

1 **Landscape of G-quadruplex DNA structural regions in breast cancer**

2  
3 **Authors:**

4 Robert Hänsel-Hertsch<sup>1,2,3,4</sup>, Angela Simeone<sup>4</sup>, Abigail Shea<sup>4</sup>, Winnie W.I. Hui<sup>4</sup>, Katherine G. Zyner<sup>4</sup>,  
5 Giovanni Marsico<sup>4</sup>, Oscar M. Rueda<sup>4</sup>, Alejandra Bruna<sup>4</sup>, Alistair Martin<sup>4</sup>, Xiaoyun Zhang<sup>5</sup>, Santosh  
6 Adhikari<sup>5</sup>, David Tannahill<sup>4</sup>, Carlos Caldas<sup>4,6,7</sup> and Shankar Balasubramanian<sup>4,5,8</sup>

7  
8 <sup>1</sup>Center for Molecular Medicine Cologne, University of Cologne, Germany.

9 <sup>2</sup>Cologne Excellence Cluster on Cellular Stress Responses in Aging-Associated Diseases (CECAD),  
10 University of Cologne and University Hospital Cologne, Cologne, Germany.

11 <sup>3</sup>Institute of Human Genetics, University Hospital Cologne, Cologne, Germany.

12 <sup>4</sup>Cancer Research UK Cambridge Institute and Department of Oncology, Li Ka Shing Centre,  
13 University of Cambridge, Cambridge, UK.

14 <sup>5</sup>Department of Chemistry, University of Cambridge, Cambridge, UK.

15 <sup>6</sup>Cambridge Breast Cancer Research Unit, Cambridge University Hospital NHS Foundation Trust,  
16 Cambridge, UK.

17 <sup>7</sup>Cancer Research UK Cambridge Centre, NIHR Cambridge Biomedical Research Centre and  
18 Cambridge Experimental Cancer Medicine Centre, Cambridge University Hospital NHS Foundation  
19 Trust, Cambridge, UK.

20 <sup>8</sup>School of Clinical Medicine, University of Cambridge, Cambridge, UK

21  
22 **Present addresses:**

23 Center for Molecular Medicine Cologne, University of Cologne, Germany

24 Robert Hänsel-Hertsch

25  
26 Cancer Research UK Cambridge Institute and Department of Oncology, Li Ka Shing Centre,  
27 University of Cambridge, Cambridge, UK.

28 Angela Simeone, Abigail Shea, Winnie W.I. Hui, Katherine G. Zyner, Oscar M. Rueda, David  
29 Tannahill, Alistair Martin, Carlos Caldas, Shankar Balasubramanian

30  
31 Inivata, Babraham Research Campus, Cambridge, UK.

32 Giovanni Marsico

33  
34 The Institute of Cancer Research, London, UK

35 Alejandra Bruna

36  
37 Department of Chemistry, University of Cambridge, Cambridge, UK.

38 Xiaoyun Zhang, Santosh Adhikari, Shankar Balasubramanian

39  
40 **Corresponding author:**

41 Shankar Balasubramanian

42 e-mail: sb10031@cam.ac.uk

44 **Abstract:**

45 Response and resistance to anticancer therapies vary due to inter- and intra-tumor  
46 heterogeneity<sup>1</sup>. Here, we map differentially enriched G-quadruplex (G4) DNA structure-  
47 forming regions ( $\Delta$ G4Rs) in 22 breast cancer patient-derived tumor xenograft (PDTX) models.  
48  $\Delta$ G4Rs are associated with the promoter of highly amplified and expressed genes, and with  
49 somatic single-nucleotide variants. Specific  $\Delta$ G4Rs reveal 7 transcription factor (TF) programs  
50 across PDTXs.  $\Delta$ G4R abundance and locations stratify PDTXs into at least three G4-based  
51 subtypes.  $\Delta$ G4Rs in most PDTXs (14/22) associated with more than one breast cancer subtype,  
52 which we also call an integrative cluster (IC)<sup>2</sup>. This suggests the frequent coexistence of  
53 multiple breast cancer states within a PDTX model; the majority displaying aggressive triple-  
54 negative IC10. Short-term cultures of PDTX models with increased  $\Delta$ G4R levels are more  
55 sensitive to small molecules targeting G4 DNA. Thus, G4 landscapes reveal additional IC-  
56 related intra-tumor heterogeneity in PDTX biopsies, improving breast cancer stratification and  
57 potentially new treatment strategies.

58

59 **Main:**

60 G-quadruplexes are four-stranded secondary structures that can form in certain G-rich DNA  
61 sequences<sup>3,4</sup>. We previously used in vitro sequencing (G4-seq) to establish where endogenous  
62 G4s could form in the human genome<sup>5</sup>. Qualitative profiling of endogenous G4 DNA in  
63 chromatin by G4-ChIP-seq revealed prominence of their formation in promoters of highly  
64 expressed cancer genes<sup>6-8</sup>. Computational predictions of G4s in eukaryotic genomes have  
65 linked G4 motifs to genomic instability<sup>9-11</sup>, suggesting that G4-selective helicases maintain  
66 genome stability during DNA replication and transcription<sup>3,4</sup>. Supporting this, we have recently  
67 reported the prevalence of endogenous DNA double-strand breakage (DSB) in G4-seq derived  
68 sequences that are found in nucleosome-depleted regions (NDRs) of highly expressed human  
69 cancer genes<sup>12</sup>. Fundamental mechanisms including DNA transcription and replication are  
70 endogenous sources for DSBs and genome instability<sup>13</sup>. Computational predictions of DNA  
71 motifs<sup>14</sup> have suggested that human G4s may be associated with pan-cancer somatic copy  
72 number aberrations (CNAs), which we previously confirmed by G4-seq<sup>15</sup>. CNA landscapes  
73 impact gene expression and shape breast cancer heterogeneity<sup>2</sup>. Our analysis of 2,000 primary  
74 breast cancers previously revealed 11 different subgroups, called integrative clusters (ICs)<sup>16-</sup>  
75 <sup>18</sup>.

76

77 To establish how G4 DNA structures may relate to breast cancer biology, we developed and  
78 applied a quantitative, comparative G4-ChIP-seq (qG4-ChIP-seq) methodology to map G4  
79 DNA structure formation in 22 breast cancer PDTX models that retain their original inter- and  
80 intra-tumor heterogeneity<sup>17</sup>. We adapted the ChIP-Rx approach<sup>19</sup> and employed *Drosophila*  
81 *melanogaster* chromatin as an internal reference to normalize the ChIP-seq data and reduce  
82 technical variability to enhance the characterization of true biological variation (**Fig. 1a** and  
83 **Methods**). Improvement in experimental reproducibility can be evaluated by analyzing the  
84 similarity between four repeated measurements of one PDTX sample vs. four repeated  
85 measurements acquired from a different PDTX sample; either from the same or a different  
86 PDTX model (**Fig. 1b-d, Extended Data Fig. 1a**). Normalization increased the reproducibility  
87 of our human cancer qG4-ChIP-seq data across technical and biological replicates. We derived

88 a coefficient termed improvement factor ( $I_F$ ) whereby  $I_F > 0$  indicates increased  
89 reproducibility, whereas  $I_F < 0$  signifies decreased reproducibility after normalization (**Fig. 1b-**  
90 **d**, see **Methods**). We applied qG4-ChIP-seq to profile the G4 landscape in estrogen receptor-  
91 positive (ER+) or triple-negative (ER-, HER2-, PR-) PDTX models representing most  
92 integrative clusters (IC 1, 5, 8, 9 and 10)<sup>17</sup>. We assessed the reproducibility of qG4-ChIP-seq  
93 by processing different parts of the same tumor on a different day with different reagents, while  
94 keeping the reference chromatin batch constant (**Fig. 1d** and **Supplementary Table 1**).  
95 Overall, across all studied PDTX models, qG4-ChIP-seq identified ~26,000 reproducibly  
96 enriched regions of which 97% comprised a G4 sequence motif (**Extended Data Fig. 1b**).  
97 Comparative qG4-ChIP-seq analysis of 22 PDTX models revealed differentially enriched G4  
98 regions (~700-17,000), hereafter called  $\Delta$ G4Rs, and constant G4 regions (~100), hereafter  
99 called CG4Rs (see **Methods** for detailed description). We found that some  $\Delta$ G4Rs are unique  
100 to a given PDTX (**Fig. 1e, f**) whilst others are common to more than one PDTX model  
101 (**Extended Data Fig. 1c**), suggesting that  $\Delta$ G4R loci may relate to differences in intrinsic  
102 cancer biology.

103

104 To explore whether the  $\Delta$ G4Rs are coupled to the underlying PDTX biology, we performed a  
105 pairwise comparison of the  $\Delta$ G4Rs in all PDTX models and stratified them according to their  
106 similarity (**Fig. 2a**). Without consideration of the annotated PDTX IC or ER status (**Extended**  
107 **Data Fig. 2a**), hierarchical clustering of the  $\Delta$ G4R similarity alone revealed the existence of  
108 three PDTX clusters (**Fig. 2a**). To explore the relationship between CNAs and  $\Delta$ G4Rs or  
109 CG4Rs, we determined CNAs in the PDTX models<sup>20</sup> by comparing the data from sequenced  
110 input libraries with the corresponding qG4-ChIP-seq data for each PDTX sample (see  
111 **Methods**). Examination of highly amplified (AMP), amplified (GAIN), neutral (NEUT),  
112 heterozygous deletions (HETD) and homozygous deletions (HOMD) across all PDTX models  
113 revealed a significant enrichment ( $P < 0.0001$ , **Fig. 2b**) of  $\Delta$ G4Rs, but not CG4Rs, in AMPs  
114 relative to the other CNA categories.  $\Delta$ G4Rs are also more abundant ( $P < 0.0001$ , **Extended**  
115 **Data Fig. 2b**) in amplified regions (AMP + GAIN) in comparison to the other CNA categories.  
116 Notably, the number of observed AMPs does not explain the  $\Delta$ G4R abundance or enrichment  
117 in AMPs since  $\Delta$ G4R and AMP levels vary independently (**Extended Data Fig. 2c**). We also  
118 explored a possible connection of G4 structure with the occurrence of single-nucleotide  
119 variants (SNVs); we previously derived SNVs for some of the PDTXs used here<sup>17</sup>. Notably,  
120  $\Delta$ G4Rs, but not CG4Rs, are significantly ( $P < 0.0001$ ) enriched in SNVs of the PDTXs relative  
121 to random permutation, implying a potential role of G4 formation in the formation of breast  
122 cancer point mutations (**Fig. 2b**). In agreement with our previous observations in cell lines<sup>6-</sup>  
123 <sup>8,21</sup>, G4 structures in the PDTXs are highly enriched in gene promoters, including 5'UTR  
124 regions (**Fig. 2c**). We find that in the tumor material derived from PDTXs,  $\Delta$ G4Rs are also  
125 significantly enriched ( $P < 0.0001$ ) in promoters of highly expressed genes when compared to  
126 medium and lowly expressed ones (**Fig. 2d**, for gene expression classification see **Methods**).  
127 Strikingly, regardless of IC or ER classification, highly expressed promoters show significantly  
128 ( $P < 0.05$ ) greater qG4-ChIP-seq signal in highly amplified (AMP) CNAs relative to other  
129 promoters (**Fig. 2e**). Thus, G4 structures are more prevalent in promoters of highly expressed  
130 and amplified genes in a way that cannot be explained by a single IC and/or ER status. To  
131 explore whether  $\Delta$ G4Rs of a particular PDTX associate with its anticipated IC gene signature,

132 we systematically overlapped promoter regions of the 10 different IC gene sets with the 22  
133 different  $\Delta$ G4Rs. Across all PDTXs,  $\Delta$ G4Rs associate more ( $P < 0.001$ ) with the signature gene  
134 promoter of IC10 than with IC1-9 (**Extended Data Fig. 2d**). These results suggest that the  
135 majority of PDTXs in our cohort display aggressive triple-negative IC10-related breast cancer  
136 gene activity. While the individual  $\Delta$ G4Rs of the 22 PDTXs generally associate with their  
137 anticipated IC status, most (14/22) models display the existence of multiple or distinct IC-  
138 related signature genes (**Fig. 2f** and **Extended Data Fig. 2e**). For example, integrative CNA  
139 and expression profiling of '+/1/HCI005' and '-/10/VHIO179' stratifies these PDTXs as IC1  
140 and IC10, yet their  $\Delta$ G4Rs predominantly associate with at least two different IC-defining  
141 promoter sets that are highly expressed (**Fig. 2f**). This suggests  $\Delta$ G4Rs provide additional  
142 information relative to CNA/expression profiling and revealed the coexistence of multiple  
143 cancer states, thus more intra-tumor heterogeneity with respect to ICs for the majority of PDTX  
144 models (**Extended Data Fig. 2e**). The analysis of 2,000 primary breast tumors revealed 45  
145 common driver regions that are characteristic for CNA-induced gene expression alterations<sup>2</sup>.  
146  $\Delta$ G4Rs, but not CG4Rs, associate with the 45 common breast cancer driver regions (**Fig. 2g**),  
147 highlighting  $\Delta$ G4Rs as a genomic marker of breast cancer driver regions.

148  
149 While pioneering factors such as FOXA1 establish nucleosome-depleted regions (NDRs), TFs  
150 bind to NDRs, thereby mediating transcriptional activity, e.g. through promoter enhancer  
151 interactions<sup>22</sup>. Importantly, TFs can co-target a particular NDR via interactions with other TFs,  
152 thus they can bind DNA independently of their primary consensus binding motif<sup>23</sup>. As  $\Delta$ G4Rs  
153 are prevalent in NDRs<sup>6</sup>, we hypothesized that any  $\Delta$ G4R association with TF binding sites  
154 (TFBS) might reveal TFs that differentially regulate breast cancer development in the PDTX  
155 models. To address this hypothesis, we extracted TF binding sites (TFBS) (see **Methods**) from  
156 breast cancer TF ChIP-seq datasets (ChIP-ATLAS)<sup>24</sup> and computed fold-enrichment over  
157 random in the different  $\Delta$ G4Rs of all the 22 PDTX models (see **Methods**). Hierarchical  
158 clustering of  $\Delta$ G4R fold-enrichments in TFBS revealed increased similarity among some  
159 PDTX models (**Extended Data Fig. 3a, b**), suggesting that some PDTX models share the same  
160 TF activities while others do not. Considering the similarity of  $\Delta$ G4R fold-enrichments in  
161 TFBS across the 22 PDTX models, we identified 7 distinct TF programs that are differentially  
162 active across the PDTXs (**Fig. 3, Extended Data Fig. 3c**). Strikingly,  $\Delta$ G4R – TFBS  
163 enrichments of 4/7 TF programs were significantly higher in either IC10/9, IC8/1, ER-negative  
164 or -positive stratified PDTX models, suggesting that these TF programs are more active in  
165 certain breast cancer subtypes. We found that differential TF expression levels of a TF program  
166 can coincide with the  $\Delta$ G4R fold-enrichments in TFBS of particular PDTXs. For example, the  
167 TLE3-GATA3 TF cluster is significantly more expressed and enriched for  $\Delta$ G4Rs in PDTXs  
168 that are ER-positive or IC 8/1 but not ER-negative or IC10/9 (**Fig. 3, Extended Data Fig. 3c**).  
169 Importantly, all TF programs are expressed (**Extended Data Fig. 3c**), suggesting that  $\Delta$ G4R  
170 fold-enrichments in TFBS may infer differential TF activity in cancer tissues.

171  
172 To characterize pharmacogenomic correlations and enable strategies for precision medicine,  
173 we established short-term cultures of PDTX-derived tumor cells (PDTC). Importantly, PDTCs  
174 preserve the intra-tumor heterogeneity of the PDTX models<sup>17</sup>. Our high-throughput drug  
175 screens, deposited in the Breast Cancer PDTX Encyclopedia (BCaPE), revealed substantial

176 differences in PDTTC response, importantly, even among PDTTCs derived from models stratified  
177 into the same integrative cluster (**Extended Data Fig. 4a**). This suggests the need to consider  
178 additional approaches to decode pharmacogenomic correlations. We previously demonstrated  
179 that human immortalized keratinocytes displayed ~7-fold more G4 regions than normal  
180 keratinocytes by G4-ChIP-seq<sup>6</sup> and exhibited a corresponding increase in sensitivity (~7-fold)  
181 to G4-ligand treatment by pyridostatin (PDS)<sup>25</sup>. This led us to hypothesize that models with  
182 higher  $\Delta$ G4R levels would respond better to G4-ligand treatment, because they are a  
183 quantitative measure of differences in the number of G4 regions. To explore this, we evaluated  
184 G4 ligand-sensitivity in PDTTC derived from models with qG4-ChIP-seq data. We evaluated  
185 two established, yet structurally distinct, small molecules with high G4 DNA selectivity; PDS<sup>25</sup>  
186 and CX-5461<sup>26</sup>. As a negative control, we synthesized an isomer of PDS (i-PDS) that shows  
187 substantially reduced G4 affinity (**Extended Data Fig. 4b, Methods and Supplementary**  
188 **Data 1**). We found that PDTTCs with an increased level of  $\Delta$ G4Rs were significantly ( $P < 0.05$ ,  
189  $r = 0.5-0.8$ ) more sensitive to PDS and CX-5461 but not control i-PDS G4-ligand treatments  
190 (**Fig. 4**). Since PDTTC  $\Delta$ G4Rs are highly enriched in amplified CNAs, we asked whether CNA  
191 amplification level alone was sufficient to predict G4 sensitivity. Notably, we found that  
192 amplified CNA levels lacked a positive correlation with PDTTC responses to all G4-ligands  
193 (**Fig. 4, Extended Data Fig. 4c**). Taken together, these findings highlight the potential of  
194  $\Delta$ G4R mapping as a predictive biomarker for G4-ligand therapy<sup>26</sup>.

195

196 By developing quantitative G4-ChIP-seq, we have obtained G4 DNA maps in chromatin from  
197 patient-derived models, which substantially advances previous qualitative methods using  
198 established cell lines in 2D culture or tissue immunohistochemistry<sup>6,27</sup>. We have generated G4-  
199 DNA maps for 22 PDTTC breast cancer models and revealed how they reflect the underlying  
200 breast cancer biology, such as the relationship with TF occupancy and highly expressed driver  
201 genes in amplified CNAs. Our matched integrative analysis of PDTTC-derived somatic  
202 mutations, CNAs and SNVs, and endogenous G4 DNA landscapes highlight a link between  
203 cancer genome instability and G4 structure formation. Overall, we discovered that G4 DNA  
204 regions are highly associated to critical drivers of triple-negative breast cancer models and/or  
205 IC9-10 relative to ER+/IC1-8 PDTTC models. While strategies are currently under development  
206 to identify cancers that respond to G4 ligand treatment based upon their BRCA1/2 status<sup>26,28,29</sup>,  
207 our results indicate that G4 profiling alone can identify sensitive cancers, which may or may  
208 not be related to their BRCA status. By integrating PDTTC  $\Delta$ G4Rs with established gene  
209 signatures of 10 different breast cancer subtypes (IC), we discovered that the majority of the  
210 22 PDTTC models have G4 patterns that associate with more than one IC, providing an added  
211 layer of intra-tumor heterogeneity. Our interrogation of breast cancer TF ChIP-seq profiles  
212 with the  $\Delta$ G4Rs has highlighted the existence of at least seven distinct TF programs that are  
213 mostly dominant in either ER+/IC1-8 or triple-negative/IC9-10 breast cancers. This supports  
214 that many TFs, instead of a single, defined TF or TF-complex, co-target and -regulate breast  
215 cancer gene activity. Quantitative profiling of G4 structures adds information to conventional  
216 copy number aberration and expression profiling, potentially increasing resolution up to  
217 ~1000-fold (~100-500 bp vs. ~100 kb), hence helping in the identification of specific drivers  
218 within large amplicons. We also provide evidence that  $\Delta$ G4Rs, in combination with established  
219 knowledge on subtypes, can refine the genomic, transcriptomic and regulatory classification of

220 breast cancer. Finally, G4 levels in cancer models are sufficient to predict response to treatment  
221 by small molecules that target G4 DNA structures, highlighting G4s as genomic features with  
222 potential for future diagnostics and therapeutics.

223

## 224 **Acknowledgements**

225 The authors would like to thank the staff in the Genomics and the Compliance and Biobanking  
226 Core Facilities at Cancer Research UK Cambridge Institute. We acknowledge support from  
227 University of Cambridge and Cancer Research UK program. We kindly thank Dr. Ben Czech  
228 of the Hannon laboratory for providing S2 *D. melanogaster* cells. The Caldas and  
229 Balasubramanian laboratories are supported by core funding from Cancer Research UK  
230 (C14303/A17197). The Balasubramanian laboratory is supported by Program grant funding  
231 from Cancer Research UK (C9681/A18618 and C9681/A29214) and a Wellcome Trust  
232 Investigator Award (209441/z/17/z). We acknowledge Dr. Marco Di Antonio for  
233 conceptualizing the design of i-PDS. Prior to the revision of this study work by Dr. Robert  
234 Hänsel-Hertsch was supported by the Balasubramanian group, afterwards additionally  
235 supported by core funding of the Center for Molecular Medicine Cologne (CMMC).

236

## 237 **Author Contributions**

238 R.H.H., C.C. and S.B. conceived this study. R.H.H. developed quantitative G4-ChIP-seq.  
239 R.H.H., A.B., O.M.R. and C.C. designed the PDTX model panel for this study. R.H.H.  
240 processed all the PDTX tissues and prepared the chromatin samples for G4-ChIP-seq. R.H.H.,  
241 W.H. and K.G.Z. performed G4-ChIP-seq. A.M., A.B., O.M.R. and C.C. performed and  
242 interpreted genomic and transcriptomic characterization of all PDTX models. A. Shea  
243 performed the G4-ligand treatment assay, which was analyzed by O.M.R.. R.H.H. and G.M.  
244 implemented a computational pipeline to measure normalization performance, which was  
245 refined by A. Simeone. X.Z. synthesized i-PDS with support of S.A., and performed G4-ligand  
246 *in vitro* experiments and analysis. R.H.H. and A. Simeone performed all the G4-ChIP-seq  
247 related computational analysis. R.H.H., C.C., and S.B. interpreted the results with input from  
248 all authors. R.H.H. prepared the figures. R.H.H., C.C., and S.B. wrote the manuscript with  
249 contributions from all authors.

250

## 251 **Competing interests**

252 S.B. is an advisor and shareholder of Cambridge Epigenetix Ltd.

253

## 254 **References:**

255

- 256 1. Flavahan, W. A., Gaskell, E. & Bernstein, B. E. Epigenetic plasticity and the hallmarks of  
257 cancer. *Science* **357**, eaal2380 (2017).
- 258 2. Curtis, C. *et al.* The genomic and transcriptomic architecture of 2,000 breast tumours reveals  
259 novel subgroups. *Nature* **486**, 346–352 (2012).
- 260 3. Rhodes, D. & Lipps, H. J. G-quadruplexes and their regulatory roles in biology. *Nucleic Acids*  
261 *Res.* **43**, 8627–37 (2015).
- 262 4. Varshney, D., Spiegel, J., Zyner, K., Tannahill, D. & Balasubramanian, S. The regulation and  
263 functions of DNA and RNA G-quadruplexes. *Nat. Rev. Mol. Cell Biol.* (2020).  
264 doi:10.1038/s41580-020-0236-x
- 265 5. Marsico, G. *et al.* Whole genome experimental maps of DNA G-quadruplexes in multiple

- 266 species. *Nucleic Acids Res.* **47**, 3862–3874 (2019).
- 267 6. Hänsel-Hertsch, R. *et al.* G-quadruplex structures mark human regulatory chromatin. *Nat.*  
268 *Genet.* **48**, 1267–72 (2016).
- 269 7. Hänsel-Hertsch, R., Spiegel, J., Marsico, G., Tannahill, D. & Balasubramanian, S. Genome-  
270 wide mapping of endogenous G-quadruplex DNA structures by chromatin  
271 immunoprecipitation and high-throughput sequencing. *Nat. Protoc.* **13**, 551–564 (2018).
- 272 8. Hänsel-Hertsch, R., Antonio, M. Di & Balasubramanian, S. DNA G-quadruplexes in the  
273 human genome: detection, functions and therapeutic potential. *Nat. Rev. Mol. Cell Biol.* **18**,  
274 279–284 (2017).
- 275 9. Paeschke, K., Capra, J. A. & Zakian, V. A. DNA replication through G-quadruplex motifs is  
276 promoted by the *Saccharomyces cerevisiae* Pif1 DNA helicase. *Cell* **145**, 678–91 (2011).
- 277 10. Cheung, I., Schertzer, M., Rose, A. & Lansdorp, P. M. Disruption of dog-1 in *Caenorhabditis*  
278 *elegans* triggers deletions upstream of guanine-rich DNA. *Nat. Genet.* **31**, 405–409 (2002).
- 279 11. Georgakopoulos-Soares, I., Morganella, S., Jain, N., Hemberg, M. & Nik-Zainal, S.  
280 Noncanonical secondary structures arising from non-B DNA motifs are determinants of  
281 mutagenesis. *Genome Res.* **28**, 1264–1271 (2018).
- 282 12. Lensing, S. V. *et al.* DSBCapture: In situ capture and sequencing of DNA breaks. *Nat.*  
283 *Methods* **13**, (2016).
- 284 13. Bouwman, B. A. M. & Crosetto, N. Endogenous DNA Double-Strand Breaks during DNA  
285 Transactions: Emerging Insights and Methods for Genome-Wide Profiling. *Genes (Basel)*. **9**,  
286 (2018).
- 287 14. De, S. & Michor, F. DNA secondary structures and epigenetic determinants of cancer genome  
288 evolution. *Nat. Struct. Mol. Biol.* **18**, 950–955 (2011).
- 289 15. Chambers, V. S. *et al.* High-throughput sequencing of DNA G-quadruplex structures in the  
290 human genome. *Nat. Biotechnol.* **33**, 1–7 (2015).
- 291 16. Pereira, B. *et al.* The somatic mutation profiles of 2,433 breast cancers refine their genomic  
292 and transcriptomic landscapes. *Nat. Commun.* **7**, 11479 (2016).
- 293 17. Bruna, A. *et al.* A Biobank of Breast Cancer Explants with Preserved Intra-tumor  
294 Heterogeneity to Screen Anticancer Compounds. *Cell* **167**, 260-274.e22 (2016).
- 295 18. Rueda, O. M. *et al.* Dynamics of breast-cancer relapse reveal late-recurring ER-positive  
296 genomic subgroups. *Nature* **567**, 399–404 (2019).
- 297 19. Orlando, D. A. *et al.* Quantitative ChIP-Seq normalization reveals global modulation of the  
298 epigenome. *Cell Rep.* **9**, 1163–1170 (2014).
- 299 20. Scheinin, I. *et al.* DNA copy number analysis of fresh and formalin-fixed specimens by  
300 shallow whole-genome sequencing with identification and exclusion of problematic regions in  
301 the genome assembly. *Genome Res.* **24**, 2022–32 (2014).
- 302 21. Mao, S.-Q. *et al.* DNA G-quadruplex structures mold the DNA methylome. *Nat. Struct. Mol.*  
303 *Biol.* **25**, 951–957 (2018).
- 304 22. Zaret, K. S. & Carroll, J. S. Pioneer transcription factors: establishing competence for gene  
305 expression. *Genes Dev.* **25**, 2227–2241 (2011).
- 306 23. Gertz, J. *et al.* Distinct Properties of Cell-Type-Specific and Shared Transcription Factor  
307 Binding Sites. *Mol. Cell* **52**, 25–36 (2013).
- 308 24. Oki, S. *et al.* ChIP-Atlas: a data-mining suite powered by full integration of public ChIP-seq  
309 data. *EMBO Rep.* **19**, e46255 (2018).
- 310 25. Rodriguez, R. *et al.* A Novel Small Molecule That Alters Shelterin Integrity and Triggers a  
311 DNA-Damage Response at Telomeres. *J. Am. Chem. Soc.* **130**, 15758–15759 (2008).
- 312 26. Xu, H. *et al.* CX-5461 is a DNA G-quadruplex stabilizer with selective lethality in BRCA1/2  
313 deficient tumours. *Nat. Commun.* **8**, 14432 (2017).
- 314 27. Biffi, G., Tannahill, D., Miller, J., Howat, W. J. & Balasubramanian, S. Elevated Levels of G-  
315 Quadruplex Formation in Human Stomach and Liver Cancer Tissues. *PLoS One* **9**, e102711  
316 (2014).
- 317 28. McLuckie, K. I. E. *et al.* G-Quadruplex DNA as a Molecular Target for Induced Synthetic  
318 Lethality in Cancer Cells. *J. Am. Chem. Soc.* **135**, 9640–9643 (2013).
- 319 29. Zimmer, J. *et al.* Targeting BRCA1 and BRCA2 Deficiencies with G-Quadruplex-Interacting  
320 Compounds. *Mol. Cell* **61**, 449–460 (2016).

321  
322 **Figure 1 | Quantitative G4-ChIP-seq of PDTX reveals differentially enriched G4 DNA**  
323 **regions. a**, Quantitative G4-ChIP-seq (qG4-ChIP-seq), exemplified with two different PDTX  
324 models (brackets: estrogen receptor status/integrative cluster status/PDTX model name). Blue  
325 and red indicates chromatin isolated from two different PDTX models, which is combined with  
326 *D. melanogaster* (reference) chromatin (black). **b**, Normalization strategy: technical and  
327 biological replicates of the same condition (nodes of the same color) get closer in space after  
328 normalization and samples of different PDTX models become more separated in space. **c**,  
329 Estimation of normalization factors using reference read coverage. To derive normalization  
330 factors, either all (Total recovery) reference sequencing reads are considered or only the ones  
331 in a predefined set of enriched regions (G4 region). Subsequent rescaling of the cumulative  
332 human cancer signal by the normalization factors is done across all experiments (see **Methods**  
333 section). **d**, Barplot of the improvement factors ( $I_F$ ) quantifying normalization performance for  
334 all 22 PDTX models considering the reads in the enriched (G4 regions) and all recovered reads  
335 (Total recovery). Improvement factor evaluates the level of increased similarity (positive  
336 values) between technical replicates (black) and biological replicates (grey) (see **Methods**). **e**,  
337 Heatmap of human cancer normalized (reference normalized and input subtracted counts per  
338 million) qG4-ChIP-seq data for  $\Delta$ G4Rs of two PDTX models -/10/AB863M (red) and  
339 +/8/STG143 (blue). **f**, Example genome browser views showing  $\Delta$ G4Rs and normalized qG4-  
340 ChIP-seq track intensities of four technical replicates qG4-ChIP-seq for two PDTX models  
341 (red -/10/AB863M, blue +/8/STG143). PDTX annotation: estrogen receptor status/integrative  
342 cluster status/PDTX model name.

343  
344 **Figure 2 | G4 DNA prevalence in the genomic and transcriptomic architecture of PDTX**  
345 **breast cancer models. a**, Heatmap visualizing similarity of  $\Delta$ G4Rs from 22 different PDTX  
346 models. Hierarchical clustering is shown (Euclidean distance, ward.d2); color intensity and the  
347 size of the circle are proportional to the correlation coefficients. **b**, Left: Distribution of PDTX  
348 median fold-enrichments for  $\Delta$ G4Rs and CG4Rs (common or unchanged qG4-ChIP-seq  
349 regions) in copy number aberrations (CNAs) relative to random permutation ( $n = 10$   
350 permutations for each of the 23 independent  $\Delta$ G4R and CG4R maps); AMP = highly amplified  
351 regions, GAIN = amplified, NEUT = unchanged or neutral regions, HETD = heterozygous  
352 deletions, HOMD = homozygous deletions. Right: Distribution of  $\Delta$ G4R and CG4R  
353 enrichments for single-nucleotide variants (SNVs) within the PDTX samples relative to  
354 random permutation ( $n = 10$  permutations for each of the 16 independent  $\Delta$ G4R and CG4R  
355 maps). Significances were calculated using a *t* test (Mann-Whitney) \*\*\*\*  $P < 0.0001$  (exact,  
356 two-tailed). **c**, Genome annotation of  $\sim 26,000$  PDTX G4-ChIP-seq regions. Black bars:  
357 proportion of G4 regions in particular genomic annotation, red bars: fold-enrichment over  
358 random ( $n = 5$  permutations) genomic regions. Data are presented as mean values  $\pm$  SD. **d**, Y-  
359 axis: high (blue), medium (red) and low (black) expressed genes. X-axis: For each PDTX,  
360 percentage of  $\Delta$ G4Rs in the expressed promoters. Brackets indicate significant differences of  
361  $\Delta$ G4Rs spanning promoters of different expression levels.  $N = 20$  PDTX  $\Delta$ G4Rs were  
362 associated with  $n = 20$  PDTX promoter expression levels. Significant differences were  
363 calculated using a Tukey multiple comparison test \*\*\*\*  $P < 0.0001$  (adjusted *P* value). **e**,  
364 Distribution of the integral of  $\Delta$ G4R signal intensities (median of cpm) in high, medium or low



365 expressed gene promoters ( $\pm 1$  kb TSS) that are in AMP, GAIN, NEUT or HETD regions. N  
366 = 22 PDTX  $\Delta$ G4R ChIP intensities measured in n = 22 PDTX promoter expression levels in  
367 the different CNA categories. Significances were calculated using a *t* test (Mann-Whitney) \*\*  
368  $P < 0.01$ , \*  $P < 0.05$  ( $P$  values are exact, two-tailed). **f**, Scatter plots of four individual PDTX.  
369 Y-axis: Overlap of gene promoters (%) for distinct gene signatures of the 10 different  
370 integrative clusters<sup>2</sup> with  $\Delta$ G4Rs. X-axis: the significance of the overlap relative to chance  
371 (Fisher). The expected IC classification for each PDTX model is highlighted in red. **g**, Fold-  
372 enrichment over random (bar plot) of  $\Delta$ G4Rs or CG4Rs in 45 common breast cancer driver  
373 regions<sup>2</sup> relative to chance. Significance (color code): empirical  $P$  value (exact, two-tailed)  
374 obtained with 1,000 randomizations. Red dashed line indicates threshold of fold-enrichment  
375 over random. Box plot elements: center line, median; box limits, lower and upper quartiles;  
376 whiskers, lowest and highest value.

377

### 378 **Figure 3 | G4 DNA regions reveal the activity of distinct transcription factor programs.**

379 Transcription factor correlation matrix heatmap ( $134 \times 134$ ). obtained by starting from the  
380 matrix of fold-enrichments over random of 22  $\Delta$ G4Rs at 134 breast cancer ChIP-seq  
381 transcription factor binding regions (TFBS) from ChIP-ATLAS database. Hierarchical  
382 clustering (ward.d2) of the correlations identifies TF sub-groups with similar correlation values  
383 across the 22 PDTX models; TF subgroups are highlighted by dashed line. Color intensity and  
384 the size of the circle are proportional to the correlation coefficients. Name of each TF subgroup  
385 relates to first and last TF within each subgroup. For each subgroup, there are boxplots of  
386  $\Delta$ G4R/TFBS fold-enrichments (blue) and of TF expression levels (TPM, red) stratified by  
387 various classifications of the PDTX models (ER+, ER-, membership to IC8/1, membership to  
388 IC10/9). N = 22 PDTX models were used to derive 22  $\Delta$ G4R maps and ChIP-ATLAS fold-  
389 enrichment values. Significances illustrated in box plots were calculated using the Mann-  
390 Whitney test \*\*\*\*  $P < 0.0001$ , \*\*  $P < 0.01$ , \*  $P < 0.05$  (exact  $P$  values, two-tailed). Box plot  
391 elements: center line, median; box limits, lower and upper quartiles; whiskers, lowest and  
392 highest value.

393

394 **Figure 4 | G4 DNA levels predict response to G4-ligands.** Scatterplots of  $\Delta$ G4Rs (left-top,  
395 left-bottom, right-top) or highly amplified regions (AMP, right-bottom) levels (x-axis) against  
396 PDTC response (Area under the curve; AUC, y-axis) of PDTC models to G4-ligands with  
397 enhanced (PDS, CX-5461) and reduced (i-PDS) G4 affinities, see also **Methods**. Error bars  
398 reflect mean, upper and lower limit AUCs. N = 9 PDTC samples. Additionally, N = 3 PDTC  
399 samples were independently investigated. Spearman correlation ( $r$ ) and significance (exact  
400 two-tailed  $P$  value for nonparametric correlation) are shown.

401

### 402 **Methods:**

#### 403 ***Quantitative G4-ChIP-seq (qG4-ChIP-seq).***

404 G4-ChIP-seq was performed as previously described<sup>7</sup> with the following adaptations for PDTX  
405 tissue. Briefly, *D. melanogaster* S2 cells were cultured in Schneider's *Drosophila* Medium  
406 (Thermo Fisher Scientific, cat no. R69007) containing 10% fetal bovine serum (FBS) Medium  
407 (Thermo Fisher Scientific, cat no. 10500064). To prepare spike-in *Drosophila* chromatin, 100  
408 million cells were i) harvested by centrifugation, ii) fixed for 10 min in a solution of media

409 containing 10 % FBS, 1% formaldehyde (Thermo Fisher Scientific, cat no. 28908) and iii)  
410 quenched for 5 min by addition of 125 mM glycine (Fisher Scientific, cat no. 11545005). The  
411 cell pellet was washed with 10 ml PBS, pelleted by centrifugation and subsequently stored on  
412 ice for the lysis procedure. The 2-step chromatin lysis procedure was performed according to  
413 the Chromatrap procedure ("Spin column ChIP kit for qPCR v6.4"). 500  $\mu$ l intact chromatin  
414 was sonicated into 100-500 bp fragments using a Bioruptor Plus (Diagenode cat. no.  
415 B02010003 with cooling) at 4°C. Sonicated chromatin was diluted with 1.5 ml lysis buffer  
416 (Chromatrap cat no. 100005) before snap-freezing as 25  $\mu$ l aliquots. PDTX chromatin was  
417 prepared essentially as described in Schmidt *et al.* (Methods 2009)<sup>30</sup>. Briefly, a snap-frozen  
418 PDTX biopsy,  $\sim$ 1 cm<sup>3</sup>, was transferred into a 50 ml falcon tube, on dry-ice, and crushed into  
419 smaller chunks on dry-ice using a scalpel followed by fixation for 20 min in 30 ml solution A,  
420 containing 1% formaldehyde, and then quenched for 5 min by adding 125 mM glycine. The  
421 supernatant of the pelleted tissue was discarded, and the pellet washed twice with 10 ml ice-  
422 cold PBS before resuspending in 1 ml PBS and transferred to a 1 ml glass Douncer (Fisher  
423 scientific, cat no. 11591295). 10 strokes were employed for each douncing step with a loose  
424 and then tight pestle, and the remaining tissue slurry was transferred to a 15 ml tube,  
425 centrifuged for 5 min at 2,500 $\times$ g and subjected to lysis according to Schmidt et al.<sup>30</sup>. Briefly,  
426 after the 10 ml LB2 treatment and nuclei pelleting step, the pellet was resuspended in 500  $\mu$ l  
427 LB3 and LB3-chromatin solution split into two Bioruptor TBX (Diagenode, cat no.  
428 C30010010-300) sonication tubes. The samples were sonicated until the desired fragment  
429 length (100-500 bp) was achieved. Finally, 50  $\mu$ l of a 10% Triton X-100 LB3 solution was  
430 mixed with the sonicated solution and aliquoted into 50  $\mu$ l aliquots before snap-freezing in  
431 liquid nitrogen. 5  $\mu$ l of PDTX chromatin was quantified by Qubit using the "broad range kit"  
432 (Thermo Fisher Scientific, cat. no. Q32853). In each qG4-ChIP-seq reaction, 225 ng of PDTX  
433 chromatin, 102 ng of spike-in *Drosophila* chromatin and 2% RNaseA (Invitrogen, cat. no.  
434 AM2271) in blocking buffer (25 mM HEPES, pH 7.5, 10.5 mM NaCl, 110 mM KCl, 1 mM  
435 MgCl<sub>2</sub> and 1% BSA (Merck, cat. no. A7030) in Milli-Q water were mixed and incubated at  
436 37°C for 30 min at 800 rpm. All PDTX chromatin batches containing a different concentration  
437 than 30 ng/ $\mu$ l were balanced to the same level, either by dilution with LB3 containing 1%  
438 Triton X-100 or by up-scaling the ChIP reaction. For PDTX chromatin with a concentration of  
439 30 ng/ $\mu$ l, 7.5  $\mu$ l of the PDTX chromatin was added to a solution containing 270  $\mu$ l blocking  
440 buffer including 2% RNase A and 7.5  $\mu$ l spike-in *Drosophila* chromatin. After RNaseA  
441 treatment, 15  $\mu$ l of 2  $\mu$ M BG4, prepared as described previously<sup>31</sup>, was added to each qG4-  
442 ChIP-seq reaction and the reaction mixture shaken at 1,400 rpm at 16°C for 1 hour. Meanwhile,  
443 65  $\mu$ l of anti-FLAG magnetic beads (Sigma-Aldrich, cat. no. M8823) were washed three times  
444 with 650  $\mu$ l of blocking buffer and resuspended in 1,300  $\mu$ l blocking buffer. The pre-washed  
445 beads were incubated at 16°C at 1,400 rpm and 300  $\mu$ l of pre-washed beads added to the  
446 reaction mixture after BG4 incubation. The reaction mixture with beads was incubated at 16°C  
447 for 1 hour at 1,400 rpm. Then, the beads were washed four times in 400  $\mu$ l cold wash buffer  
448 (100 mM KCl, 0.1% Tween 20 and 10 mM Tris, pH 7.4 in Milli-Q water) in the cold room and  
449 twice at 37°C for 15 min at 1,400 rpm, followed by one cold wash on magnetic stand. The  
450 enriched chromatin on beads was resuspended in 75  $\mu$ l TE buffer and 1  $\mu$ l Proteinase K  
451 (Invitrogen, cat. no. AM2546) added. 6  $\mu$ l Proteinase K was added to input sample which refers

452 to a qG4-ChIP-seq reaction mixture without BG4 and beads. The reaction mixture was  
453 incubated at 65°C for 3 hours at 1,400 rpm and purified using QIAGEN MinElute Kit  
454 (QIAGEN, cat. no. 28206).

455

456 **Library preparation and sequencing.** For 40 µl library preparation reaction, 3-5 ng of the ChIP  
457 or input DNA (Qubit high sensitivity kit, Thermo Fisher Scientific, cat. no. Q32854), 20 µl 2×  
458 tagmentation buffer (Illumina, cat. no. 15027866), 1.25 µl Tn5 enzyme (Illumina, cat. no.  
459 18027865) and nuclease-free water was incubated at 37°C for 20 min at 800 rpm. The reaction  
460 mixture was purified using QIAGEN MinElute Kit (QIAGEN, cat. no. 28206) according to the  
461 manufacturer's instruction and eluted in 20 µl EB buffer. To amplify the library, 20 µl of the  
462 DNA was then mixed with 25 µl NEB Next High Fidelity 2× PCR Master Mix (New England  
463 Biolabs, cat. no. N0541S), 2.5 µl Nextera index kit i5 primer (Illumina, cat. no. 15055290) and  
464 2.5 µl Nextera index kit i7 primer (Illumina, cat. no. 15055290). The PCR program was as  
465 follows: 72°C for 5 minutes, 98°C for 30 seconds, followed by 8 cycles of 98°C for 10 seconds,  
466 63°C for 30 seconds and 72°C for 1 minute. Libraries were quantified using a Bioanalyzer  
467 (Agilent) to estimate the average library size and concentration determined via Qubit HS. The  
468 library concentration was corrected for the library size using the following relationship: 1 ng/µl  
469 = 3nM = 500 bp. Samples were subjected to single-end sequencing with a read length of 75 bp  
470 on an Illumina NextSeq instrument.

471

472 **Mapping, peak calling and peak processing.** Fastq files were trimmed from adapters using  
473 *cutadapt* (options: -q 20 -O 3 <http://dx.doi.org/10.14806/ej.17.1.200>, ver: 1.16) and aligned<sup>32</sup>  
474 to a combined genome consisting of hg19 (*Homo sapiens*), dm6 (*D. melanogaster*) and mm10  
475 (*Mus musculus*) with *bwa-mem* (ver. 0.7.17-r1188). Bam files were generated from the  
476 alignment with *samtools view* (options: -Sb -F2308 -q 10, ver: 1.8) and subsequently split by  
477 organisms to obtain 3 bam files for each sample. Duplicated reads were marked and removed  
478 using *picard MarkDuplicates* (ver: 2.20.3). For all organisms, the total sequencing coverage  
479 (total recovery) was quantified as the total number of unique reads aligning to the respective  
480 genome. Standard peak calling was performed for each sample using *MACS2* (ver. 2.1.2) with  
481 default options. For each human PDTX model, peak regions were considered positive if  
482 confirmed in 2 out of 4 technical replicates (multi2) with *bedtools v2.27.1 multiinter* (see  
483 **Supplementary Table 2**). All human confirmed G4-ChIP-seq peak files (multi2) of the 22  
484 models were merged (*bedtools merge*) and regions more than 99 bp long retained to generate  
485 a single G4 DNA consensus of 26,103 G4 regions. Finally, the coverage of the samples was  
486 quantified using a consensus human set (*bedtools coverage*).

487

488 **Reference normalization factor estimation and human ChIP signal normalization.** For each  
489 PDTX biopsy, four technical qG4-ChIP experiments were performed and sequenced alongside  
490 one input chromatin (control), see also Life Sciences Reporting Summary. In each experiment,  
491 a similar amount of reference (*D. melanogaster*) chromatin from the same batch was added.  
492 To estimate PDTX normalization factors, reference coverage was determined at a pre-defined  
493 consensus consisting of 1,367 intervals (see **Supplementary Data 2**). The reference consensus  
494 set was defined from, and covers, G4-enriched regions observed in more than 110 pull-down

495 experiments. The normalization factor of each ChIP sample has been defined as the ratio  
 496 between the maximum observed coverage (across all ChIP samples) and the individual sample  
 497 coverage. Note that only ChIP experiments were used for this step (i.e. inputs are excluded and  
 498 forced to 1). In turn, the outcome of the normalization approaches were tested using either the  
 499 total recovery or the recovery at the G4 reference consensus regions. The normalization factors  
 500 were then exported and used as input for a customized R script performing the normalization  
 501 of the human signal. For each G4-ChIP-seq experiment, human signal (i.e. read coverage  
 502 within human G4 consensus) was quantified by performing input subtraction and normalization  
 503 with their respective reference reads and human library sizes. To assess if the normalization  
 504 step has globally improved the experimental reproducibility, a quantitative parameter, the  
 505 Improvement Factor  $I_F$ , was devised that measures both the increase (i.e. improvement) in data  
 506 similarity between experiments corresponding to the same technical and biological samples  
 507 and the increase dissimilarity between different samples. Specifically, the improvement factor  
 508 of each biological sample has been estimated as:  
 509

$$I_F = \sum_i^N D_{ratio_{inter}} - \sum_i^N D_{ratio_{intra}}$$

$$D_{ratio} = \frac{Eucl\_dist\_after\_dm\_norm}{Eucl\_dist\_before\_dm\_norm}$$

Where:

- N: ChIP samples
- $Eucl\_dist\_after\_dm\_norm$ : Euclidean similarity matrix computed on input subtracted, library size adjusted, drosophila normalized data and rescaled to its maximum value;
- $Eucl\_dist\_before\_dm\_norm$ : Euclidean similarity matrix computed on input subtracted, library size adjusted data and rescaled to its maximum value;
- $D_{ratio_{inter}}$ : similarity values among samples belonging to the same technical or biological group;
- $D_{ratio_{intra}}$ : similarity values among samples not belonging to the same technical or biological group;

135 individual samples (ChIP + Input) were processed from 22 different PDTX models. Some PDTX models have more than one biological sample (**Supplementary Table 1**).

**Guidelines to normalize G4-ChIP-seq data.** During the optimization of the normalization procedure, we identified some general empirical criteria that can guide in deciding whether the reference (*D. melanogaster*) G4-ChIP-seq data can be used to normalize the human G4-ChIP-seq data, and whether it reduces technical noise and therefore has a beneficial outcome for the reproducibility of the replicated experiments.

1. Sequencing depth of the reference data per G4-ChIP-seq library should be around 5 M reads (after alignment and duplicate removal).

- 535 2. The number of detected peaks in the reference data of the 4 G4-ChIP-seq replicates  
536 (merge of confirmed peaks) should be in the range of several hundreds - 1,000. If no  
537 peaks are detected, it is not reliable to use the reference signal for normalization  
538 purposes.
- 539 3. The fraction of reads in the consensus reference peaks should exceed 0.5% of the total,  
540 ideally 1%. Consensus reference peaks are high-confidence regions that were  
541 consistently detected across many experiments (> 100) and are provided with this study.
- 542 4. Each G4-ChIP-seq library must have a fraction of reads at the consensus reference  
543 peaks at least 2x greater than the respective input library.
- 544 5. Technical and biological IF (average) should be positive, which indicates that  
545 reference normalization has improved the experimental reproducibility of the human  
546 G4-ChIP-seq replicates.

547  
548 ***ΔG4Rs and CG4Rs.*** After the normalization step, differential G4-binding analysis was  
549 employed to identify differentially enriched G4 regions (ΔG4Rs), as described<sup>6,7</sup>. Both  
550 normalization and differential analysis are integrated into our workflow (see  
551 <https://github.com/sblab-bioinformatics/qG4-ChIP-seq-of-breast-cancer-PDTEX/>).

552 Differential G4-binding was carried out with edgeR<sup>33</sup>. Initially, library size and *Drosophila*-  
553 normalized (human) read coverage within human G4 consensus regions were computed. Then,  
554 a generalized linear model with default parameters (negative binomial log-linear distribution  
555 of read counts) were used to assess regions with differential binding signal.

556 Specifically, the differential binding analysis compared each PDTEX model to all the others.  
557 For each comparison, differential DNA G4 regions ΔG4R (i.e., regions specifically present in  
558 a given PDTEX model) were defined as those satisfying the following criteria:  $\log_2(\text{CPM}) \geq 0.6$   
559 and  $\text{FDR} < 0.05$ . Constant G4 regions CG4Rs were defined as those that did not show any  
560 significant differential binding in any model i.e. regions that did not pass the filter in any of the  
561 individual comparisons.

562  
563 ***PDTEX gene expression data.*** Gene expression profiling of the individual PDTEX models, except  
564 for STG316 for which part of the primary, patient-derived, tumor was used, was acquired via  
565 RNA-seq. For the AB863M model, a PDTEX and a primary tumor sample were separately  
566 processed to generate RNA-seq data. Normalized TPMs have been quantified as explained in  
567 Georgopoulou & Callari et al., in preparation (EGA accession: EGAS00001001913) in all  
568 PDTEX models except for PAR1006, PAR1022. For each model, expressed genes were  
569 stratified in 3 groups: high-, medium- and low-expression if they were belonging to the top,  
570 middle or bottom expression tertile, respectively.

571  
572 ***ΔG4R PDTEX stratification.*** Similarity across all the 22 PDTEX ΔG4Rs was estimated using  
573 *bedtools jaccard*. Jaccard indexes of all pairwise comparisons resulted into a  $22 \times 22$  matrix  
574 that have been explored with the Shiny-based web application  
575 <https://asntech.shinyapps.io/intervene><sup>34</sup>. After loading the data, a pairwise intersection  
576 heatmap has been generated with the following settings: plot type: corrplot; correlation  
577 coefficient: Spearman; Agglomerative method: ward.d2; N. of cluster:3; distance matrix  
578 computation: Euclidean.

579

580 **Somatic copy number aberration regions identification.** Copy number segmentation was  
581 performed using the R package QDNAseq<sup>20,35</sup> on input (genomic background of qG4-ChIP-  
582 seq) BAM files sub-set to 5 million reads. A customized R script binned the genome into 100-  
583 kb windows, extracted the read-counts (*binReadCounts*), applied the QDNAseq filters,  
584 calculated (*estimateCorrection*), applied GC correction (*correctBins*), and then normalized and  
585 smoothed outliers. Finally, the copy-number profile of each PDTX model was segmented and  
586 exported. The copy number alterations regions were classified according to the following  
587 filtering criteria:

- 588 ○ highly amplified regions AMP:  $\log_2(\text{fold ratio}) > 0.75$ ;
- 589 ○ amplified regions GAIN :  $0.25 < \log_2(\text{fold ratio}) \leq 0.75$ ;
- 590 ○ neutral regions NEUT :  $-0.3 < \log_2(\text{fold ratio}) \leq 0.25$
- 591 ○ heterozygous deletions HETD:  $-1.4 < \log_2(\text{fold ratio}) \leq -0.3$
- 592 ○ homozygous deletions HOMD:  $\log_2(\text{fold ratio}) \leq -1.4$ .

593

594  **$\Delta G4R$  and  $CG4R$  enrichment in CNA regions relative to random.** The fold-enrichment of  
595  $\Delta G4R$  and  $CG4R$  was empirically estimated over randomly permuted genomic regions. First  
596  $\Delta G4Rs$  and  $CG4Rs$  were 10 times randomly shuffled across the genome (*bedtools shuffle*);  
597 then the number of PDTX  $\Delta G4Rs$  and  $CG4Rs$  overlapping each of the CNA type was counted  
598 in the actual case and in the randomized case. For all CNA types in each PDTX model, the  
599 fold-enrichments were estimated as the ratio of the actual case over each of the ten random  
600 cases, see **Supplementary Table 3**. The distribution of all PDTX models' median fold-  
601 enrichments were then visualized in all individual CNA regions as a combined boxplot (**Fig.**  
602 **2b**).

603

604  **$\Delta G4R$  and  $CG4R$  enrichment for single nucleotide variants SNV relative to random.** As in  
605 the case of CNA, the fold-enrichment of SNV at  $\Delta G4R$  and  $CG4R$  was empirically estimated.  
606 After 5 random shufflings of the  $\Delta G4R$  across the genome, we computed the fold-enrichment  
607 as the actual number of overlaps of G4 regions with SNVs over the average random case (i.e.  
608 average of number of overlaps obtained in each randomization) (see **Supplementary Table**  
609 **4**). For all PDTX models, the analysis was conducted by comparing the model specific  $\Delta G4R$   
610 and SNVs maps.  $CG4Rs$  were compared to all PDTX SNV individually.

611

612 **Genomic and G4-motif annotation and enrichment analysis of PDTX qG4-ChIP-seq peaks.**  
613 PAVIS<sup>36</sup> was used to annotate 26,103 PDTX G4 human consensus regions. Fold-enrichment  
614 analysis was performed as described<sup>6</sup>. The consensus peaks were randomly shuffled across the  
615 genome 5 times. Fold-enrichments were computed as the ratio between the fraction of overlaps  
616 with each genomic feature in the actual case versus the corresponding average random  
617 fractions. G4 motifs were predicted and the presence in the PDTX G4 human consensus regions  
618 measured as previously reported<sup>6</sup>.

619

620 **Promoter -  $\Delta G4R$  - gene expression.** Promoter transcription start site (TSS) coordinates, 1 kb  
621 ( $\pm$ ) from TSS, were generated for 22,483 genes using hg19<sup>6</sup>

622 [https://www.genecodegenes.org/human/release\\_19.html](https://www.genecodegenes.org/human/release_19.html). The fraction of  $\Delta$ G4Rs overlapping  
623 high-, medium-, low-expression gene promoters was estimated. Significance was tested using  
624 the Tukey multiple comparisons test (GaphPad Prism7).

625

626 **Promoter - G4 intensity - Gene expression - CNA analysis.** For this, the human G4 drosophila-  
627 normalized intensity at  $\Delta$ G4R overlapping promoters was considered. The distribution of this  
628 signal was visualized after stratifying promoter by CNA alteration (promoters overlapping to:  
629 AMP, GAIN, NEUT, HETD) and gene-expression groups (promoters belonging to: High-,  
630 Med.- Low- expression group).

631

632  **$\Delta$ G4R and CG4R - Association to upregulated genes from integrative cluster signature IC.**  
633 Promoter coordinates of differentially upregulated genes (Adjusted  $P$  value  $< 0.05$ ;  $\log_2(\text{fold-}$   
634  $\text{change}) > 0.6$ ) of each integrative cluster IC1-10<sup>2</sup> were extracted. For each PDTX model the  
635 association of  $\Delta$ G4Rs and CG4Rs to upregulated promoters was quantified by computing the  
636 corresponding  $P$  value ( $-\log_{10} P$  value) from the fisher test (intervene pairwise option fisher).  
637 For high significant associations, resulting in  $P$  values of 0,  $-\log_{10} P$  value was set to 300. In  
638 addition, the fraction of the IC promoters having a  $\Delta$ G4R overlapping was estimated (intervene  
639 pairwise option fraction). Fractions were transformed into percentage overlap and visualized  
640 together with  $P$  values as scatter plots (**Fig. 2f, Extended Data Fig. 2e**).

641

642  **$\Delta$ G4R and CG4R – Association to 45 common driver regions.** The genomic coordinates of 45  
643 common breast cancer driver regions were taken from Curtis et al.<sup>2</sup> and lifted to hg19 (UCSC  
644 *liftover* tool) for assessment of the fold-enrichments of each PDTX  $\Delta$ G4Rs and CG4Rs at those  
645 genomic locations using the Genomic Association Tester (GAT,  
646 <https://gat.readthedocs.io/en/latest/contents.html>).

647

648 **Transcription factor binding site (TFBS) -  $\Delta$ G4R enrichment analysis.** The genomic fold-  
649 enrichment of each  $\Delta$ G4R over transcription factors binding profiles from breast cancer, and  
650 breast immortalized cells, was determined using the ChIP-ATLAS enrichment web tool  
651 ([https://chip-atlas.org/enrichment\\_analysis](https://chip-atlas.org/enrichment_analysis)), with the following parameters: Antigen class: TFs  
652 and others; Cell type Class: Breast; Threshold for Significance: 500; Select your data:  
653 individual  $\Delta$ G4R in bed format; Select permutation to be compared: 100 random permutations).  
654 22 result tables were obtained with each containing 12 tab-separated columns from which the  
655 following parameters selected: #2 Antigen name, #9 LogPvalue; #11 FoldEnrichment (FE).  
656 Rows were excluded where FE was “Inf”. Selected enrichments with LogPvalue  $< -3$  for each  
657 “antigen name” were averaged by their fold enrichments to give a table with 2 columns:  
658 “antigen name” and “relative averaged fold-enrichment”. The 22 tables were then entered into  
659 a FE matrix with 134 TFs on the rows (where a TF has a FE value in at least 1/22 cases) and  
660 22 columns representing FE in each of the 22 PDTX models. Next, we computed (a) the  
661 Spearman correlation of FE matrix 134 x 134 to assess the similarity between TF FE and (b)  
662 the Spearman correlation on the transposed FE matrix (22  $\times$  22) to assess the similarity between  
663 PDTXs. The correlation matrix (a) for TF was additionally analyzed via hierarchical clustering  
664 (ward.d2). Seven subgroups of TFs were identified. For each subgroup of TF, boxplots were

665 generated for fold-enrichments and TF expression levels stratified by IC classification and/or  
666 ER status of the PDTX models (ER+, ER-, membership to IC8/, membership to IC10/9).

667

668 ***PDTX prepared into cell suspension (PDTC) and high throughput G4-ligand screen.*** PDTCs  
669 were prepared from cryopreserved xenograft fragments using a Tumour Dissociation Kit  
670 (MACS Miltenyi Biotec, cat no. 130-095-929) following the protocol for tough tumors. PDTCs  
671 were filtered through a 40  $\mu\text{m}$  strainer and washed by centrifugation with complete growth  
672 media: RPMI-1640, supplemented with serum-free B27, EGF (20 ng/ml), FGF (20 ng/ml),  
673 Penicillin-Streptomycin (50 U/ml) and Gentamicin (5  $\mu\text{g/ml}$ ). Cells were plated to  
674 approximately 1.5 million cells/ml in 384-well plates. PDTCs were cultured for 24 hours and  
675 the PDTC compound screen was performed as described by Bruna et al.<sup>17</sup>. 9 different PDTX  
676 models (AB521M, HCI005, HCI009, STG139M, STG143, STG201, STG316, STG331,  
677 VHIO098) were treated with different concentrations of 3 different small molecules (i-PDS,  
678 PDS, CX-5461) for 14 days. 3 technical replicates were performed, and 3 models were  
679 analyzed within an independent screen. i-PDS and PDS were employed at 10, 3, 1, 0.3, 0.1,  
680 0.03 and 0.01  $\mu\text{M}$ . Due to solubility, CX-5461 treatments were 100x lower in comparison to i-  
681 PDS/PDS. Cell viability was assessed at day 0 and after 14 days of G4-ligand treatment using  
682 CellTiter-Glo 3D (Promega, cat no. G968). To correlate G4 ligand PDTC response with qG4-  
683 ChIP-seq signatures, area under the curves (AUC) were extracted from PDTC G4-ligand dose-  
684 response curves, fitted using isotonic regression, and scattered against  $\Delta\text{G4R/CNA}$  or  
685  $\text{CG4R/CNA}$  signatures.

686

687 ***Fluorescence quench equilibrium dissociation binding assay for PDS and i-PDS.*** The assay  
688 was performed as reported elsewhere<sup>37</sup>. The chemical synthesis of i-PDS is described in the  
689 supplementary information (**Supplementary Data 1**). Cy5-labelled oligonucleotides were  
690 analyzed as previously described<sup>37</sup> (see **Supplementary Table 5**).

691

692 ***Animal experiments and human research participants.***

693 The research was done with the appropriate approval by the National Research Ethics Service,  
694 Cambridgeshire 2 REC (REC reference number: 08/H0308/178), which were all obtained  
695 under the appropriate Institutional Review Boards and transferred to Cambridge under  
696 Materials Transfer Agreements. All animal experiments were conducted in compliance with  
697 the rigorous Home Office framework of regulations (Project License 707679). Full names of  
698 the ethics committee: Revd. Dr. Derek Fraser, Mrs. Beth Midgley, Mr. Adam Garretty. The  
699 mouse strain NOD.Cg-Prkdcscid Il2rgtm1Wjl/SzJ was used as PDTX avatar. Sex of mice:  
700 female. Age of mice: 3 month. Housing conditions for mice: 21  $^{\circ}\text{C}$ , Humidity: 55%  $\pm$  10%,  
701 light/dark cycle 12 h on, 12 h off. All patients were women with breast cancer. Patients were  
702 recruited by the Cambridge Cancer Centre. The covariate-relevant population characteristics  
703 of the breast cancer patients from the Cambridge Cancer Center (e.g. age, genotypic  
704 information) are reported in Supplementary Table 1.

705

706 **Data Availability**



707 The qG4-ChIP-seq data reported in this paper are available at GEO (NCBI repository),  
708 accession number GSE152216. Gene expression (RNA-seq) data of the PDTX models are  
709 available at the European Genome-Phenome Archive, accession number EGAS00001001913.  
710

#### 711 **Code availability**

712 Sample sheets describing the detailed experimental design are available at  
713 <https://github.com/sblab-bioinformatics/qG4-ChIP-seq-of-breast-cancer-PDTX/>. Details of  
714 data analysis have been deposited at [https://github.com/sblab-bioinformatics/qG4-ChIP-seq-](https://github.com/sblab-bioinformatics/qG4-ChIP-seq-of-breast-cancer-PDTX/)  
715 [of-breast-cancer-PDTX/](https://github.com/sblab-bioinformatics/qG4-ChIP-seq-of-breast-cancer-PDTX/). An overview of all software tools for the processing of sequencing  
716 data is available (see **Supplementary Table 6**).  
717

#### 718 **Methods-only References:**

- 719 30. Schmidt, D. *et al.* ChIP-seq: Using high-throughput sequencing to discover protein–DNA  
720 interactions. *Methods* **48**, 240–248 (2009).
- 721 31. Biffi, G., Tannahill, D., McCafferty, J. & Balasubramanian, S. Quantitative visualization of  
722 DNA G-quadruplex structures in human cells. *Nat. Chem.* **5**, 182–6 (2013).
- 723 32. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.  
724 (2013).
- 725 33. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for  
726 differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139–140  
727 (2010).
- 728 34. Khan, A. & Mathelier, A. Intervene: a tool for intersection and visualization of multiple gene  
729 or genomic region sets. *BMC Bioinformatics* **18**, 287 (2017).
- 730 35. Chin, S.-F. *et al.* Shallow whole genome sequencing for robust copy number profiling of  
731 formalin-fixed paraffin-embedded breast cancers. *Exp. Mol. Pathol.* **104**, 161–169 (2018).
- 732 36. Huang, W., Loganantharaj, R., Schroeder, B., Fargo, D. & Li, L. PAVIS: a tool for Peak  
733 Annotation and Visualization. *Bioinformatics* **29**, 3097–3099 (2013).
- 734 37. Le, D. D., Di Antonio, M., Chan, L. K. M. & Balasubramanian, S. G-quadruplex ligands  
735 exhibit differential G-tetrad selectivity. *Chem. Commun.* **51**, 8048–8050 (2015).  
736