

1 Pervasive Chromosomal Instability and 2 Karyotype Order During Tumour Evolution

3 **Authors:**

4 Thomas BK Watkins*¹, Emilia L Lim*^{1,2}, Marina Petkovic^{†3}, Sergi Elizalde^{†4}, Nicolai J Birkbak^{1,5,6},
5 Gareth A Wilson¹, David A Moore^{2,7}, Eva Grönroos¹, Andrew Rowan¹, Sally M Dewhurst⁸, Jonas
6 Demeulemeester^{9,10}, Stefan C Dentre^{9,11,12}, Stuart Horswell¹³, Lewis Au^{1,14}, Kerstin Haase⁹,
7 Mickael Escudero¹³, Rachel Rosenthal^{1,2,15}, Maise Al Bakir¹, Hang Xu¹⁶, Kevin Litchfield¹, Wei
8 Ting Lu¹, Thanos P Mourikis¹⁷, Michelle Dietzen¹⁷, Lavinia Spain^{1,14}, George D Cresswell¹⁸,
9 Dhruva Biswas^{1,15}, Philippe Lamy⁵, Iver Nordentoft⁵, Katja Harbst^{19,20}, Francesc Castro-
10 Giner^{21,22}, Lucy R Yates^{23,24}, Franco Caramia²⁵, Fanny Jaulin²⁶, Cécile Vicier²⁷, Ian PM
11 Tomlinson²⁸, Priscilla K Brastianos²⁹, Raymond J Cho³⁰, Boris C Bastian³⁰, Lars Dyrskjøt⁵, Göran
12 B Jönsson^{19,20}, Peter Savas^{25,31}, Sherene Loi^{25,31}, Peter J Campbell²³, Fabrice Andre³², Nicholas
13 M Luscombe^{18,33,34}, Neeltje Steeghs³⁵, Vivianne CG Tjan-Heijnen³⁶, Zoltan Szallasi^{37,38,39}, Samra
14 Turajlic^{1,14}, Mariam Jamal-Hanjani^{1,40}, Peter Van Loo⁹, Samuel F Bakhoun^{41,42}, Roland F
15 Schwarz^{3,43,44} #, Nicholas McGranahan¹⁷ #, Charles Swanton^{1,2,40} #

16

17

18 * Joint first authors with equal contribution

19 † Joint second authors with equal contribution

20 # Joint corresponding authors

21

22 1 Cancer Evolution and Genome Instability Laboratory, The Francis Crick Institute, 1 Midland Rd,
23 London, NW1 1AT, UK24 2 Cancer Research UK Lung Cancer Centre of Excellence, University College London Cancer
25 Institute, Paul O'Gorman Building, 72 Huntley Street, London, WC1E 6BT, UK26 3 Berlin Institute for Medical Systems Biology, Max Delbrueck Center for Molecular Medicine, Berlin,
27 Germany

28 4 Department of Mathematics, Dartmouth College, Hanover, New Hampshire, USA

29 5 Department of Molecular Medicine (MOMA), Aarhus University Hospital, Aarhus, Denmark

30 6 Bioinformatics Research Centre (BiRC), Aarhus University, Aarhus, Denmark

31 7 Department of Cellular Pathology, University College London Hospitals, London, UK

32 8 Laboratory for Cell Biology and Genetics, Rockefeller University, New York, NY, USA

- 33 9 Cancer Genomics Laboratory, The Francis Crick Institute, 1 Midland Rd, London, NW1 1AT, UK
- 34 10 Department of Human Genetics, University of Leuven, Herestraat 49, 3001 Leuven, Belgium
- 35 11 Oxford Big Data Institute, University of Oxford, Oxford, UK
- 36 12 Experimental Cancer Genetics, Wellcome Trust Sanger Institute, Hinxton, UK
- 37 13 Department of Bioinformatics and Biostatistics, The Francis Crick Institute, London NW1 1AT, UK
- 38 14 Renal and Skin Units, the Royal Marsden Hospital NHS Foundation Trust, London SW3 6JJ, UK
- 39 15 Bill Lyons Informatics Centre, University College London Cancer Institute, Paul O'Gorman Building,
40 72 Huntley Street, London, WC1E 6BT, UK
- 41 16 Stanford Cancer Institute, Stanford, Palo Alto, CA 94304
- 42 17 Cancer Genome Evolution Research Group, University College London Cancer Institute, University
43 College London, London, UK
- 44 18 Bioinformatics and Computational Biology Laboratory, The Francis Crick Institute, London, UK
- 45 19 Division of Oncology and Pathology, Department of Clinical Sciences Lund, Faculty of Medicine,
46 Lund University, Scheelegatan 2, Medicon Village, 22185, Lund, Sweden
- 47 20 Lund University Cancer Centre, Lund University, Lund, Sweden
- 48 21 Department of Biomedicine, Cancer Metastasis Laboratory, University of Basel and University
49 Hospital Basel, CH-4058, Basel, Switzerland
- 50 22 Swiss Institute of Bioinformatics (SIB), Lausanne, Switzerland.
- 51 23 Wellcome Trust Sanger Institute, Hinxton CB10 1SA, UK
- 52 24 Department of Clinical Oncology, Guy's and St Thomas' NHS Foundation Trust, London, SE1 9RT,
53 UK
- 54 25 Division of Research, Peter MacCallum Cancer Centre, University of Melbourne, Melbourne,
55 Victoria, Australia
- 56 26 INSERM U1279, Gustave Roussy, 114 rue Edouard Vaillant, 94805, Villejuif, France.
- 57 27 Department of Medical Oncology, Institut Paoli-Calmettes, Aix-Marseille University, Marseille,
58 France.
- 59 28 Edinburgh Cancer Research Centre, IGMM, University of Edinburgh, Crewe Road South,
60 Edinburgh EH4 2XU
- 61 29 Stephen E. and Catherine Pappas Center for Neuro-Oncology, Divisions of Hematology/Oncology
62 and Neuro-Oncology, Departments of Medicine and Neurology, Massachusetts General Hospital,
63 Harvard Medical School, 55 Fruit Street, Boston, MA, 02114, USA.
- 64 30 Department of Dermatology, University of California, San Francisco, San Francisco, CA, USA
- 65 31 Sir Peter MacCallum Department of Oncology, University of Melbourne, Melbourne, Victoria,
66 Australia
- 67 32 Gustave Roussy Department of Medical Oncology, Faculté de Médecine Paris-Sud XI, Université
68 Paris-Saclay, Villejuif, France
- 69 33 UCL Genetics Institute, Department of Genetics, Evolution & Environment, University College
70 London, UK
- 71 34 Okinawa Institute of Science & Technology, Okinawa, Japan
- 72 35 Department of Medical Oncology, School of GROW, Maastricht University Medical Center,
73 Maastricht, the Netherlands

- 74 36 Department of Medical Oncology, Netherlands Cancer Institute, the Netherlands.
- 75 37 Danish Cancer Society Research Center, Copenhagen, Denmark
- 76 38 Computational Health Informatics Program, Boston Children's Hospital, USA
- 77 39 2nd Department of Pathology, SE-NAP Brain Metastasis Research Group, Semmelweis University,
- 78 Budapest, Hungary
- 79 40 Department of Medical Oncology, University College London Hospitals, London NW1 2BU, UK
- 80 41 Human Oncology and Pathogenesis Program, Memorial Sloan Kettering Cancer Center, New York,
- 81 NY 10065, USA
- 82 42 Department of Radiation Oncology, Memorial Sloan Kettering Cancer Center, New York, NY
- 83 10065, USA
- 84 43 German Cancer Consortium (DKTK), partner site Berlin
- 85 44 German Cancer Research Center (DKFZ), Heidelberg
- 86
- 87
- 88
- 89
- 90

91

92 **Abstract**

93 Cancer chromosomal instability (CIN) results from dynamic changes to chromosome number
94 and structure. The resulting diversity in somatic copy number alterations (SCNA) may provide
95 the variation necessary for cancer evolution. Multi-sample phasing and SCNA analysis of 1421
96 samples from 394 tumours across 24 cancer types revealed ongoing CIN resulting in pervasive
97 SCNA heterogeneity. Parallel evolutionary events, causing disruption to the same genes, such
98 as *BCL9*, *ARNT/HIF1B*, *TERT* and *MYC*, within separate subclones were present in 35% of
99 tumours. Most recurrent losses occurred prior to whole genome doubling (WGD), a clonal
100 event in 48% of tumours. However, loss of heterozygosity at the human leukocyte antigen
101 locus and loss of 8p to a single haploid copy recurred at significant subclonal frequencies, even
102 in WGD tumours, likely reflecting ongoing karyotype remodeling. Focal amplifications
103 affecting 1q21 (*BCL9*, *ARNT*), 5p15.33 (*TERT*), 11q13.3 (*CCND1*), 19q12 (*CCNE1*) and 8q24.1
104 (*MYC*) were frequently subclonal and exhibited an illusion of clonality within single samples.
105 Analysis of an independent series of 1024 metastatic samples revealed enrichment for 14 focal
106 SCNAs in metastatic samples, including late gains of 8q24.1 (*MYC*) in clear cell renal carcinoma
107 and 11q13.3 (*CCND1*) in HER2-positive breast cancer. CIN may enable ongoing selection of
108 SCNAs, manifested as ordered events, often occurring in parallel, throughout tumour
109 evolution.

110

111

112

113

114

115

116

117

118

119 Introduction

120 Chromosomal instability (CIN) results from the occurrence and tolerance of chromosome
121 segregation errors during cell division. CIN has been linked to poor prognosis¹⁻⁴ and leads to
122 somatic copy number alterations (SCNAs) which may act as a substrate for selection⁵⁻⁷.

123

124 However, the prevalence of ongoing CIN later in tumour evolution⁸ and the temporal order of
125 clonal and subclonal SCNAs in relation to whole genome doubling (WGD) events and
126 metastatic dissemination remains unclear.

127

128 Ongoing CIN and SCNA heterogeneity occur across cancer types

129

130 We applied a multi-sample phasing SCNA analysis method (Figure S1,S2A,B, Methods 2.2-6)
131 to 1421 cancer samples from 394 patients across 24 cancer subtypes (range 2-16, median 3
132 samples/tumour Figure S3A,B, Table S1), to obtain SCNA heterogeneity at haplotype
133 resolution. We used MEDICC⁹ to estimate copy number states of the most recent common
134 ancestor (MRCA) for each tumour - reflecting SCNAs acquired prior to subclonal
135 diversification. 1111 / 1421 samples were from treatment naive primary tumours, 51 were
136 from post-treatment primary tumours, 7 samples were obtained at local relapse, and 252
137 were of metastatic origin. In each case, there were at least two samples per tumour with 126
138 tumours having temporally separated samples.

139

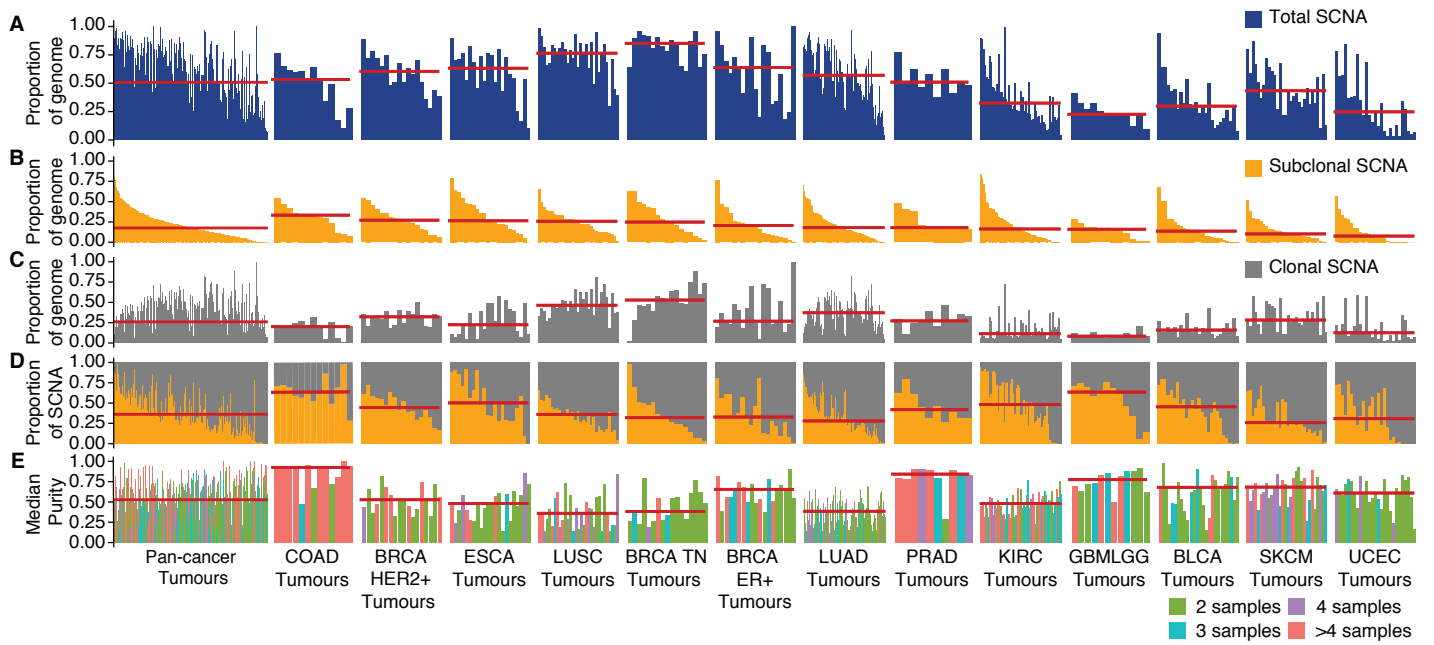
140 To explore CIN during cancer evolution, we quantified the total proportion of the genome
141 affected by SCNAs and the proportion of clonal, early, SCNAs compared to subclonal, late,
142 SCNAs (Figure 1A-D). Clonal SCNAs were identified in every tumour (Figure 1C). 99% of
143 tumours (389/394) harboured at least one subclonal SCNA (Figure 1B). A median of 24% of
144 the genome was subject to clonal SCNAs and 17% subclonal SCNAs. 43% of tumours exhibited
145 >20% of the genome subject to subclonal SCNAs, highlighting that ongoing CIN is pervasive.
146 However, this is likely an underestimate of CIN as only a minority of each tumour is sequenced.

147 Consistent with this, we observed a significant correlation between the number of samples
148 per tumour and SCNA heterogeneity (Figure S4). Moreover, triple-negative breast cancer
149 (BRCA TN), esophageal adenocarcinoma (ESCA) and, clear cell renal cell carcinoma (KIRC)
150 demonstrated a significant association between median purity (Figure 1E) and the proportion
151 of the genome affected by subclonal SCNAs (Figure S5), indicating that purity may impede
152 estimation of SCNA clonality.

153
154 The timing of SCNAs varied across cancers (Figure 1A-B, Figure S6). Despite a comparable total
155 SCNA burden between lung adenocarcinoma (LUAD) and HER2-positive breast cancer (BRCA
156 HER2+) (57% vs. 60%, $P=0.81$, $ES=0.05$), LUAD exhibited a larger proportion of the genome
157 subject to clonal SCNA, whilst BRCA-HER2+ harboured a higher frequency of subclonal SCNA
158 (LUAD: 28% vs BRCA HER2+: 44%, $P=8.1\times 10^{-3}$, $ES=0.59$, analysis also controlled for sample
159 number, see Figure S4B).

160
161 Consistent with increased proliferation in CIN tumours, total, clonal and subclonal SCNA
162 burden correlated with increased cell cycle gene expression in 58 NSCLCs with RNA-
163 sequencing and with increased mitotic index score in 84 NSCLCs with digitised diagnostic slides
164 (Figure S7,S8 Methods 4.3-4, Table S2). Furthermore, in these 84 tumours, estimates of
165 tumour size derived from diagnostic PET–CT scans were found to correlate with total and
166 subclonal SCNA burden, however these associations did not remain significant when
167 controlling sample number (Figure S9). Finally, anisonucleosis, a measure of variation in
168 nuclear size (Methods 4.4), prognostic in NSCLC^{10, 11}, was associated with increased total and
169 clonal SCNA, but not subclonal SCNAs (Figure S10, Table S2).

170
171
172
173
174
175
176
177
178
179
180
181
182



183 **Figure 1 - Overview of somatic copy number heterogeneity across cancer types**

184 A) For each tumour, the proportion of the genome that is affected by SCNA (both clonal and subclonal) is
 185 indicated. Cancer types examined with tumour samples from 10 or more patients included: colorectal
 186 adenocarcinoma (COAD, n=13), HER2+ breast cancer (BRCA HER2+, n=18), esophageal adenocarcinoma (ESCA,
 187 n=22), lung squamous cell carcinoma (LUSC, n=31), triple-negative breast cancer (BRCA TN, n=17), ER+ breast
 188 cancer (BRCA ER+, n=19), lung adenocarcinoma (LUAD, n=84), prostate adenocarcinoma (PRAD, n=10), clear cell
 189 renal cell carcinoma (KIRC, n=54), glioma (GBMLGG, n=12), bladder urothelial carcinoma (BLCA, n=26), melanoma
 190 (SKCM, n=30), and endometrial carcinoma (UCEC, n=27). Cancer types and tumours are ordered by the median
 191 percentage of the genome affected by subclonal SCNA - this order is maintained throughout the figure. Red lines
 192 indicate the median of the distribution. B-C) Barplots indicating the percentage of the genome affected by
 193 subclonal (B) and clonal (C) SCNA. D) The proportion of SCNA that are subclonal and clonal is displayed. Red line
 194 indicates median proportion of SCNA that are subclonal. E) The median purity and number of samples from each
 195 tumour.

196

197 55% of tumours exhibited whole genome doubling (WGD) (Methods 2.7), a clonal event in
 198 87% of cases (Figure S11A). WGD was associated with an increased burden of clonal and
 199 subclonal gains and losses compared to non-WGD tumours (clonal $P=1.36 \times 10^{-34}$, $ES=1.15$;
 200 Subclonal $P=4.67 \times 10^{-9}$, $ES=0.6$, Figure S11B, Methods 2.3). Through multi-sample phasing we
 201 investigated mirrored subclonal allelic imbalance³, resulting from SCNAs disrupting the same
 202 genomic region deriving from distinct haplotypes within separate tumour subclones (Methods
 203 2.2-6). WGD tumours were enriched for mirrored subclonal allelic imbalance events compared
 204 to non-WGD tumours (Methods 2.2-2.6, $P=4.23 \times 10^{-7}$, $ES=0.6$, Figure S11C). In tumours with
 205 subclonal WGD, we observed a higher frequency of SCNAs in subclones affected by WGD
 206 compared to non-WGD sister clones ($P=9.5 \times 10^{-3}$, $ES=0.59$, paired t-test, Figure S11D),
 207 accounting for germline and prior somatic alterations as confounding variables.

208 ***Evolution of the SCNA landscape***

209 To address whether the SCNA landscape is shaped by neutral evolution or selection, we
 210 considered whether the propensity for each chromosome arm to be gained or lost during
 211 tumour evolution was related to the density of tumour suppressor genes (TSGs) and
 212 oncogenes (OGs) encoded on each chromosome arm, as captured by the OG-TSG score⁵.
 213 Consistent with ongoing selection on cellular karyotypes, the OG-TSG score significantly
 214 correlated with the burden of arm-level alterations in the MRCA (Figure 2A) as well as
 215 subclonal arm-level alterations (Figure 2B and Figure S12A-C). No relationship between
 216 average clonal or subclonal chromosome copy number change and chromosome size was
 217 observed, suggesting SCNA detection is unlikely to contribute to this relationship (Figure
 218 S12D-G).

219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261

To understand subclonal SCNA dynamics within each tumour, we adapted our previous model that predicts population karyotypes over time^{12, 13}. We used arm-level copy number profiles from each tumour's MRCA as the starting point and compared how different iterations of the model predict the observed subclonal tumour karyotypes (Figure 2C, S13A,B, Methods 3.1-2). We compared three conditions; first, where karyotypes with higher oncogenic propensity or tumour suppressive propensity were favoured, or unfavoured, respectively, using the relative OG-TSG scores⁵ (weighted model); second, where chromosome arms were treated equally (neutral model); and, third where OG-TSG scores were randomly permuted (scrambled model)) to achieve the same complexity as the weighted model. On average, the weighted model predicted the trajectory of subclonal SCNA more accurately, outperforming the two other models, as evidenced by significantly reduced deviance scores (Figure 2C,D, Figure S13C-G) irrespective of the rate of chromosome missegregation or the number of cell divisions (generations)(Figure S14).

262 **Figure 2 - Selection shapes the SCNA landscape**

263 A) Scatter plot showing a positive correlation between average clonal copy number present in the MRCA and the
 264 OG-TSG score in n=394 tumours. The grey shaded area represents the 95% confidence interval. Rho and p are
 265 from a Spearman correlation test. B) Scatter plot showing a positive correlation between OG-TSG score and
 266 average change in SCNA (gain or loss) from MRCA in n=394 tumours. The grey shaded area represents the 95%
 267 confidence interval. Rho and P are from a Spearman correlation test. C) Schematic showing the three conditions
 268 under which karyotype evolution was modelled: chromosome arms incorporating OG-TSG scores (weighted
 269 model); chromosome arms treated equally (neutral model); OG-TSG scores randomly permuted (scrambled
 270 model). D) For each context (non-WGD n=194, WGD n=171, and subclonal WGD n=29), the percentage of
 271 tumours in which each model condition best recapitulates the empirically observed data is displayed in the bar
 272 chart.

273
 274

275 Collectively, these data suggest that CIN enables ongoing selection driven by relative dosage
 276 imbalance of OGs and TSGs and that WGD may support ongoing genome remodeling later
 277 during tumour evolution, permitting further selection. However, the observed pattern of
 278 SCNA acquisition in 41% of our cohort in which the neutral or scrambled models outperform
 279 the weighted model might reflect neutral karyotype evolution or the need for cancer-type
 280 specific chromosome arm weightings^{14, 15}. Notably, we see more evidence for subclonal
 281 selection in WGD tumours which may be consistent with WGD being a transformative event
 282 during subclonal evolution (Figure 2D, S13F,G)^{12, 13}.

283

284 ***Clonal SCNA recur across cancer types and losses are predominantly early***

285

286 To decipher SCNA timing during evolution, we used GISTIC2.0 to identify recurrent SCNAs in
 287 at least two cancer types (Methods 2.12-14, Figure S15A-M, Table S3). We designed these as
 288 consensus peak regions and assigned each into distinct evolutionary timing categories: early,
 289 intermediate, or late (Figure 3A,B, Methods 2.15). SCNAs overlapping early peak regions may
 290 be implicated in tumourigenesis or result from specific constraints to tumourigenesis. SCNAs
 291 overlapping intermediate or late peak regions may be involved in tumour maintenance and
 292 progression. Recurrent clonal and subclonal arm-level gain or loss SCNAs for each cancer type
 293 were identified using permutations (Methods 2.8-9, Table S4).

294

295 We observed differences in evolutionary timing between peak regions associated with gains
 296 (gain-peaks) and those with losses (loss-peaks). Loss-peaks were significantly more likely to be
 297 early compared to gain-peaks ($P=6.8 \times 10^{-8}$, $ES=0.57$, Figure S16). Similarly, a higher proportion
 298 of recurrent arm-level losses were clonal compared to arm-gains ($P=2.8 \times 10^{-9}$, $ES=0.77$, Figure

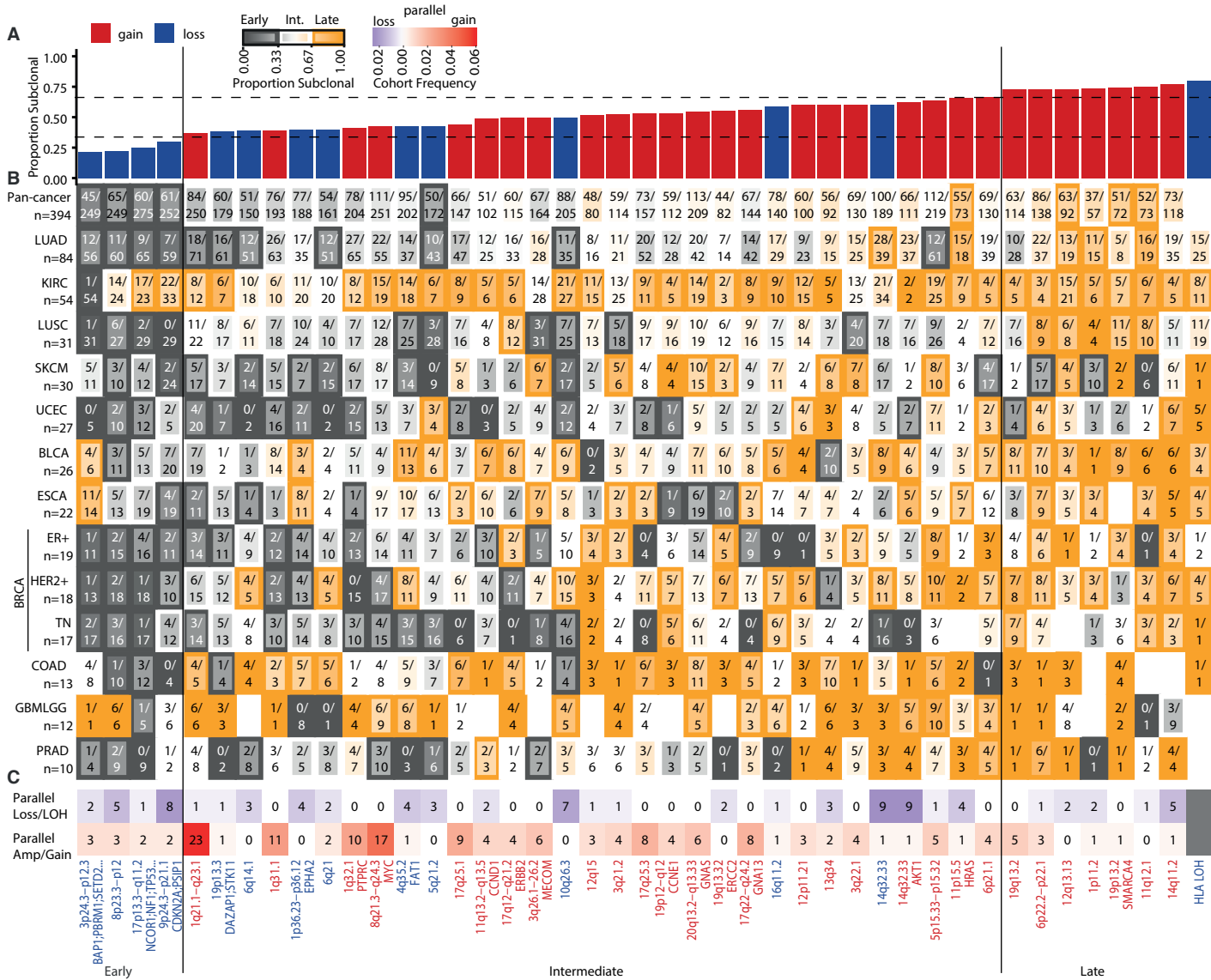
299 S17). Gain peak regions were enriched for known OGs, while loss peak regions were enriched
300 for known TSGs (Figure S18A). Early loss peak regions were also enriched for chromosomal
301 fragile sites (Figure S18B), suggesting some may lack functional significance.

302
303 Clonal SCNA frequencies affecting early peak regions exceeded the frequency of clonal
304 somatic driver point mutations and indels (Figure S19). The loss-peak 17p13.3–q11.2,
305 encompassing *TP53*, was classified as early in 9/13 cancer types and only classified as late in
306 KIRC (74% subclonal). In three cancer types (BRCA HER2+, LUSC and BRCA TN) >90% of cases
307 exhibited clonal LOH at 17p13.1, suggesting loss is required for tumourigenesis. Across
308 cancers, *TP53* LOH, was clonal in 90% of WGD cases in which it was observed, possibly
309 permitting tolerance for WGD¹⁶. In KIRC, only 3p26.3–p12.1, encompassing *VHL*, was early
310 (clonal LOH in 98% of KIRCs) (Figure S15H). Other high frequency clonal peaks within cancer
311 types included, the gain-peak 17q12–q21.2, encompassing *ERBB2* in BRCA HER2+ (61%
312 frequency), 3p LOH in LUSC (97% frequency), and gain of 7p11.2, encompassing *EGFR*, in LUAD
313 (64% frequency).

314
315 We reasoned that loss occurring prior to WGD must lead to LOH with complete loss of the
316 minor allele. Conversely, single losses occurring after WGD will not lead to LOH. On average,
317 across the cohort, 92% of clonal losses overlapping early loss-peaks involved LOH, suggesting
318 recurrent clonal loss events usually occur prior to WGD.

319
320 The timing of other peak-regions were promiscuous between cancer types (Figure 3B). For
321 example, the loss-peak 4q35.2, encompassing *FAT1*, was early in BRCA TN (89% prevalence,
322 80% clonal), intermediate in BRCA ER+ (58% prevalence, 64% clonal) and late in BRCA HER2+
323 (61% prevalence, 73% subclonal) (Figure 3B).

324
325
326
327
328
329
330
331
332
333
334



335 **Figure 3 – Timing, Recurrence and Parallel Evolution of Subclonal SCNAs**

336 A) Barplot of consensus peaks (Methods 2.13) ordered by median percentage of subclonal occurrence across
 337 cancer types. Bars representing gain-peaks are coloured in red and loss-peaks are coloured in blue. Vertical lines
 338 indicate separation of consensus peaks into pan-cancer categories of early, intermediate and late, according to
 339 tertiles of median proportion of SCNA that is subclonal (horizontal dashed lines). B) Heatmap of the percentage
 340 subclonal occurrence of all consensus peaks in each cancer type. Numerator within each cell indicates, in that
 341 cancer type, the total number of subclonal occurrences of that peak region and the denominator indicates the
 342 total number of both clonal and subclonal occurrences of that consensus peak in that cancer type. Shading of
 343 each cell in the heatmap indicates the percentage subclonal occurrence of a consensus peak within a cancer type
 344 with orange indicating higher subclonality and grey indicating higher clonality. The border of each cell indicates
 345 the classification of that consensus peak in a cancer type as either early (thick dark grey border), intermediate
 346 (no border) or late (thick dark orange border). C) Heatmap showing the number of instances of parallel evolution
 347 of loss/LOH (blue) and gain/amplification (red) affecting the consensus peak regions. Shading of each cell
 348 indicates the number of occurrences of parallel evolution and the number within the cell states the number of
 349 such parallel evolutionary events.

350

351

352 **Evolution of Subclonal SCNAs**

353

354 We next addressed which specific subclonal SCNAs are recurrent during tumour evolution.

355

356 The highest frequency gain-peaks, including 1q21.1-q21.3 (encoding *BCL9*, *ARNT/HIF1B*) and
 357 5p15.33-p15.32 (encompassing *TERT*), varied in timing across cancers. For example, in LUAD,
 358 80% of 5p15.33-p15.32 gains were clonal, while the majority were subclonal in KIRC (76%
 359 subclonal), BRCA ER+ (89% subclonal) and GBMLGG (90% subclonal) (Figure 3B). In LUSC, the
 360 timing of *TERT* gains was related to both its focality and amplitude; the majority of low-level
 361 gains (>ploidy & <2x ploidy, Methods 2.3) were both clonal and arm-level (8/14) while high-
 362 level *TERT* amplifications were often subclonal and focal (10/11). This may reflect
 363 augmentation of gene dosage during evolution, with low-level *TERT* gain selected clonally,
 364 followed by a high-level amplification selected in a subset of cancer cells later in tumour
 365 evolution.

366

367 The gain-peak of 19p12–q12 (encompassing *CCNE1*) was late or intermediate in 10/13 cancer
 368 types. High-level amplifications of *CCNE1* (>2x ploidy), previously associated with WGD^{6, 17},
 369 occurred exclusively in WGD tumours. *CCNE1* amplification was subclonal in 9/20 tumours
 370 with clonal WGD, suggesting it may be selected both before and after WGD.

371

372 Parallel evolution of SCNA events, reflecting distinct subclones within individual tumours
373 converging on a similar evolutionary solution, was observed in 139/394 (35%) tumours within
374 the cohort (Figure 3C,S20). Allele-specific expression tracked parallel evolutionary events
375 originating from distinct haplotypes in samples with matched multi-sample RNA-seq (Figure
376 S21, $\rho = 0.89$, $P=1.75 \times 10^{-15}$, Spearman correlation).

377
378 Consistent with positive selection, parallel gains were significantly more focal than non-
379 parallel subclonal gains ($P=5.9 \times 10^{-3}$, $ES=0.1$). The most prominent parallel gains included those
380 overlapping 1q21.3-q44 encompassing *BCL9* and *ARNT/HIF1B*, 5p15.33 encompassing *TERT*
381 and 8q24.1 encompassing *MYC* (Figure 3C, Figure S20). The most common parallel loss events
382 included 14q (14q32.33/*ASPP1* and 14q11.2/*NDRG2*), 10q and 9p (Figure S20).

383
384 Subclonal LOH after a clonal WGD event occurs through more than one loss event of the same
385 allele after the doubling event (Figure S22). The HLA locus (6p21.3) represented a clear peak
386 of subclonal LOH in WGD samples, affecting 22% of the cohort, indicative of two loss events
387 of the same alleles after genome doubling within the subclone (Figure S23). HLA LOH was
388 prevalent as a subclonal event in KIRC, BRCA, BLCA, NSCLC-other, UCEC and ESCA (Figure S24,
389 Methods 2.11) in addition to NSCLC as previously reported¹⁸. One exception was SKCM, which
390 is characterised by high mutational burden and benefits from checkpoint inhibitor blockade
391 ¹⁹. SKCM exhibited the lowest frequency of HLA-LOH (0% clonal, 4% subclonal) in the cohort.
392 6p24.2 also harbours the melanoma metastasis gene *NEDD9*²⁰, identified as the most
393 prevalent recurrent clonal arm-level gain event in SKCM, which may constrain subsequent HLA
394 loss (Figure S15L).

395
396 In a diploid cancer cell, any loss results in LOH. If this cell doubles, the LOH will be maintained,
397 with the remaining allele being duplicated, leading to a total copy number of two.
398 Interestingly, in the case of clonal 8p23.3-p12 loss, we observed a peak region of haploid LOH
399 in WGD tumours, with only a single copy (Figure S22). This haploid, single copy, LOH strongly
400 suggests a loss event of one of the two remaining copies after WGD. Loss of 8p23.3-p12 was
401 most prominent in breast cancer where it has been linked to a chromosome-dosage effect,
402 influencing lipid metabolism and metastatic potential²¹.

403

404 We next investigated whether we could identify an association between the presence of SCNA
405 overlapping our consensus peaks and overall survival in cancer types matched to those in our
406 multi-sample cohort from the TCGA (Figure S25). Few individual consensus SCNA were
407 associated with survival, suggesting the binary presence or absence of individual SCNAs is
408 rarely prognostic.

409

410

411 ***Late emerging subclones frequently seed metastases***

412

413 Next, we explored associations of SCNAs with metastasis. Consistent with previous work²², we
414 observed a greater percentage of the genome affected by SCNAs in metastatic (n=137
415 patients) compared with primary tumour samples (n=373 patients) (Figure S26A, $P=1.5\times 10^{-5}$,
416 $ES=0.4$). This remained significant when controlling for cancer type and considering both
417 paired and unpaired primary-metastasis tumour comparisons (Figure S26B) with LOH events
418 displaying the greatest increase from primary to metastasis compared to gain or losses
419 without LOH (Figure S26C). No significant increase in ploidy was observed between primary
420 and metastatic samples in the cohort or in any individual cancer types examined.

421

422 Consistent with an evolutionary bottleneck, SCNAs were found to be more frequently clonal
423 in metastatic compared to primary samples (Figure S26D). Indeed, in all 14 cases where we
424 had multi-sample primary tumours and a matched metastatic sample, we identified SCNAs
425 which were fully clonal in the metastasis and present as minor subclones within the primary.
426 In all cancer types with multiple primary-metastatic pairs, in the majority of tumours, most
427 LOH events were found to be shared between primary and metastatic samples, suggesting a
428 relatively late divergence of the metastatic clone assuming LOH events occur at a constant
429 rate throughout cancer evolution (Figure S27, Methods 4.9).

430

431 To evaluate the relative importance of specific SCNAs in metastasis, we focused on recurrent
432 SCNAs and performed a combined analysis using both paired analyses on 118 matched
433 primary-metastatic samples and unpaired analyses on 2631 TCGA primary samples, and 1024
434 Hartwig Medical Foundation metastatic samples²³ in the four cancer types with sufficient
435 primary-metastatic pairs (BRCA HER2+, BRCA ER+, LUAD, KIRC). However, distinct patterns of

436 SCNA metastatic dissemination were observed in different cancer types. In BRCA ER+, BRCA
437 HER2+ and LUAD, the majority of the arm-events that were enriched in metastasis relative to
438 primary samples were early (Figure S28). Conversely, in KIRC, which had the lowest proportion
439 of shared LOH between primary and metastatic samples, most recurrent arm-events enriched
440 in metastatic samples were classified as intermediate or late events (Figure S28), suggesting
441 these arm-events are associated with metastatic potential in a limited number of cells within
442 the primary tumour.

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

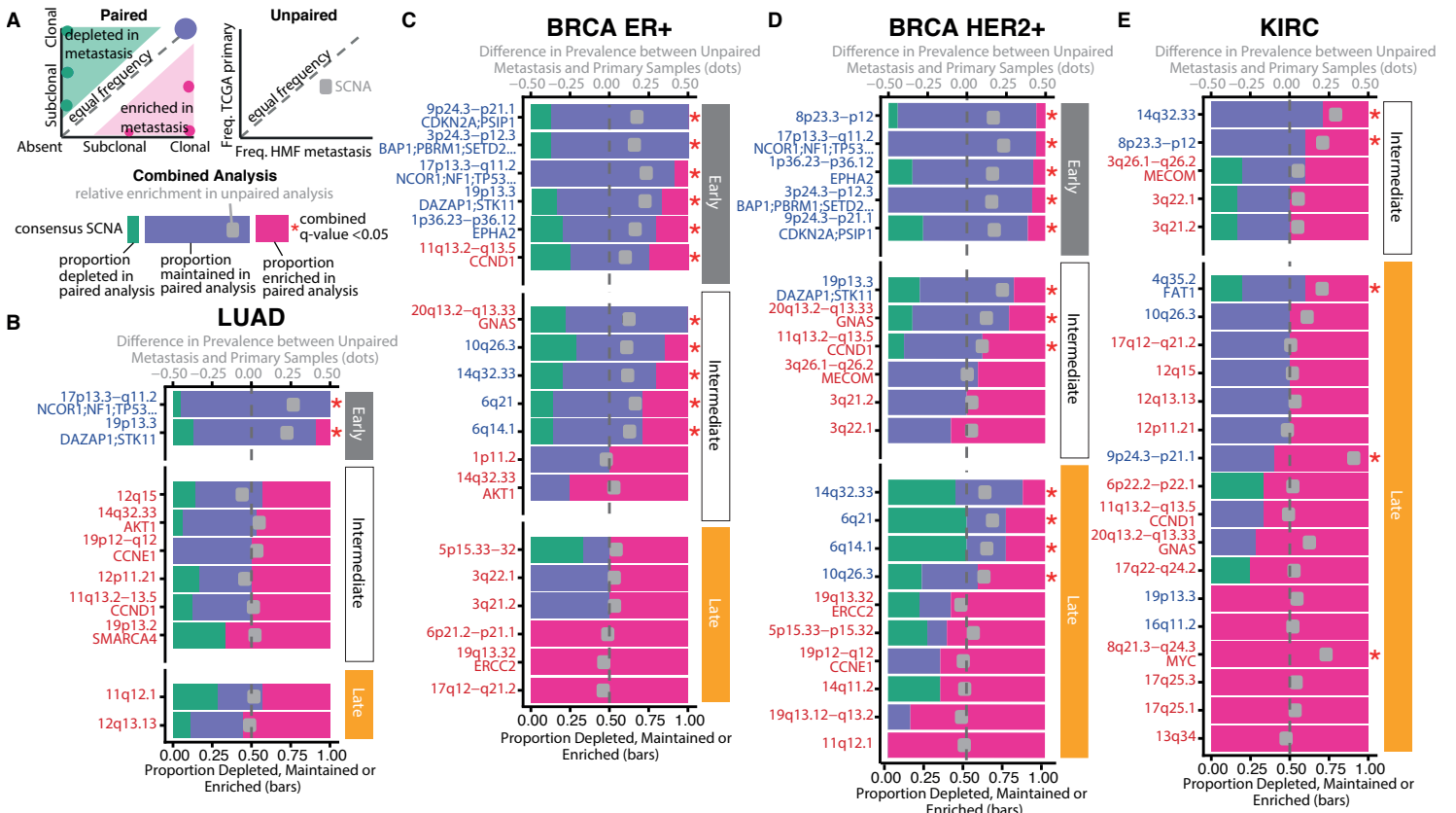
476

477

478

479

480



481 **Figure 4: Analysis of consensus peak regions in metastatic LUAD, BRCA ER+, BRCA HER2+, and KIRC.** A)
 482 Schematic describing the paired (left graph), unpaired analysis (right graph), and combined (barplot below) of
 483 consensus peak regions. The schematic barplot summarises the left graph for each peak consensus region and
 484 indicates the proportion of paired primary-metastasis cases where a SCNA overlapping the consensus peak
 485 region was enriched (pink), depleted (green) or maintained (blue) in metastatic samples. These data were
 486 assessed using a two-sided binomial test. The grey square in the schematic bar plot indicates the difference
 487 between proportions of metastatic (Hartwig Medical Foundation) and primary (TCGA) samples that harbour the
 488 event in the unpaired primary-metastasis analysis (two-sided test of equal or given proportions) - a positive
 489 number indicates that the event was more prevalent in the metastatic (Hartwig Medical Foundation) samples,
 490 while a negative number indicates that the event was more prevalent in the primary (TCGA) samples. The red
 491 stars indicate if an event was significantly enriched in metastatic samples as determined by a combined analysis
 492 of paired (multi-sample) and unpaired (Hartwig Medical Foundation and TCGA) data using Fisher's method after
 493 multiple testing correction using the Benjamini-Hochberg method. The event timing classifications (Early,
 494 Intermediate or Late) were determined based on proportion of subclonal occurrence (Methods 2.15). Only losses
 495 (blue text) or gains (red text) which are either significant ($q < 0.05$) or exhibit $\geq 40\%$ enrichment are shown. We
 496 restricted our analysis to cancer subtypes with ≥ 10 primary-metastasis paired samples (LUAD paired $n=28$,
 497 unpaired $n=844$ TCGA; 315 Hartwig (B), BRCA ER+ paired $n=17$, unpaired $n=1015$ TCGA; 620 Hartwig (C), BRCA
 498 HER2+ paired $n=13$, unpaired $n=1015$ TCGA; 620 Hartwig (D), and KIRC paired $n= 10$, unpaired $n=772$ TCGA; 89
 499 Hartwig (E)).

500

501

502 The early clonal loss-peak at chromosome 1p36.23–p36.12, which encompasses *EPHA2*, and
 503 early clonal loss-peak at chromosome 17p13.3-q11.2, encoding *TP53*, were found to be
 504 enriched in metastatic samples compared to primary samples in BRCA ER+ and BRCA HER2+
 505 (Figure 4). In LUAD, we observed two early loss consensus peak regions significantly enriched
 506 in metastasis (17p13.3-q11.2 [*TP53*], and 19p13.3 encompassing *STK11*), consistent with these
 507 events engendering phenotypes permissive for dissemination early in tumour evolution.

508

509 In contrast, other consensus peak-regions enriched in metastases were classified as
 510 intermediate or late (Figure 4). Examples include losses of 14q32.33, 6q21 (encompassing
 511 *PRDM1*), 6q14.1 and 10q26.3 (encompassing *MGMT*) in BRCA HER2+, and losses of 4q35.2
 512 (encompassing *FAT1*), 9p24.3-p21.1 and gain of 8q21.3-q24.3 in KIRC. Gain of 8q21.3-q24.3,
 513 encompassing *MYC*, was highly enriched in our combined analysis as well as exclusively
 514 identified in the metastatic samples of our matched primary-metastatic pairs in KIRC.
 515 Intriguingly, loss of 9p24.3-p21.1, which encompasses *CDKN2A*, was a late metastasis-
 516 associated event in KIRC, while in ER+ and HER2+ BRCA, where loss of 9p24.3-p21.1 was also
 517 significantly associated with metastasis, it was predominantly early. Similarly, 11q13.2-q13.5,
 518 which encompasses *CCND1*, was an early event in BRCA ER+ and intermediate in BRCA HER2+
 519 and associated with metastasis in both cancer types.

520

521 Taken together, these highlight the importance of both early and ongoing SCNA acquisition
522 during tumour evolution and their potential importance during the metastatic transition.

523

524 **Discussion**

525

526 Clonal and subclonal SCNAs are pervasive across cancers and show a propensity for order,
527 potentially reflecting the ongoing optimization of fitness landscapes throughout cancer
528 evolution. WGD is a transformative event in tumourigenesis, associated with clonal and
529 subclonal SCNA acquisition. LOH events affecting tumour suppressor genes, including *TP53*,
530 preceded WGD and recurrent gains (eg *CCNE1*) frequently followed WGD and were more likely
531 subclonal.

532

533 The subclonal landscape of SCNA is sculpted by both positive and negative selection, as well
534 as neutral evolution. In a minority of tumours, our results are consistent with subclonal
535 karyotypic evolution may predominantly reflecting neutral growth^{14, 15}. However, particularly
536 in tumours with WGD, SCNA evolution was better recapitulated using models incorporating
537 both positive and negative selection (Figure 2D). Positive selection was further evidenced by
538 recurrent peaks of subclonal amplifications, enriched for established oncogenes, subclonal
539 losses resulting in LOH, even after genome doubling, and parallel evolution of SCNAs. These
540 data are consistent with previously documented parallel and convergent evolution of SCNAs³.
541 ²⁴⁻²⁶ Finally, recurrent focal subclonal SCNAs, encompassing oncogenic events including *CCND1*
542 and *MYC* were enriched at metastatic sites suggesting a potential role in metastasis.
543 Consistent with this, *MYC* was recently described as a SCNA driver of brain metastasis in
544 LUAD²⁷. While certain SCNA were enriched in metastasis, in the majority of tumours, most
545 LOH events were shared between primary and metastatic samples, suggesting a late
546 divergence of the metastatic clone, and or negative selection against extensive loss after the
547 emergence of the MRCA²⁸.

548

549 Our work is not without limitations. Detection of recurrent SCNAs is not solely indicative of
550 selection for functional advantage and may result from other processes driving tumour
551 progression such as DNA repair dysfunction or the presence of adjacent fragile sites. Indeed,

552 the higher frequency of recurrent SCNA compared to driver point-mutations need not reflect
553 selection. However, we only found an association of fragile sites with early loss peak regions.
554 Extrachromosomal DNA may also contribute to the subclonal SCNA amplification events
555 observed²⁹. The number of tumour samples, their sequencing depths and the lack of an
556 extensive paired primary-metastasis cohort or single cell sequencing analysis influence the
557 degree to which subclonal heterogeneity can be deciphered, suggesting the extent of diversity
558 is underestimated. The lack of uniform clinical data collection and central pathology review
559 prevented detailed analysis of clinically relevant parameters. We are endeavoring to address
560 these deficiencies to time metastatic dissemination events and clonal expansions within
561 TRACERx³.

562
563 In conclusion, our work highlights the importance of ongoing chromosomal instability during
564 cancer evolution and metastasis. As our understanding of the propensity for different
565 chromosomes to mis-segregate³⁰ and extent to which chromosomal alterations may be
566 deleterious or advantageous to the cancer cell improves³¹, it will be possible to refine the
567 parameters of selection models and improve the ability to detect novel SCNA drivers, which
568 may drive metastatic dissemination and death.

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584 References

585

586

587 1. McGranahan, N., et al., *Cancer chromosomal instability: therapeutic and diagnostic*
588 *challenges. 'Exploring aneuploidy: the significance of chromosomal imbalance' review*
589 *series*. EMBO Rep, 2012.

590 2. Schwarz, R.F., et al., *Spatial and temporal heterogeneity in high-grade serous ovarian*
591 *cancer: a phylogenetic analysis*. PLoS Med, 2015. **12**(2): p. e1001789.

592 3. Jamal-Hanjani, M., et al., *Tracking the Evolution of Non-Small-Cell Lung Cancer*. N Engl
593 J Med, 2017. **376**(22): p. 2109-2121.

594 4. Hieronymus, H., et al., *Tumor copy number alteration burden is a pan-cancer*
595 *prognostic factor associated with recurrence and death*. Elife, 2018. **7**.

596 5. Davoli, T., et al., *Cumulative haploinsufficiency and triplosensitivity drive aneuploidy*
597 *patterns and shape the cancer genome*. Cell, 2013. **155**(4): p. 948-62.

598 6. Zack, T.I., et al., *Pan-cancer patterns of somatic copy number alteration*. Nat Genet,
599 2013. **45**(10): p. 1134-1140.

600 7. Turajlic, S., et al., *Deterministic Evolutionary Trajectories Influence Primary Tumor*
601 *Growth: TRACERx Renal*. Cell, 2018. **173**(3): p. 595-610 e11.

602 8. Bolhaqueiro, A.C.F., et al., *Ongoing chromosomal instability and karyotype evolution*
603 *in human colorectal cancer organoids*. Nat Genet, 2019. **51**(5): p. 824-834.

604 9. Beerenwinkel, N., et al., *Cancer evolution: mathematical models and computational*
605 *inference*. Syst Biol, 2014.

606 10. von der Thusen, J.H., et al., *Prognostic significance of predominant histologic pattern*
607 *and nuclear grade in resected adenocarcinoma of the lung: potential parameters for a*
608 *grading system*. J Thorac Oncol, 2013. **8**(1): p. 37-44.

609 11. Kadota, K., et al., *Comprehensive pathological analyses in lung squamous cell*
610 *carcinoma: single cell invasion, nuclear diameter, and tumor budding are independent*
611 *prognostic factors for worse outcomes*. J Thorac Oncol, 2014. **9**(8): p. 1126-39.

612 12. Laughney, A.M., et al., *Dynamics of Tumor Heterogeneity Derived from Clonal*
613 *Karyotypic Evolution*. Cell Rep, 2015. **12**(5): p. 809-20.

614 13. Elizalde, S., A.M. Laughney, and S.F. Bakhoun, *A Markov chain for numerical*
615 *chromosomal instability in clonally expanding populations*. PLoS Comput Biol, 2018.
616 **14**(9): p. e1006447.

617 14. Sottoriva, A., et al., *A Big Bang model of human colorectal tumor growth*. Nat Genet,
618 2015. **47**(3): p. 209-16.

619 15. Williams, M.J., et al., *Identification of neutral tumor evolution across cancer types*. Nat
620 Genet, 2016. **48**(3): p. 238-244.

621 16. Fujiwara, T., et al., *Cytokinesis failure generating tetraploids promotes tumorigenesis*
622 *in p53-null cells*. Nature, 2005. **437**(7061): p. 1043-7.

623 17. Bielski, C.M., et al., *Genome doubling shapes the evolution and prognosis of advanced*
624 *cancers*. Nat Genet, 2018. **50**(8): p. 1189-1195.

625 18. McGranahan, N., et al., *Allele-Specific HLA Loss and Immune Escape in Lung Cancer*
626 *Evolution*. Cell, 2017.

627 19. Vogelstein, B., et al., *Cancer genome landscapes*. Science, 2013. **339**(6127): p. 1546-
628 58.

629 20. Kim, M., et al., *Comparative oncogenomics identifies NEDD9 as a melanoma metastasis*
630 *gene*. Cell, 2006. **125**(7): p. 1269-81.

- 631 21. Cai, Y., et al., *Loss of Chromosome 8p Governs Tumor Progression and Drug Response*
632 *by Altering Lipid Metabolism*. *Cancer Cell*, 2016. **29**(5): p. 751-766.
- 633 22. Bakhoun, S.F., et al., *Chromosomal instability drives metastasis through a cytosolic*
634 *DNA response*. *Nature*, 2018. **553**(7689): p. 467-472.
- 635 23. Priestley, P., et al., *Pan-cancer whole-genome analyses of metastatic solid tumours*.
636 *Nature*, 2019. **575**(7781): p. 210-216.
- 637 24. Lackner, C., et al., *Convergent Evolution of Copy Number Alterations in Multi-Centric*
638 *Hepatocellular Carcinoma*. *Sci Rep*, 2019. **9**(1): p. 4611.
- 639 25. Jakubek, Y.A., et al., *Large-scale analysis of acquired chromosomal alterations in non-*
640 *tumor samples from patients with cancer*. *Nat Biotechnol*, 2020. **38**(1): p. 90-96.
- 641 26. Zaccaria, S.R., B. J., *Characterizing the allele- and haplotype-specific copy number*
642 *landscape of cancer genomes at single-cell resolution with CHISEL*. *bioRxiv*
643 <https://doi.org/10.1101/837195>, 2019.
- 644 27. Shih, D.J.H., et al., *Genomic characterization of human brain metastases identifies*
645 *drivers of metastatic lung adenocarcinoma*. *Nat Genet*, 2020. **52**(4): p. 371-377.
- 646 28. William Cross, M.M., Salpie Nowinski, George Cresswell, Abhirup Banerjee, Marc
647 Williams, Laura Gay, Ann-Marie Baker, Christopher Kimberley, Hayley Davis, Pierre
648 Martinez, Maria Traki, Viola Walther, Kane Smith, Giulio Caravagna, Sasikumar
649 Amarasingam, George Elia, Alison Berner, Ryan Changho Choi, Pradeep Ramagiri,
650 Ritika Chauhan, Nik Matthews, Jamie Murphy, Anthony Antoniou, Susan Clark, Jo-Anne
651 Chin Aleong, Enric Domingo, Inmaculada Spiteri, Stuart AC McDonald, Darryl Shibata,
652 Miangela M Lacle, Lai Mun Wang, Morgan Moorghen, Ian PM Tomlinson, Marco
653 Novelli, Marnix Jansen, Alan Watson, Nicholas A Wright, John Bridgewater, Manuel
654 Rodriguez-Justo, Hemant Kocher, Simon J Leedham, Andrea Sottoriva, Trevor A
655 Graham, *Stabilising selection causes grossly altered but stable karyotypes in metastatic*
656 *colorectal cancer*. *bioRxiv*, 2020. doi: <https://doi.org/10.1101/2020.03.26.007138>.
- 657 29. Turner, K.M., et al., *Extrachromosomal oncogene amplification drives tumour*
658 *evolution and genetic heterogeneity*. *Nature*, 2017. **543**(7643): p. 122-125.
- 659 30. Worrall, J.T., et al., *Non-random Mis-segregation of Human Chromosomes*. *Cell Rep*,
660 2018. **23**(11): p. 3366-3380.
- 661 31. Lopez, S., et al., *Interplay between whole-genome doubling and the accumulation of*
662 *deleterious alterations in cancer evolution*. *Nat Genet*, 2020. **52**(3): p. 283-293.
- 663