ is shown in color if
547 metastases were sequenced from that site. The adjacent wedged semi-circle depicts the
548 clinical timelines for each patient. Each wedge corresponds to one month; blue wedges
549 indicate the total lifetime of the patient and red wedges periods of therapy. Phylogenetic
550 trees for each patient are shown and methodology is in Supplementary Note and Extended
551 Data Fig. 1a-b; pink = truncal events shared by all samples, purple = branch events shared by
552 more than one sample, yellow = leaves, events unique to a sample. The circle at the end of
553 a trunk, branch or leaf represents a clone or subclone. Each clone or subclone is annotated
554 to show which samples it is present in. E1-E4 = primary esophageal tumor, L1-L4= lymph
555 nodes, D1-8 =distant metastases, B = Barrett's Esophagus. A subclone annotated with E1, L2

556 for example indicates that this subclone is seen only in samples E1 and L2. The CCF of each
557 subclone/clone (barring the MRCA) is in Supplementary Table 5 and 6. The length of the
558 branches of the tree are reflective of the number of SNVs in the subclone/clone. The scales
559 adjacent to each case are relative, given the variable number of SNVs per case. Trees are
560 annotated with potential driver events, black: missense variants, red: amplifications. Gray
561 dots outlined with a black dashed line denote the first subclone/clone to metastasize that
562 would be classified as non-curative based on anatomical location. Red dots mark the
563 stellate pattern on the phylogenetic tree.

564 **Figure 3 Mutational Signatures**

565 a. Contributions of mutational signature in 18 cases (n=122) across the cohort. The bar chart
566 displays samples on a per case basis (X-axis) and depicts the number of SNVs contributing to
567 each signature (Y-axis). b. Mutational signatures pre-and post- diaspora across all samples
568 (n=122) in 18 cases.

569 Mutations were separately assigned to signatures and the proportion of mutations within
570 each case assigned to each signature is shown. Dark lines = median, Boxes = 25th and 75th
571 quartiles, whiskers extend to the most extreme point within 1.5× interquartile range of the
572 box edge. Signatures 1 mutations have a significantly lower representation in post-diaspora
573 mutations, while signature 3 mutations have significantly high. c. Mutational signature
574 analysis of ageing signature (signature 1) pre-and post-diaspora in all cases (n=8) with local
575 and distant spread ($p < 1.18 \times 10^{-90}$ across all cases) Chi squared test was used to determine
576 the p value. Survival is shown in months from the point of diagnosis *=cases which
577 underwent surgery.

578

579 **Figure 4 1x WGS and similarity matrix clustering of 248 further tissue samples from six** 580 **cases**

581 1x WGS was performed at an average depth of 1x to track subclones and clones previously
582 discovered using 50x WGS for further tissue samples (n=248). Pearson correlation
583 similarity matrix clustering was performed on all samples for each case (plotted against
584 each other) with red indicating sample similarity ($r=1$) and blue indicating dissimilarity ($r=-$
585 1). Sample sites used in this part of the study are shown in Supplementary Table 9 and the
586 entire organ is highlighted if solid organ sites were sequenced. For example, liver
587 metastases were only seen in P4, P6, P8, P10. Similarly, P2 had lymph nodes only (only

588 colored dots are seen which represent lymph nodes, no solid organs are highlighted).
589 Clustering was performed based on the presence of subclones and clones already
590 detected using 50x WGS and distinct clusters were identified for each case as
591 demonstrated by the adjacent key per case (each group is both colored and numbered).
592 Samples are displayed on the adjoining body maps for which the color coding corresponds
593 to the genomic clustering in the adjacent heatmap. Sites with multiple samples are
594 magnified and the division of samples shown. Maps of the primary tumor with
595 representation of metastatic subclones are shown with each case, with the colors of the
596 subclones being the same as those in the matrix and body map. Areas shaded red in the
597 primary tumor represent subclones that were not detected in the metastatic samples that
598 underwent 1x WGS and were instead confined to areas of the primary tumor.

599

600 **Figure 5 Temporal and spatial tracing of metastatic subclones in plasma**

601 a. Plasma ctDNA 1x WGS and digital droplet PCR (ddPCR) analysis for *TP53* mutant allele
602 fraction (MAF) for P10 and P6. The MAF of *TP53* (%) is shown on the Y-axis and days from
603 diagnosis are shown on the X-axis. The shaded areas represent time periods of therapy. 1x
604 WGS at select time-points was performed and the clonal composition of these samples
605 are shown by the presence of colored clusters. The color of each corresponds to the color
606 of the corresponding node on the adjacent 50x phylogenetic tree with the presence of
607 colored clusters which correlate with the 50x tree. Moreover, copy number traces for
608 each time point are shown for select chromosomes. b. The presence or absence of
609 amplifications and deletions in plasma compared to tissue, detected from 1x WGS for 8
610 cases. Tissue refers to all samples collected at autopsy and at earlier time-points. c.
611 Stacked bar charts to demonstrate the presence or absence of clusters across all plasma
612 samples, including truncal and branch clusters using 1x WGS.

613

614 **Figure 6 Diaspora model of metastatic spread and associated features**

615 Panel a depicts clonal diaspora with colored circles representing clones and subclones. *=
616 evidence of selection. Panel b explains the five features seen in diaspora (one is defining,
617 and the other are associated with diaspora) and whether these are present (✓) or absent
618 (x) in each case. *✓ implies that the feature is present, and that the evidence was from
619 1x WGS.

620 **Methods**

621 **Statistics**

622 Unless otherwise stated, statistical analyses were performed using R, version 3.3.3.
623 Clustering of mutations was carried out using a previously published Bayesian Dirichlet
624 Process method, DPCLust (<https://github.com/Wedge-Oxford/dpclust>), which calculates
625 CCFs of each SNV, taking into account tumor purity and copy number aberrations as
626 previously described⁴⁹. Analysis of structural variants used generalized linear models,
627 implemented with the R package MASS. Grouping of 1x WGS samples was performed with
628 the GENE-E package (<https://software.broadinstitute.org/GENE-E/download.html>).
629 Wilcoxon signed rank tests and Chi-squared tests were used as described in the main text.
630 Simulations were used to ascertain the robustness of DPCLust to violations of the infinite
631 sites assumption and its sensitivity to detect small deviations from stellate patterns.
632 Simulations were also used to confirm the correlation between the number of mutations
633 detected from 1x WGS and CCF determined from 50x WGS, as described in Online Methods.
634 dN/dS analysis was performed using the previously published package dndscv⁵⁰
635 (<https://github.com/im3sanger/dndscv>).

636

637 **Patient recruitment and Sample collection**

638 EAC patients were recruited from Addenbrooke's Hospital, Cambridge University Hospitals
639 NHS Trust with the explicit aim to study the clonal evolution of metastases as a sub-study
640 within OCCAMS (Oesophageal Clinical And Molecular Stratification). When it was clear that
641 extensive sampling of metastases could not be achieved without multiple invasive
642 procedures, the PHOENIX autopsy study was set up (Phylogenetic of Oesophageal
643 Neoplasia – An Investigation of Clonal Expansion under REC 07/H0305/52, and REC
644 EE/0043) with a prospective study design. Due diligence was undertaken to ensure
645 compliance with ethical regulations at all times. Patients were eligible if they were at least
646 18 years of age and had received a confirmed diagnosis of EAC following central pathology
647 review. Patients were only approached for the PHOENIX study following a palliative
648 diagnosis of metastatic EAC, with the full involvement of the multidisciplinary team.
649 Samples from the PHOENIX autopsy study were obtained within 6 hours of death and all
650 post-mortems were carried out at Papworth Hospital NHS Trust, United Kingdom.

651 Samples from Cambridge OCCAMS patients were obtained during diagnostic
652 oesophagogastroduodenoscopy (OGD), at endoscopic ultrasound (EUS) and/or from the
653 surgical resection specimen. Where possible, multiple samples were taken from spatially
654 distinct sites of the primary tumor or metastases. In two cases, brain metastases were
655 sampled at a clinically indicated craniotomy. Blood or normal squamous esophageal
656 samples, at least 5cm distant from the tumor, were used as a germline reference.

657 All tissue samples were snap-frozen in liquid nitrogen immediately after collection and
658 stored at -80°C. Cancer samples were deemed suitable for DNA extraction only after
659 consensus review of an H&E stained frozen section, from the same sample that would be
660 sent for sequencing, by two expert pathologists who confirmed tumor cellularity at $\geq 70\%$.

661 Samples with overall $\geq 70\%$ cellularity underwent dissection of the whole surface area with
662 a scalpel, whereas marked areas of $< 70\%$ underwent macrodissection or laser capture
663 micro- dissection aided by methylene blue staining visualized on the PALM-Zeiss
664 microscope (Zeiss, Oberkochen, Germany). An H&E stained slide was obtained before and
665 after extraction to confirm tumor cellularity of the microdissected section.

666 DNA was extracted from frozen tissues using the All PrepDNA/RNA Mini Kit (Qiagen,
667 Hilden, Germany) and from blood samples using the Nucleon™ Genomic Extraction kit
668 (Gen-Probe, San Diego, USA) according to the manufacturer's instructions. Some samples
669 were preserved in paraffin blocks after initially being stored in formalin. DNA from these
670 samples was extracted using the QiAmp FFPE Kit (Qiagen). Plasma extraction (for ctDNA)
671 was performed using the QiASymphony platform (Qiagen) as per the manufacturer's
672 instructions. All samples were eluted in 60 μ l of AE buffer and quantified using the High
673 Sensitivity Qubit (Thermo Fisher Scientific, MA, USA).

674 We included 388 samples, predominantly from PHOENIX, and some additional samples
675 from surgery and endoscopy (part of esophageal ICGC).

676 All samples were collected according to a strict SOP with quality control measures as already
677 described. All demographic and clinical data was anonymized and stored on a central study
678 database (OpenClinica and Labkey). The clinical characteristics of the patients are provided
679 in Supplementary Table 1 and 2. In terms of specifics of sample collection at autopsy, the
680 primary tumor was opened down the midline of the esophagus and the greater curve of the
681 stomach to expose the lumen. The tumor was divided in 12 areas with sampling as shown.

682 The size of tumors varied per case, but the division of sampling was always kept identical to
683 preserve reproducibility. In terms of the strategy for genomic sequencing (as per Figure 1),
684 up to 3 lymph nodes were chosen for 50x WGS in the areas shown (cervical, regional and
685 para-aortic) and up to 24 lymph nodes in each case (8 further lymph nodes per cervical,
686 regional and para- aortic areas (as per the Japanese Classification of nodal staging⁵¹) were
687 chosen for the 1x WGS part of the study. At least one metastasis per solid organ was chosen
688 for 50x WGS and for the 1x WGS part up to 8 samples were taken per organ for further
689 analysis. In addition, 8 samples from metastatic sites which had previously been sequenced
690 for 50x WGS were further sequenced for 1x WGS to assess the effects of metastatic
691 heterogeneity.

692

693 **Whole genome sequencing and data analysis strategy**

694 We used the Illumina HiSeq platform to perform WGS on multiple regions collected from
695 each primary tumor, lymph node and/or solid organ metastasis (Figure 1a,b, Supplementary
696 Table 3, 4). All DNA extractions and WGS conformed with ICGC quality control standards and
697 required $\geq 70\%$ cellularity and a matched germline sample. WGS was performed at high
698 depth (median coverage 66.3, IQR 56.1-87.2) to discover mutations in 122 samples from 18
699 patients (Supplementary Table 3, 4). In addition, low depth WGS (median coverage 1, IQR 1-
700 5) was performed to track these mutations spatially in up to 48 solid tissue samples per
701 case, (total=248) and 8 ctDNA samples at autopsy. Temporal tracking was performed in
702 cases with archival biopsy material, and where historical bloods were available
703 (Supplementary Table 12, Figure 5, Extended Data Fig. 6). For each patient the number of
704 subclones and the cancer cell fraction within each subclone was inferred using an extension
705 of a previously described Bayesian Dirichlet process¹¹ and we applied a set of previously
706 described rules to derive a phylogenetic tree (Additional Methods⁵²). All sequencing data
707 have been deposited in the European Genome-Phenome Archive under accession number
708 EGAD00001005434. *TP53* analysis in cell free tumor DNA (ctDNA) was performed using
709 Digital PCR on the Bio-rad platform (Bio-rad, California) using validated *TP53* assays
710 (Supplementary Table 14).

711

712 **Mutation clustering and phylogenetic tree construction**

713 The workflow used to perform mutation clustering and phylogenetic tree construction is
714 depicted in Extended Data Fig. 1a and illustrated with an example case, S3, in Extended
715 Data Fig. 1b. For each patient, we inferred the number of subclones and the fraction of
716 tumor cells within each subclone by using a previously described Bayesian Dirichlet process
717 (BDP) to cluster mutations according to their mutation copy number⁴⁹. We extended this
718 process into n dimensions for patients with n related samples, where the number of
719 mutant reads obtained from multiple related samples were modelled as independent
720 binomial distributions. The BDP uses Markov chain Monte Carlo (MCMC) to sample the CCF
721 values of the subclones in each sample. MCMC is run for 1000 iterations and outputs, for
722 each iteration, the sampled position of each cluster, p_{i_h} and the weight of each cluster, V_h ,
723 which is an estimate of the proportion of mutations assigned to that cluster. The first 200
724 iterations are considered as a ‘burn-in’ and are not used in subsequent steps. In order to
725 obtain the set of subclones present within a tumor and their CCF values, the following
726 procedure was followed:

- 727 • Using the aforementioned MCMC sampling of CCF values from all n samples, for
728 every possible triplet of samples, obtain posterior density estimates of CCF using
729 the function `kde` in the R package `ks`, with input parameters $x = p_{i_h}$, `bandwidth =`
730 `0.1`, `w = V_h`. Set `gridsize` such that density estimates are obtained to a resolution of
731 `0.02`. Identify local peaks in the posterior mutation density as locations higher
732 than any other gridpoint within a range of 2 gridpoints. For each local peak, define
733 a region representing a ‘basin of attraction’, defined by a set of planes running
734 through the `_point` of minimum density between each pair of cluster positions.
735 Assign each mutation to the cluster in whose basin of attraction they are most
736 likely to fall, using CCF values from MCMC sampling.
- 737 • Across the set of all possible triplets, identify sets of mutations that are assigned
738 to the same cluster in every triplet. Estimate the CCF of each cluster as the mean
739 CCF of the mutations assigned to that cluster. Estimate the 95% confidence
740 intervals as the `[0.025,0.975]` quantiles of the mean p_{i_h} values of the mutations
741 assigned to each cluster within MCMC sampling.

742 Finally, again using the aforementioned MCMC sampling of CCF values from all n samples,
743 for every pair of samples, plot the mutation density, estimated using the function `kde` in

744 the R package *ks*, with input parameters $x = \pi_{i_h}$, bandwidth = 0.1, $w = V_{i_h}$.

745 Taking a conservative approach, clusters were identified as subclonal only if the 95%
746 confidence intervals of the posterior estimate of the proportion of cells excluded the value
747 1. Clusters containing less than 1% of all mutations identified in a tumor were not included
748 in phylogenetic reconstruction.

749 Occasionally, copy number states are incorrectly called in small regions of some cancer
750 genomes. As a consequence, mutations falling in these regions have inaccurate estimates
751 of CCF and can cause artefact clusters. Such clusters may be identified after mutation
752 clustering since they contain a small percentage of mutations (less than 2.5%), the
753 mutations within them are located in localized regions of the genome, and, often, they
754 cannot be placed on the phylogenetic tree because they have discordant CCF values. We
755 excluded these clusters from phylogenetic tree construction. The number of clusters
756 excluded in total was seven (5 in P2, 1 in P3, 1 in P10). Two samples had low tumor content
757 (36% in P3_E1, 14% in S5_T1). As a result, CCF estimates for subclones found in these
758 samples are imprecise and led to violations of the sum rule (see below). The CCF values of
759 the relevant clusters were manually corrected to enable them to be placed on the
760 phylogenetic tree, as follows: P3_E1 only cluster adjusted from 1 to 0.85; S5_E1 truncal
761 cluster adjusted from 0.85 to 1.

762 To determine the most likely phylogenetic tree, we applied two rules, previously
763 described⁵². Briefly, the ‘sum rule’ (which is an extension of the pigeonhole principle
764 described in Ref 11), asserts that if a subclone A is ancestral to both subclones B and C and
765 if the summed CCFs of B and C exceed the CCF of A in any sample, the relationship
766 between the subclones must be linear. The ‘crossing rule’ is applied to tree construction
767 from multiple samples. It asserts that if the CCF of B is higher than the CCF of C in sample X
768 and the CCF of B is lower than the CCF of sample C in sample Y then B and C must be in
769 separate branches of the phylogenetic tree, i.e. they are not collinear. For all clonally
770 related samples, the same underlying phylogenetic tree must exist. This exerts much
771 greater stringency to the inferred ordering of subclonal clusters present in more than one
772 sample and defines their position on the phylogenetic tree unequivocally. Note that P9
773 contains two independent cancers derived from Barrett’s esophagus and adenocarcinoma
774 regions. CCF values are reported relative to the dominant cancer, so in P9_D4, which
775 contains both cancers, the two cancers are reported with CCFs of 100% and 69%. This

776 apparent violation of the sum rule results from the mathematical convenience of
777 normalizing to the dominant cancer.

778 It should be noted that the sum rule and crossing rule only strictly apply when the infinite
779 sites assumption (ISA) is obeyed. The ISA states that each mutation only occurs once during
780 the lifetime of a tumor and that mutations never revert to normal. A recent study⁵³ has
781 shown, through analysis of targeted sequencing of single cells, that the ISA is not always
782 followed in real data, for two reasons:

- 783 • Copy number alterations (CNAs), specifically losses and loss of heterozygosity,
784 have the effect of removing mutations in the deleted region, resulting in the
785 apparent 'reversion' of a mutation.
- 786 • The same mutation may occur on more than one occasion, particularly if the
787 mutation is a driver mutation.

788 In our study, we take account of CNAs when calculating the CCF of each mutation. In
789 regions that have undergone gain of one or both alleles, a mutation may be present on
790 more than one chromosome copy, up to the number of copies of the most amplified
791 chromosome copy. Conversely, if one or both chromosome copies have undergone loss in
792 a particular sample, a mutation may be lost in that sample. In the situation where a
793 mutation is unobserved in a sample and that sample has a copy number state lower than
794 that observed in another sample in which the mutation is observed, we do not call the
795 mutation as absent. Rather, we cluster it based on its CCF in the remaining samples,
796 treating its CCF in the target sample as unknown.

797

798 **Identification of cancer cell fraction**

799 For each mutation we calculated the mutation copy number as previously described, using
800 the mutant allele burden, tumor cellularity and locus specific copy number in the tumor
801 and matched normal⁴⁹. The mutation copy number reflects the percentage of tumor cells
802 within a sample carrying that mutation, and permits the cross-comparison of the mutation
803 in related samples despite differences in tumor purity and/ or copy number profiles.
804 Mutations present on multiple copies of a chromosomal segment will have a mutation
805 copy number greater than 1. To group mutations according to the percentage of cells
806 containing it, or cancer cell fraction (CCF), the number of chromosomes carrying the
807 mutation must be determined. For all mutations within amplified regions with a major

808 allele copy number, the observed fraction of mutated reads was compared to the expected
809 fraction of mutated reads resulting from a mutation present assuming a binomial
810 distribution³⁷.

811

812 **Annotation of the trees with mutations**

813 We annotated each tree with oncogenic or putative oncogenic alterations including
814 substitutions and copy number changes. For substitutions, cluster assignment information
815 from a multidimensional Dirichlet process was used.

816 For rearrangements and copy number changes, branch assignment was achieved by
817 considering the set of samples containing the variant and the subclonal fraction of the
818 associated copy number segment where applicable. All potential driver alterations were
819 annotated. For substitutions, structural variants and copy number events, these included a
820 set of genes compiled from the TARGET database from the Broad Institute and multiple
821 sequencing datasets for OAC^{14-16,18,19}.

822

823 **Shallow Whole Genome Sequencing for Subclone Identification**

824 For shallow whole genome sequencing, samples were sequenced to a median depth of
825 $\sim 1x$. It was not therefore feasible to call mutations de novo for these samples, but we were
826 able to count the number of mutations from each subclone that reported a mutant read in
827 $1x$ WGS sequencing. We performed simulations of $1x$ WGS data in order to ascertain the
828 correlation between the number of mutations identified and the CCF of each subclone.
829 First, we simulated subclones with CCF values between 0.01 and 1.00, assuming 1000
830 mutations per subclone, sequencing depth drawn from a Poisson distribution with
831 expected value 1, and binomial sampling of WT and mutant reads. The correlation
832 between the number of mutations detected and the CCF of the subclone was very high
833 (Pearson $r = 0.992$, Extended Data Fig. 4). In order to test whether subclones containing
834 fewer mutations also had good correlations between CCF and number of detected
835 mutations, we performed further simulations of subclones containing between 50 and
836 1,000 mutations and ascertained that the correlation remained very high (> 0.997) for
837 cluster sizes as small as 200 (Extended Data Fig. 5). Of the 169 subclones identified in our
838 study, only two contained fewer than 200 mutations, indicating that the number of

839 mutations detected is a good proxy for the CCF of a subclone.
840 SNVs from libraries sequenced to a minimum of 1x following filtering, were allocated to
841 subclones previously identified at 50x WGS. Mapping quality and base quality of 10 were
842 used. This resulted in tabulated counts for SNVs being allocated to subclones identified at
843 50x WGS for each sample. Normalization was performed according to the number of SNVs
844 assigned to each subclone from 50x WGS, and to the total number of SNVs in that sample
845 in order to account for potential differences in coverage, using the following equation:

$$846 \text{CCF}_{\text{cluster}} = n_{\text{cluster}}/n_{\text{truncal}} \times H_{\text{truncal}}/H_{\text{cluster}}$$

847 in which n_{cluster} and n_{truncal} are the numbers of loci in the target cluster and the truncal
848 cluster that have mutant reads in the target sample and H_{cluster} and H_{truncal} are the number
849 of mutations identified from 50x WGS in the target and truncal clusters. For each 1x WGS
850 sample, this provides an estimate of the CCF of each subclone within that sample.

851 In all cases, near equal coverage was obtained and in cases of low cellularity further
852 sequencing was performed in order to achieve this. After normalization, the GENE-E
853 package (<https://software.broadinstitute.org/GENE-E/download.html>) was used to cluster
854 the 1x WGS samples according to the similarity of their CCF profiles using Pearson
855 correlation.

856

857 **Data Availability**

858 Sequencing data that support the findings of this paper have been deposited in the
859 European Genome-phenome Archive with the accession code EGAD00001005434.

860

861 **Code Availability**

862 All code required to reproduce the analysis outline in this manuscript can be found in the
863 main and supplementary methods. There are no restrictions to the accessibility of this code.

864

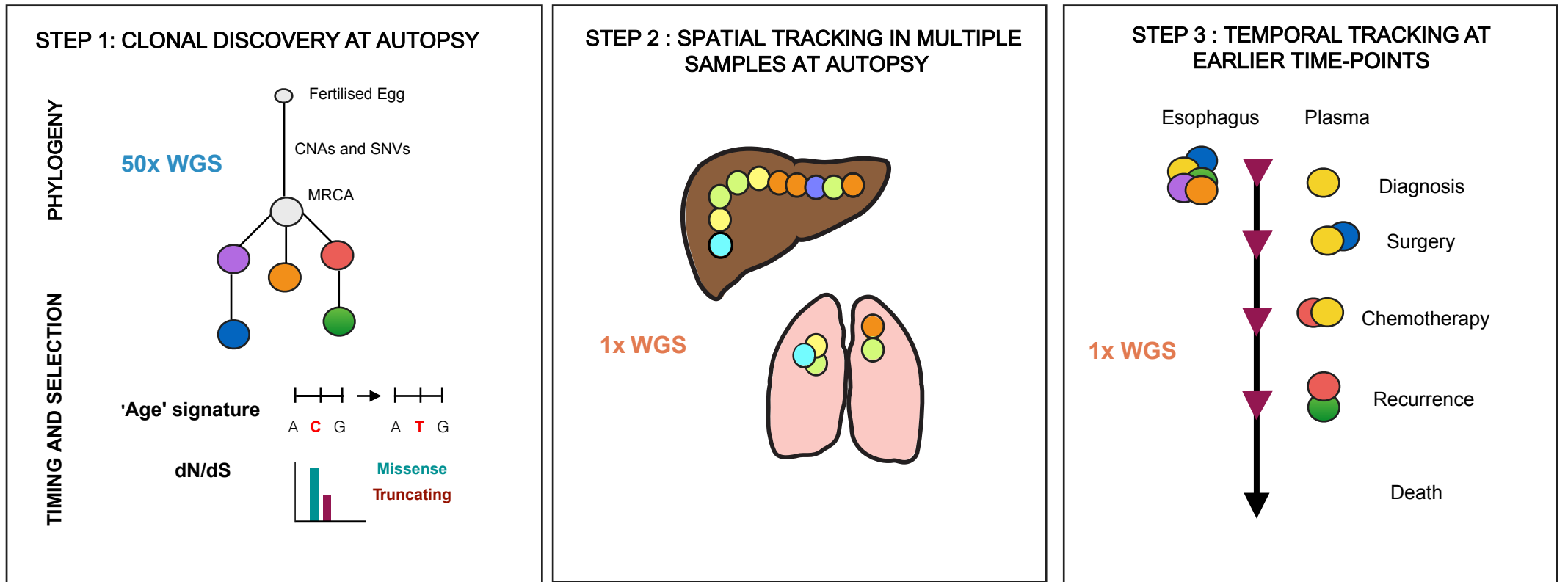
865 **Method-only references**

- 866 49. Bolli, N. *et al.* Heterogeneity of genomic evolution and mutational profiles in multiple
867 myeloma. *Nat Commun* **5**, 2997 (2014).
868 50. Martincorena Inigo, R.K.M., Gerstung Moritz, Dawson Kevin J, Haase Kerstin, Van
869 Loo Peter, Davies Helen, Michael R. Stratton Michael R, Campbell Peter J. Universal
870 Patterns Of Selection In Cancer And Somatic Tissues. *Cell* (2017).
871 51. Japanese Gastric Cancer, A. Japanese classification of gastric carcinoma: 3rd English
872 edition. *Gastric Cancer* **14**, 101-12 (2011).

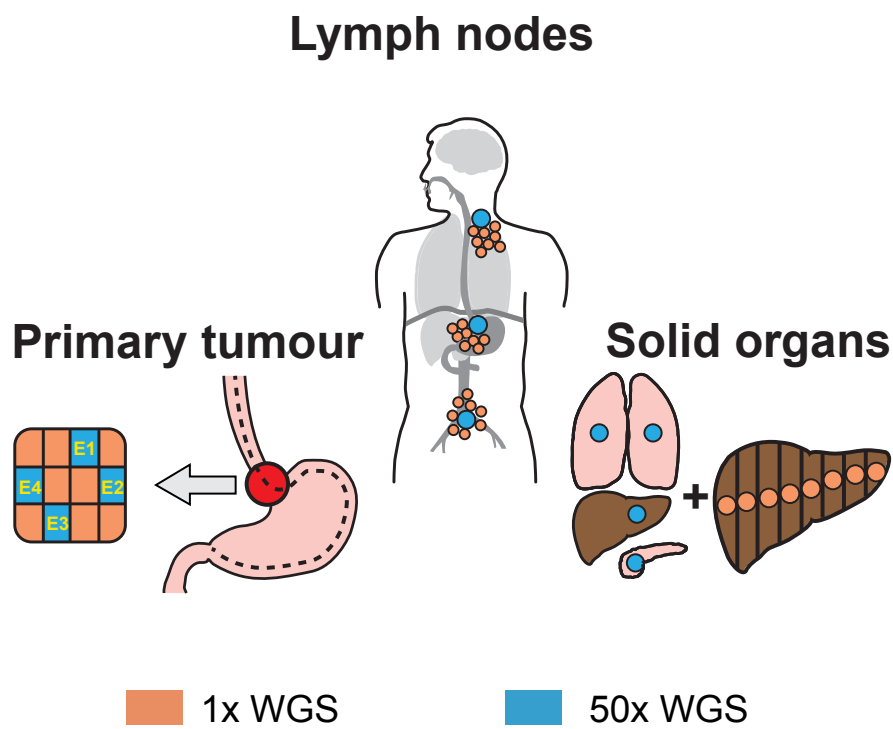
- 873 52. Jiao, W., Vembu, S., Deshwar, A.G., Stein, L. & Morris, Q. Inferring clonal evolution
874 of tumors from single nucleotide somatic mutations. *BMC Bioinformatics* **15**, 35
875 (2014).
- 876 53. Kuipers, J., Jahn, K., Raphael, B.J. & Beerenwinkel, N. Single-cell sequencing data
877 reveal widespread recurrence and loss of mutational hits in the life histories of tumors.
878 *Genome Res* **27**, 1885-1894 (2017).
879

Figure 1

a OVERALL STRATEGY TO IDENTIFY CLONAL EVOLUTION IN METASTATIC EAC



b SAMPLING STRATEGY



c SEQUENCING STRATEGY

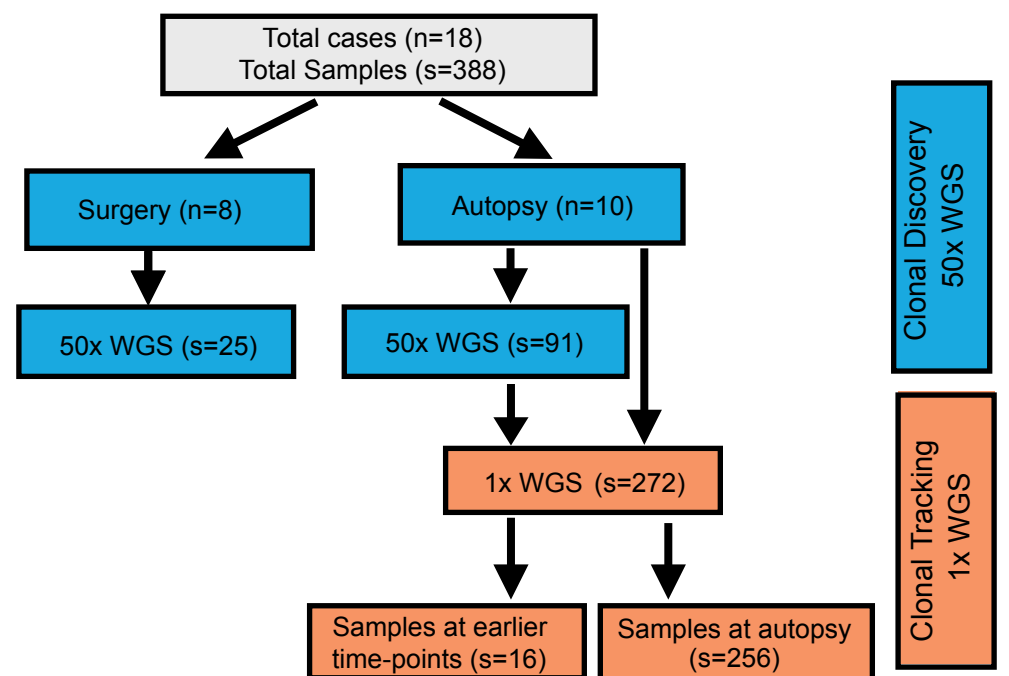


Figure 2

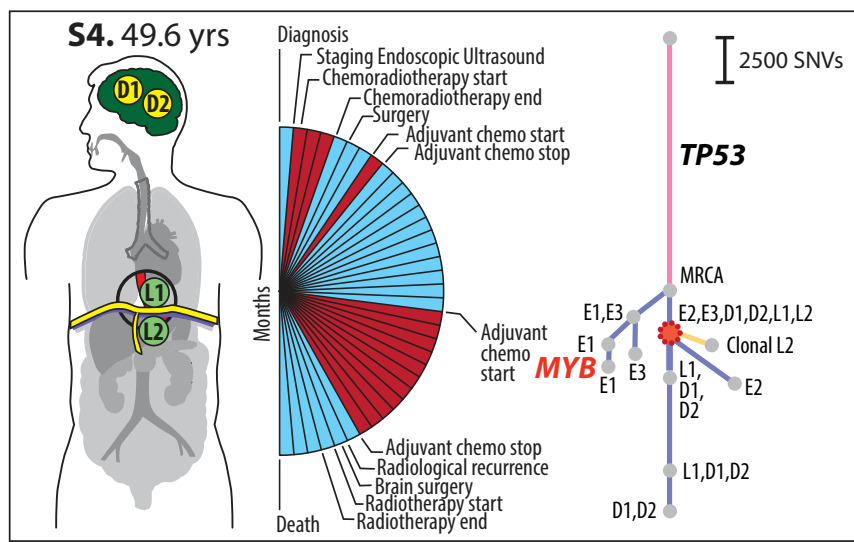
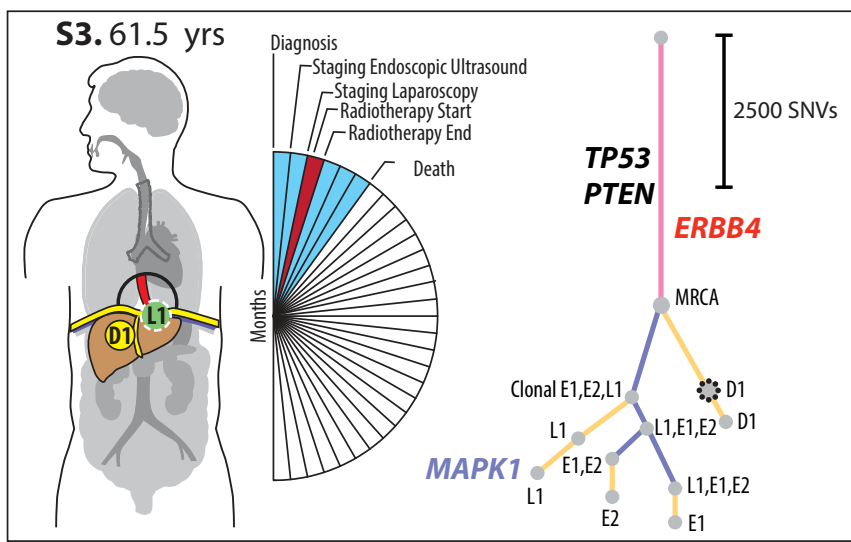
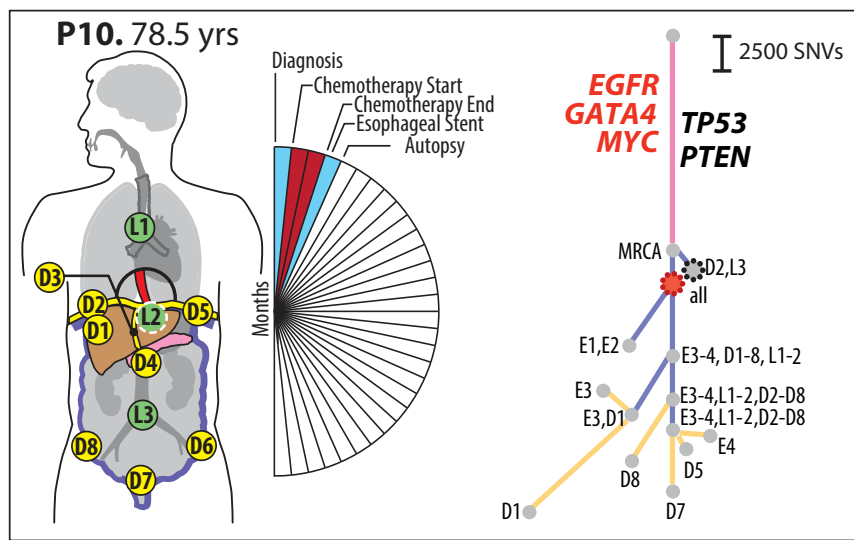
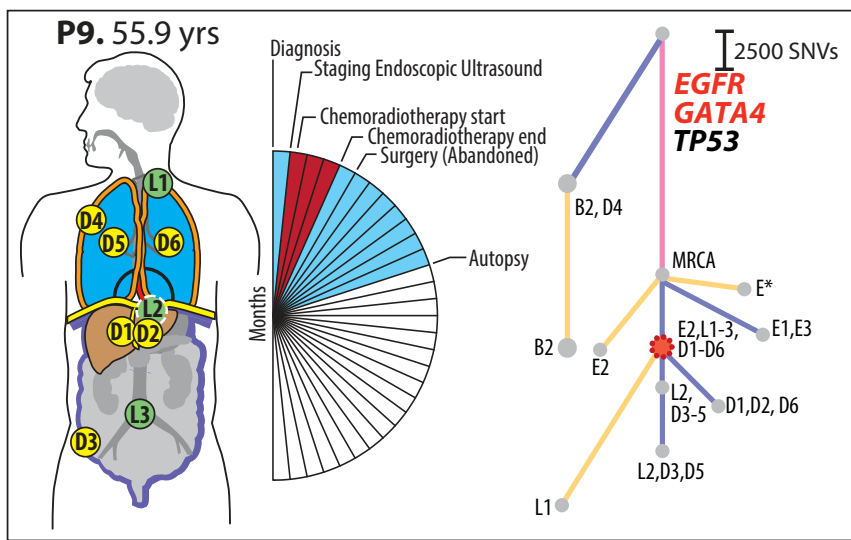
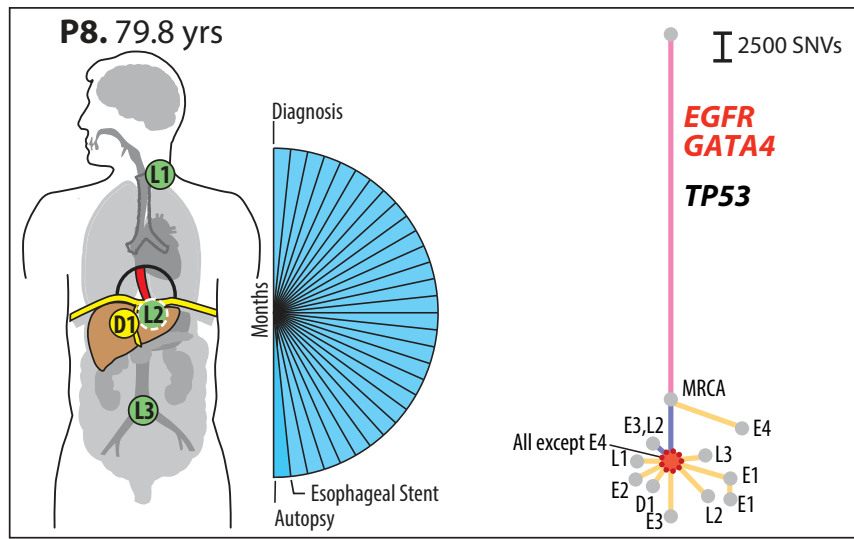
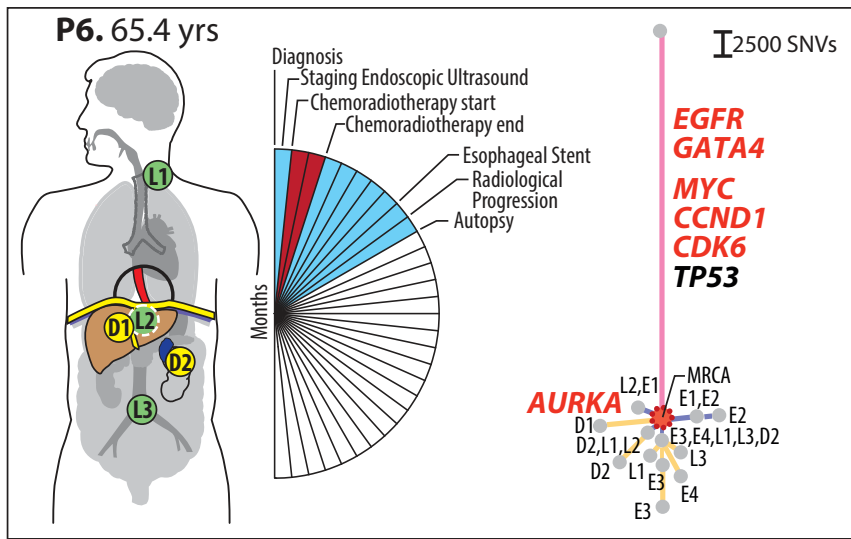
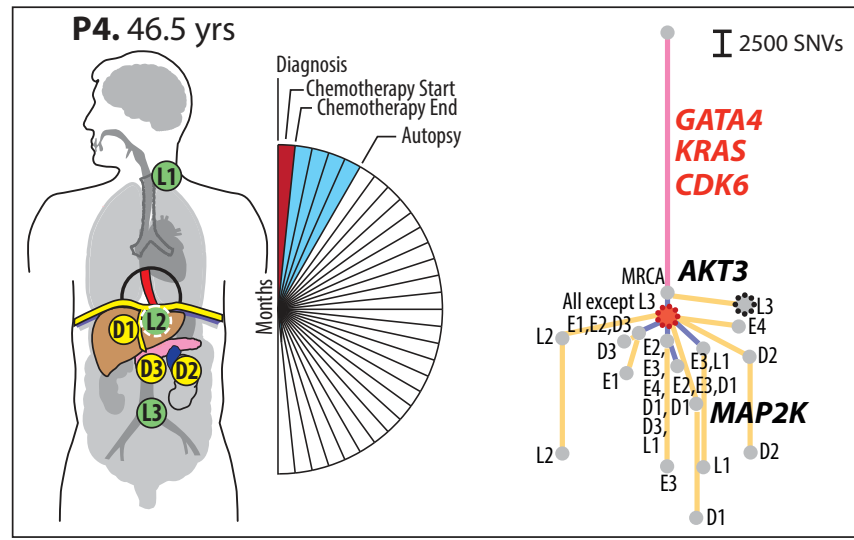
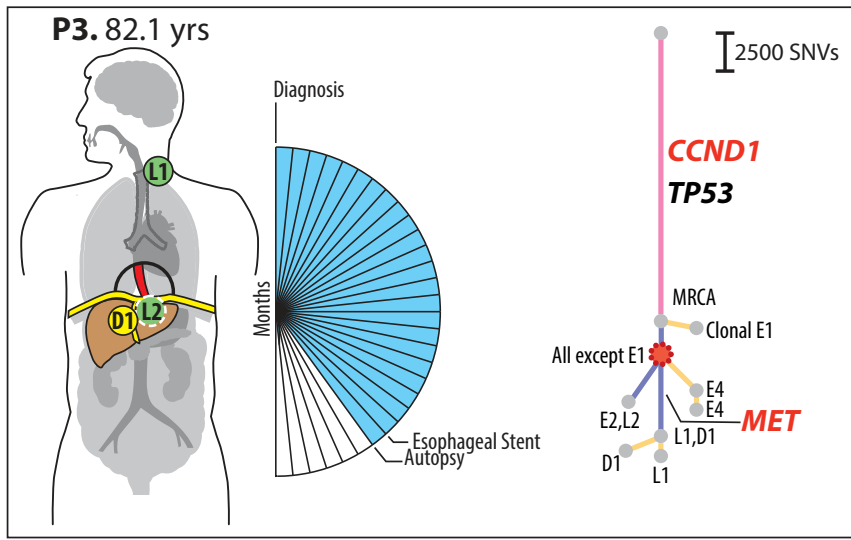
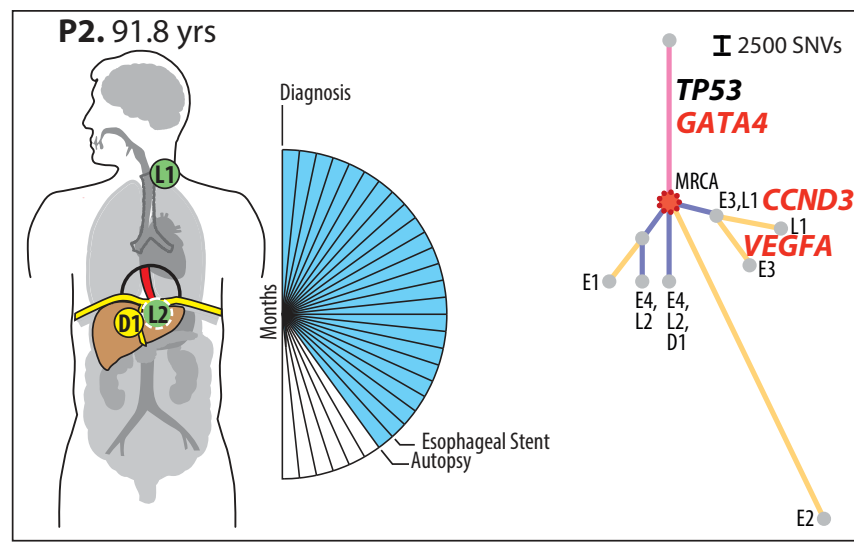
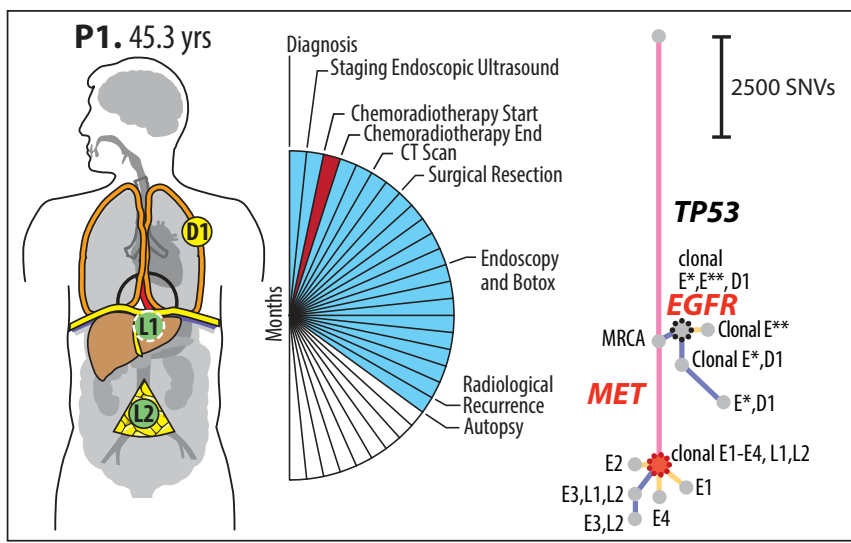
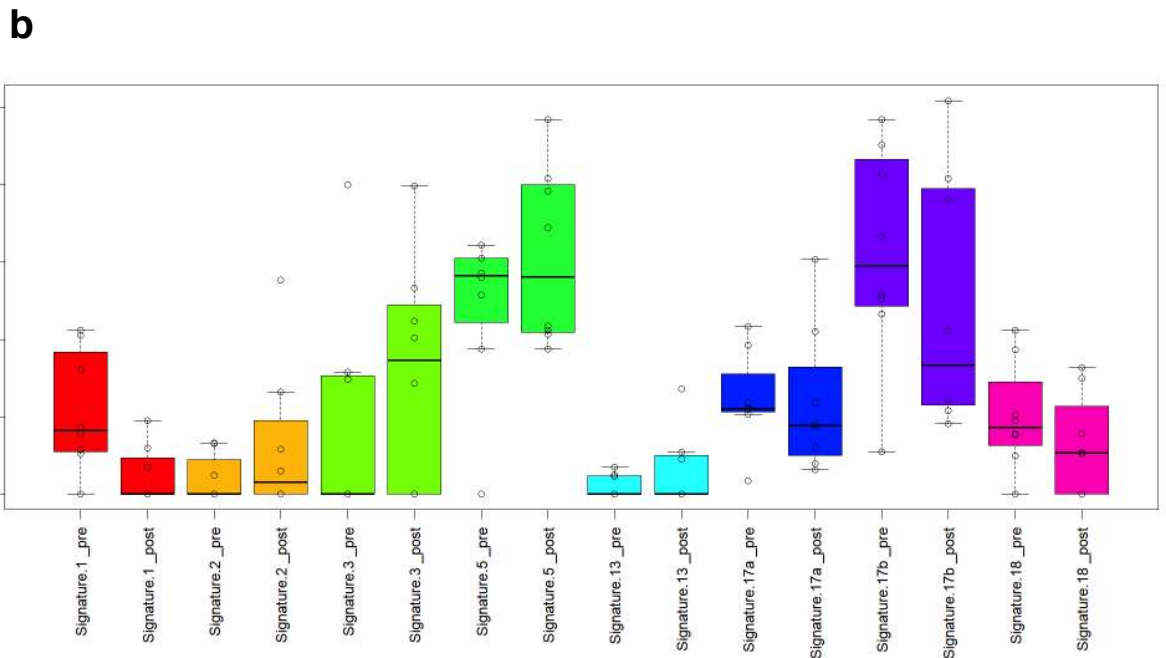
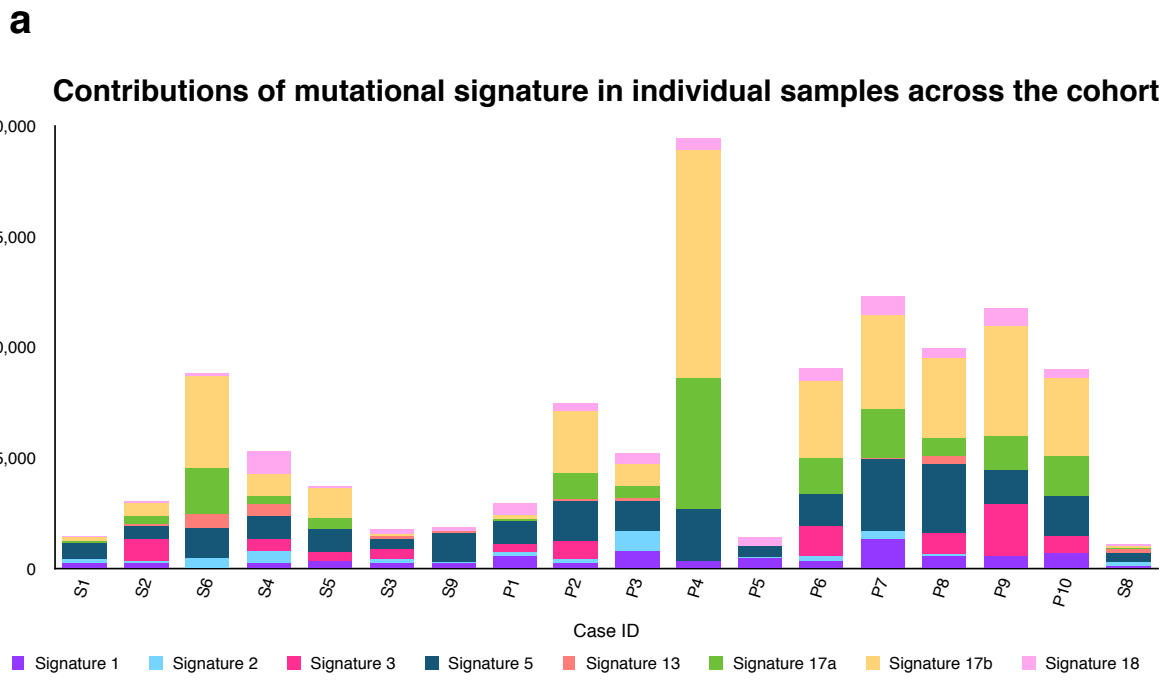


Figure 3



c

Case	Pre-Diaspora	Post-Diaspora	p value	Survival (months)
P1			9.10E-08	20*
P2			9.09E-05	12
P3			3.19E-05	14
P4			1.56E-95	5
P6			3.10E-14	5
P8			5.87E-38	30
P9			5.82E-72	12
P10			6.58E-161	4
S4			1.11E-14	37*

● Percentage of Clocklike
● Signature Percentage Non-Clocklike Signature

Figure 4

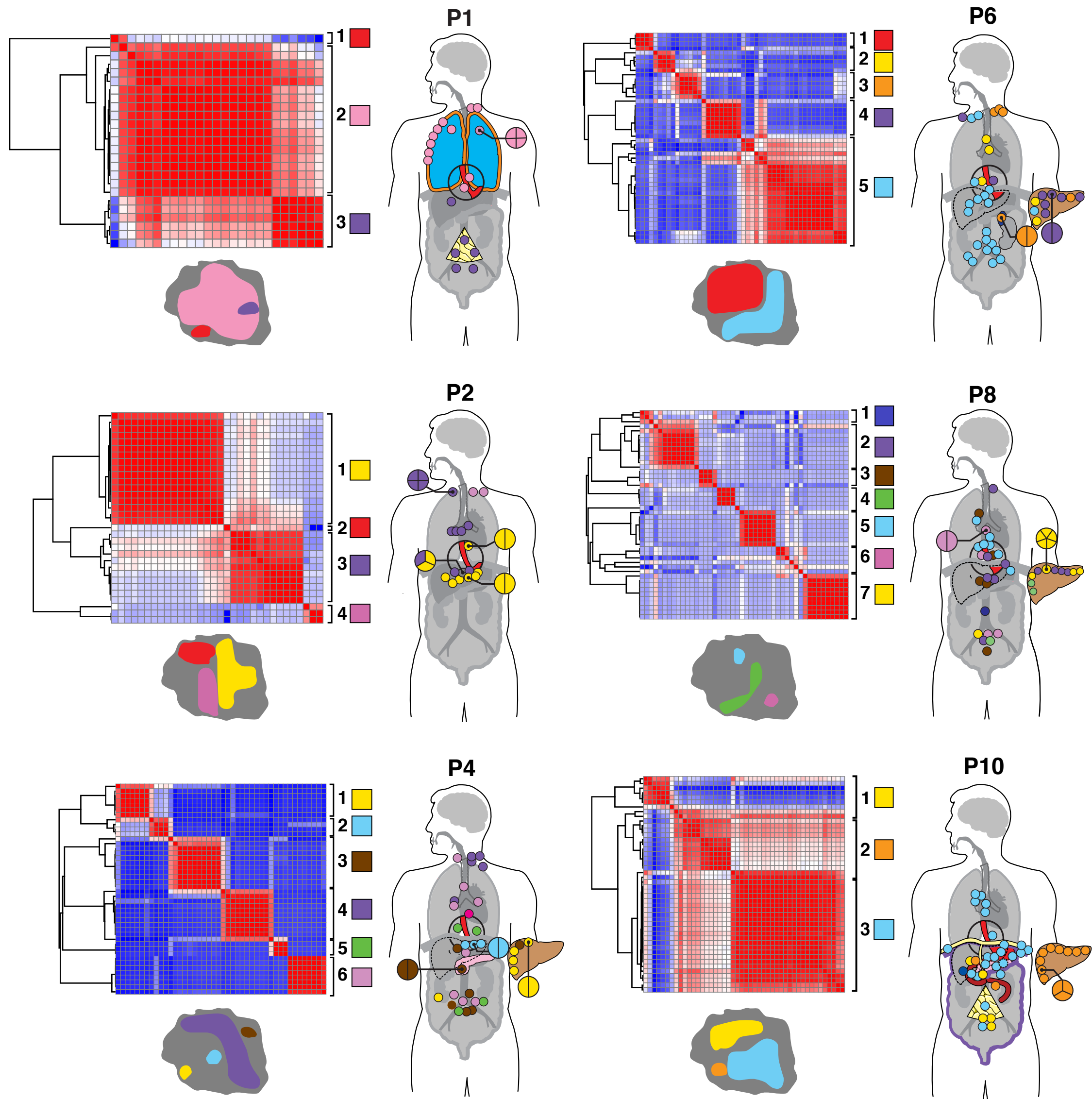
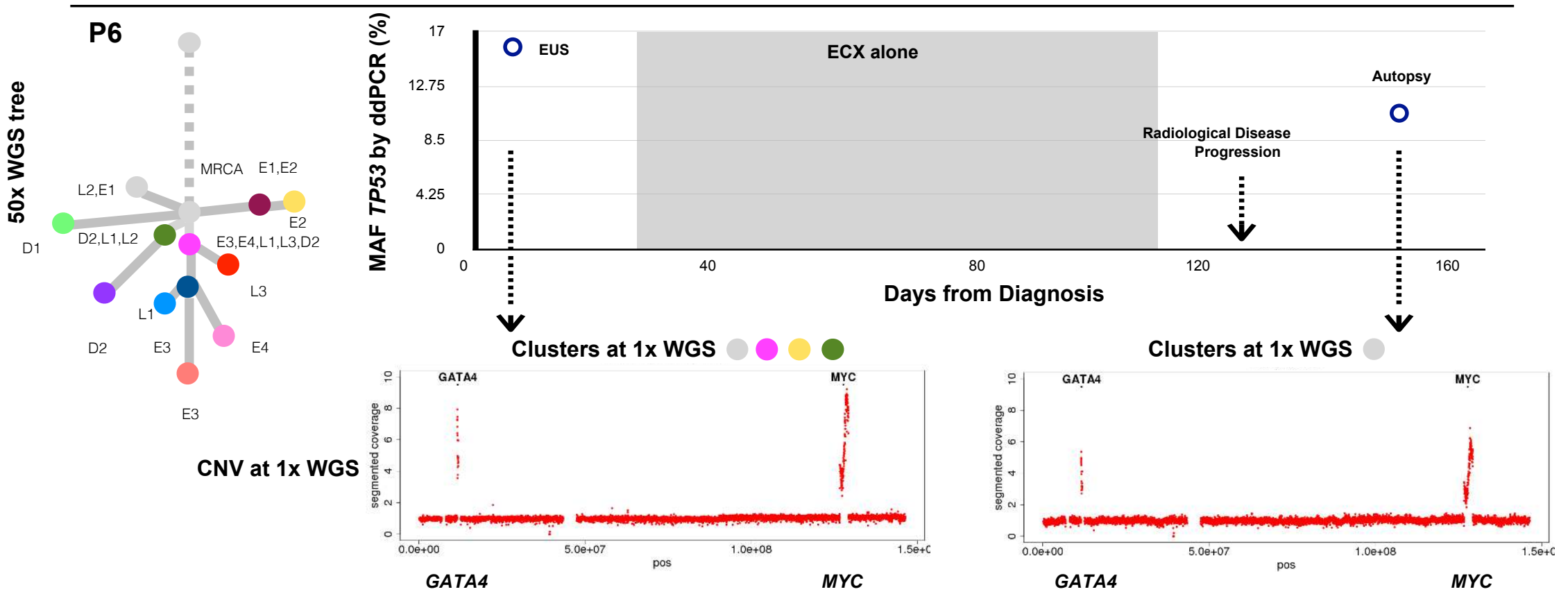
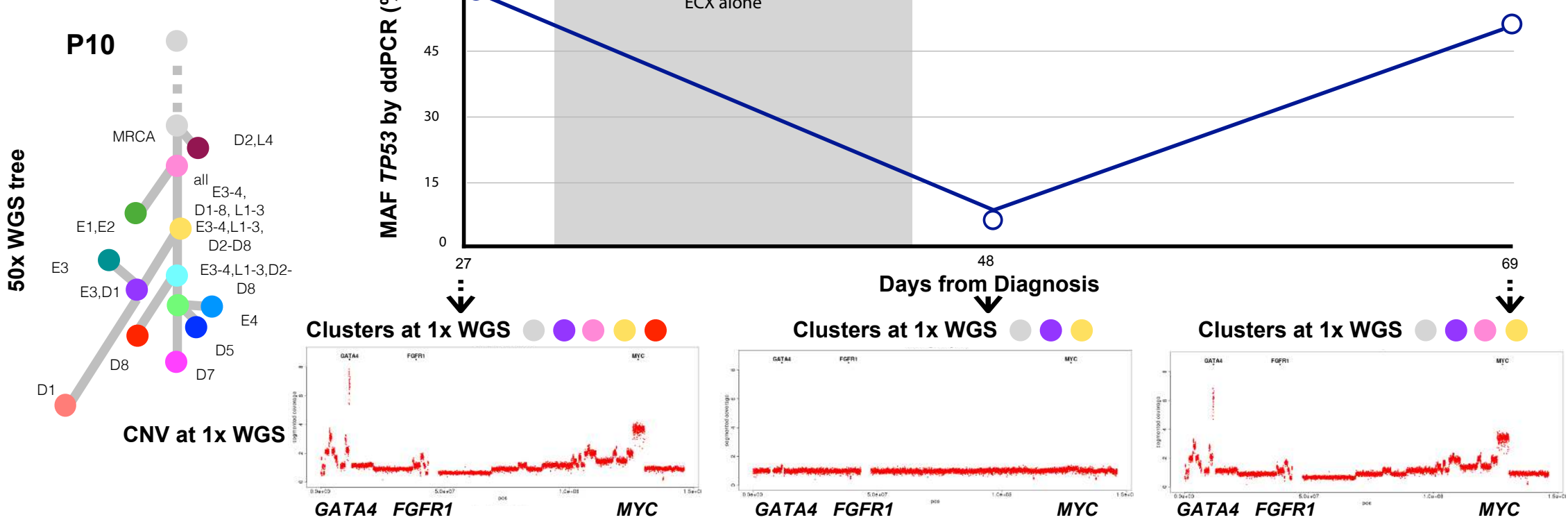


Figure 5 **a**

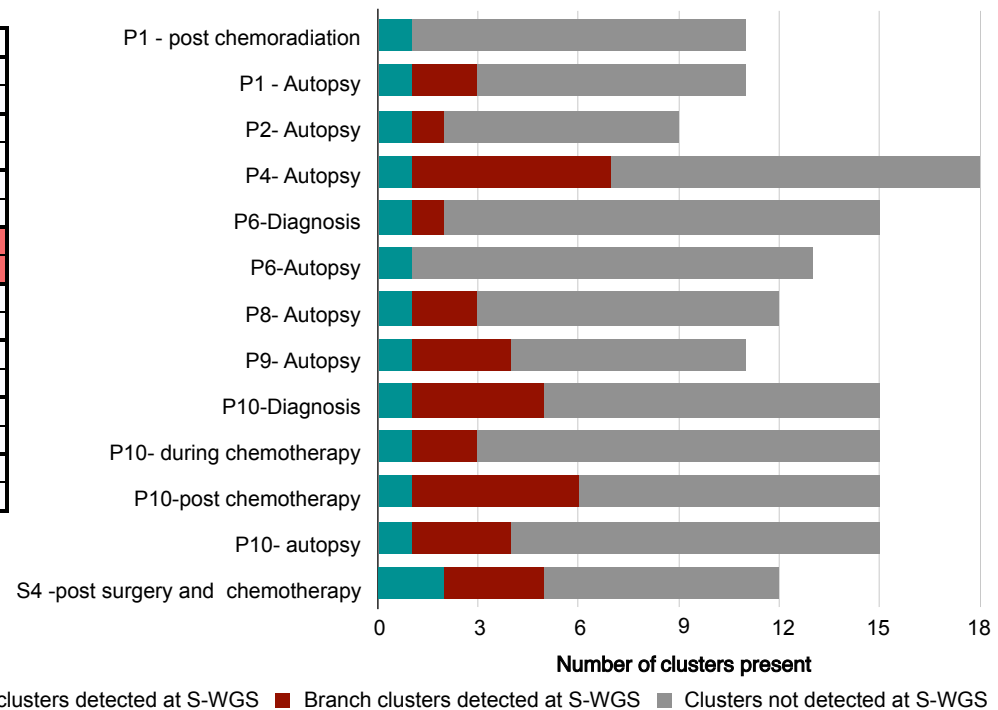


b

Case ID		<i>CDK6</i>	<i>CDKN2a</i>	<i>CCND1</i>	<i>CCND3</i>	<i>EGFR</i>	<i>GATA4</i>	<i>KRAS</i>	<i>MYC</i>	<i>MET</i>	<i>PRKC1</i>
P1	Plasma										
	Tissue										
P2	Plasma										
	Tissue										
P3	Plasma										
	Tissue										
P4	Plasma										
	Tissue										
P6	Plasma										
	Tissue										
P8	Plasma										
	Tissue										
P9	Plasma										
	Tissue										
P10	Plasma										
	Tissue										

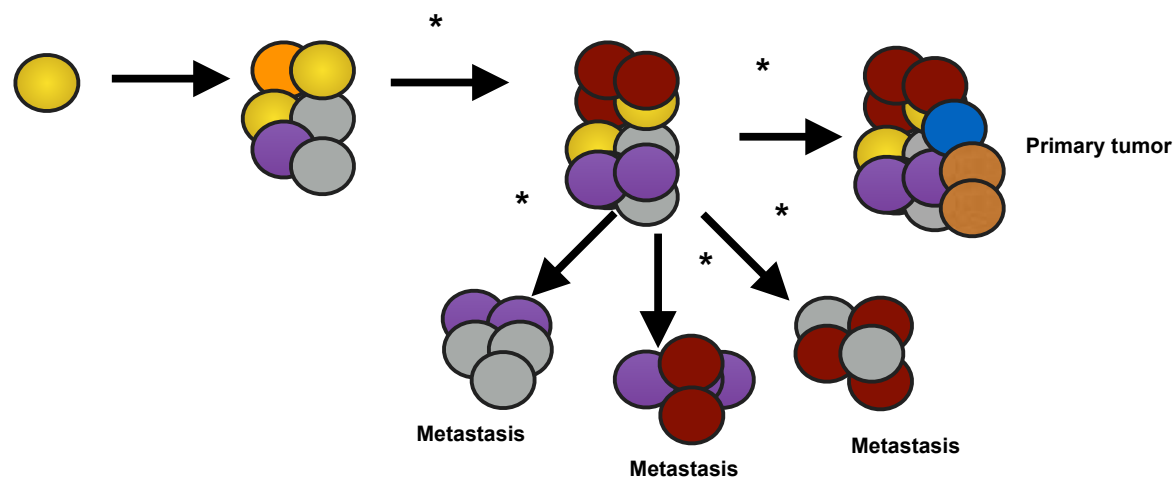
Legend: Loss (blue), Unaltered (white), Amplification (red)

c



a

Diaspora Model of Metastatic Spread



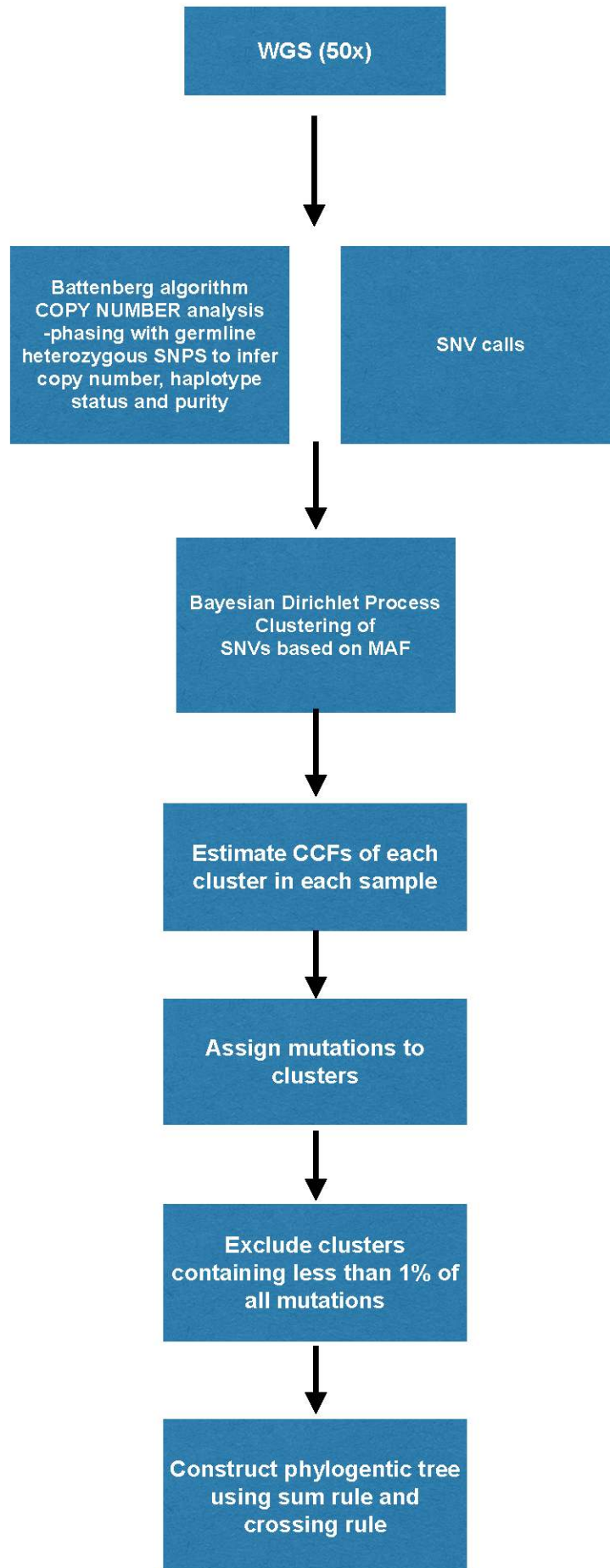
b

Features of Diaspora

Case	DEFINING	ASSOCIATED			
	Multiple subclones from primary spread to multiple metastatic sites	Stellate pattern of three or more subclones derived from the same ancestor found in metastatic sites	Lack of Signature 1 mutations, indicating rapid accumulation of mutations and near-synchronous spread	Spread of at least one subclone to organs of different types, including both lymph nodes and distant organs	Evidence for selection of subclones within the diaspora, indicative of an evolutionary niche (driver amplifications)
P1	✓	✗	✓	*✓	✗
P2	✓	✓	✓	✓	✓
P3	✓	✗	✗	✓	✓
P4	✓	✓	✓	✓	✓
P6	✓	✓	✓	✓	✓
P8	✓	✓	✗	*✓	✗
P9	✓	✓	✗	✓	✗
P10	✓	✗	✓	✓	✗
S3	✗	✗	✗	✗	✓
S4	✓	✓	✓	✓	✓

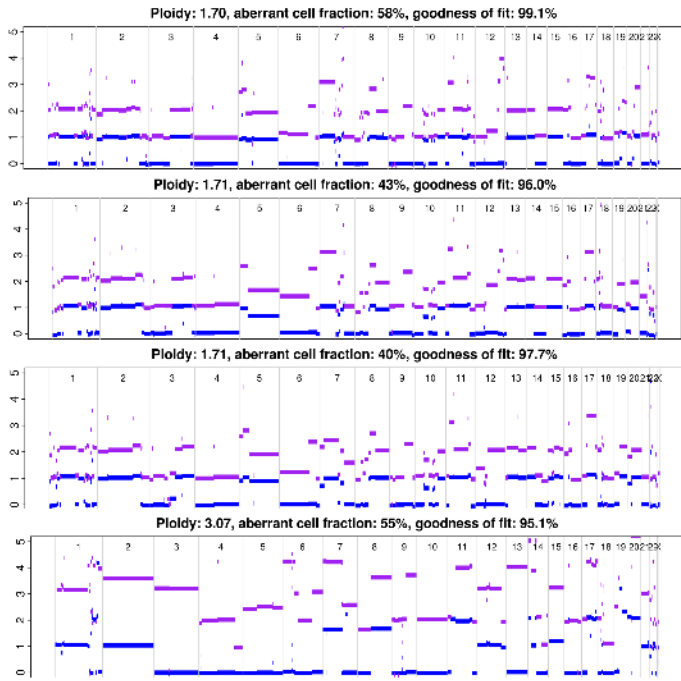
Extended Data 1

Overall Methodology

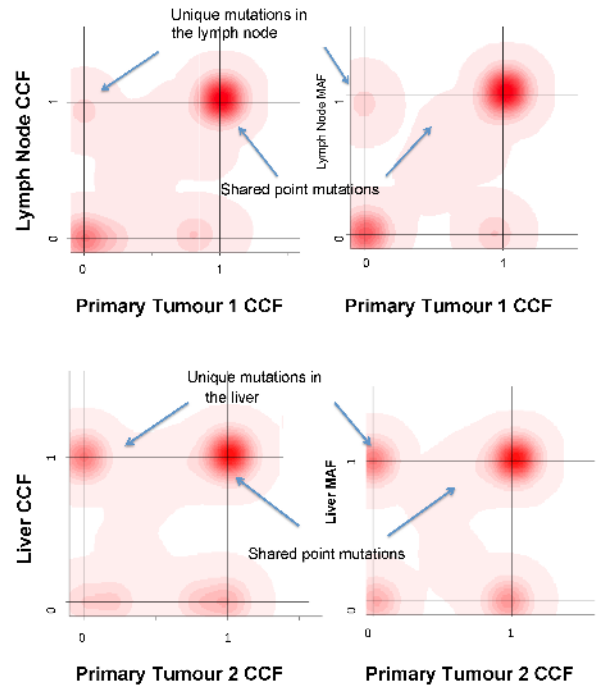


ExtendedData 2

1 Battenberg Algorithm



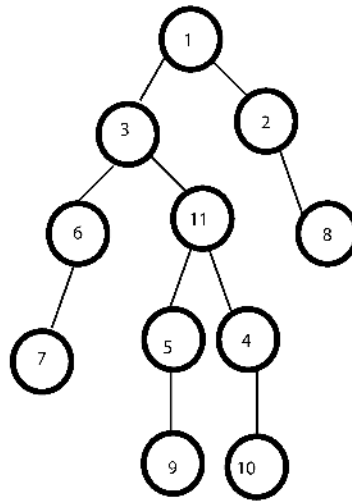
2 SNV clustering using Dirichlet Process



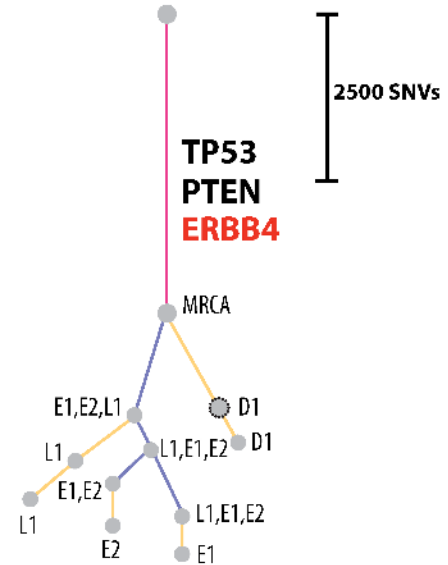
3 Clustering Results

Cluster No.	SNVs in cluster	Samples cluster seen in (CCF%)
1	5697	All (100%)
2	1913	D1 (96%)
3	1663	E1, E2, L1
4	1382	E1 (81%), E2 (4%), L1 (3%)
5	1332	E1 (21%), E2 (94%)
6	1139	L1 (94%)
7	1096	L1 (39%)
8	322	D1 (40%)
9	302	E2 (36%)
10	300	E1 (28%)
11	237	E1, E2 (95%), L1 (3%)

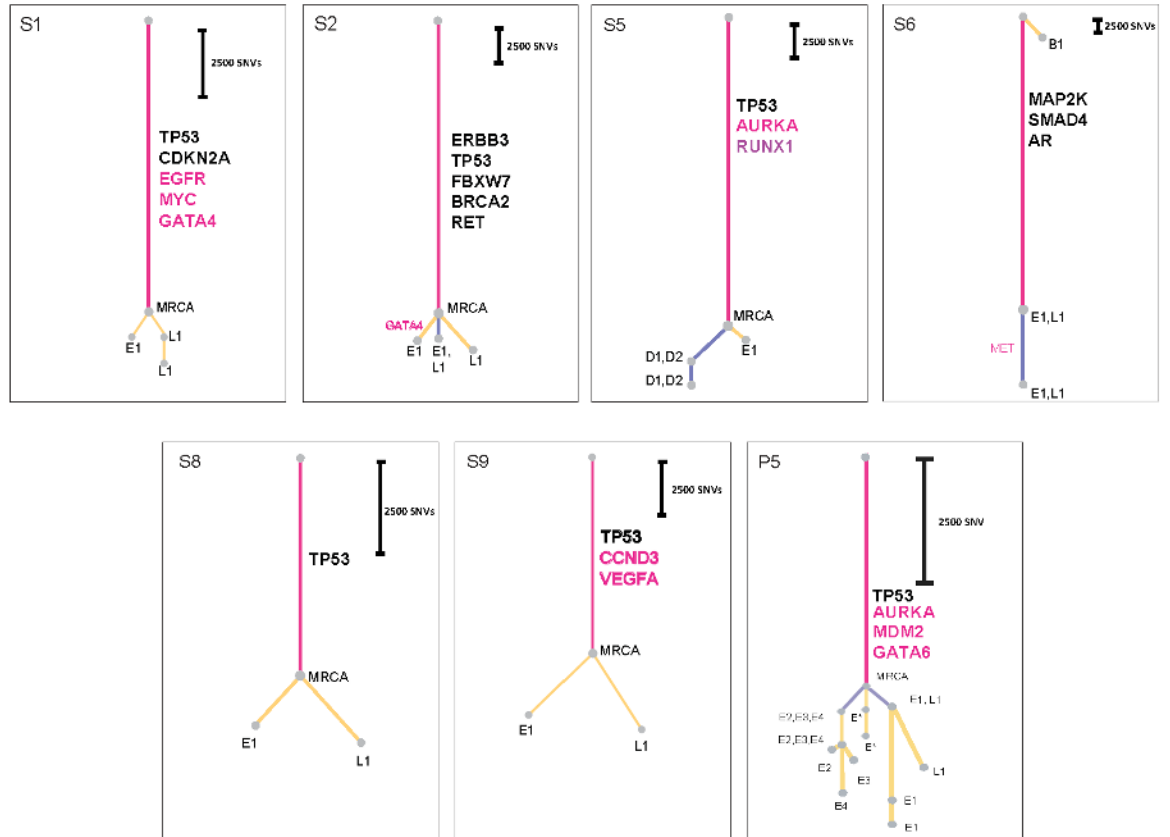
4 Unscaled Tree



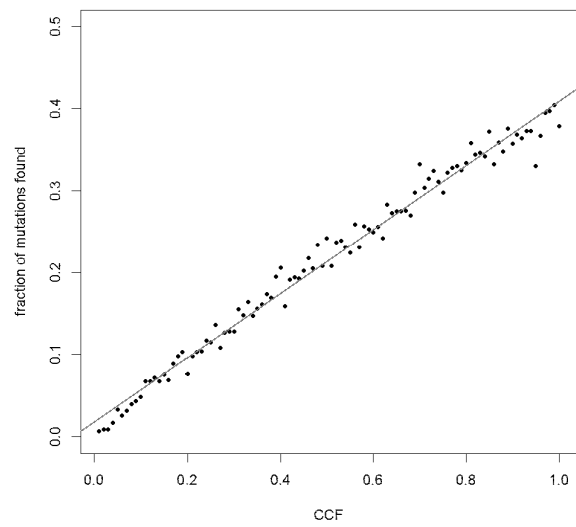
5 Final Tree with scaled branch lengths and Gene Annotation



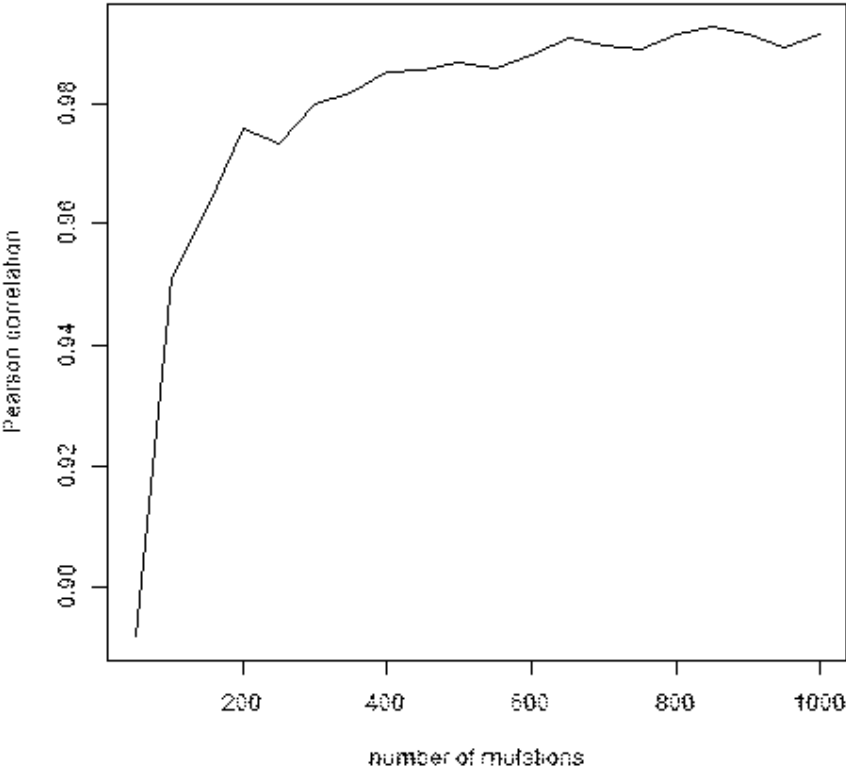
Extended Data 3



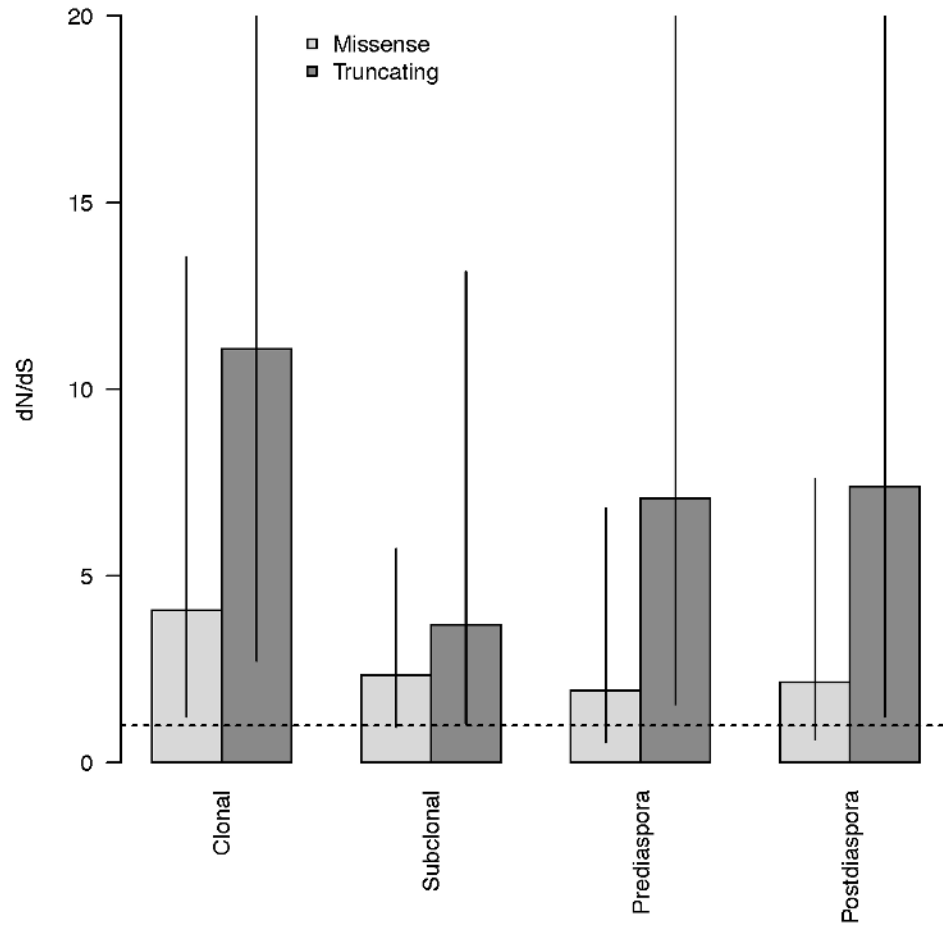
Extended Data5



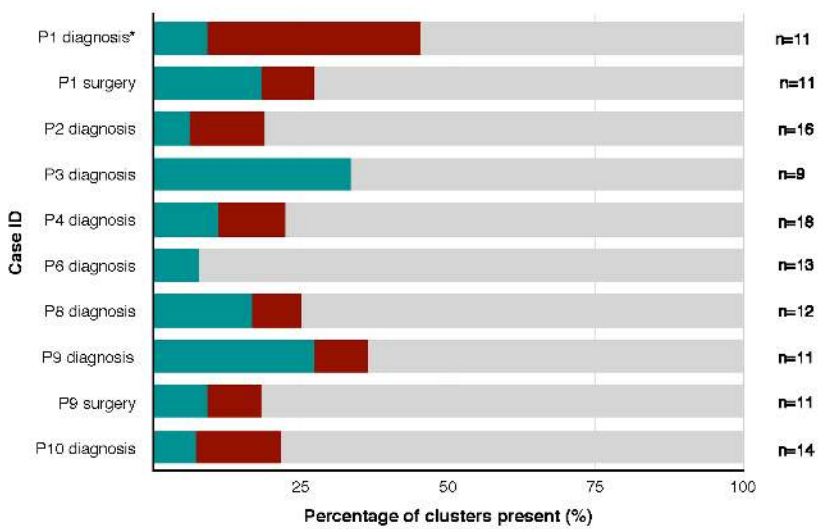
Extended Data6



Extended Data8



Extended Data9



Extended Data10

