

1 **High-resolution sequencing of DNA G-quadruplex secondary structures in the human**
2 **genome**

3 Vicki S. Chambers^{1*}, Giovanni Marsico^{2*}, Jonathan M. Boutell³, Marco Di Antonio^{1,2}, Geoffrey
4 P. Smith³, Shankar Balasubramanian^{1,2,4}

5 1. Department of Chemistry, University of Cambridge, Cambridge CB2 1EW, UK.

6 2. Cancer Research UK, Cambridge Research Institute, Li Ka Shing Centre, Cambridge CB2
7 0RE, UK.

8 3. Illumina Cambridge Ltd., Chesterford Research Park, Little Chesterford, Saffron Walden,
9 Essex, CB10 1XL, UK.

10 4. School of Clinical Medicine, University of Cambridge, Cambridge CB2 0SP, UK.

11 * These authors contributed equally to this work.

12

13

14 During active transcription and replication chromatin architecture is altered, allowing formation
15 of DNA secondary structures¹. G-quadruplexes (G4s) have emerged as important regulatory DNA
16 structures and have been associated with genomic instability, genetic diseases and cancer
17 progression²⁻⁴. Experimental evidence for G4 prevalence in the entire human genome is still
18 lacking. We present a high-resolution sequencing-based method that detected 716,310 distinct
19 G4s in the human genome, more than predicted by computational methods⁵⁻⁷, including structural
20 variants previously uncharacterised in a genomic context^{8,9}. We observed high G4-density in
21 functional regions, such as 5' UTRs and splicing sites, and in genes not predicted to have such
22 structures (*BRCA1* and *BRCA2*). We found a significant association of G4 formation with
23 oncogenes and tumor suppressors, and with Somatic Copy-Number Alterations (SCNAs) that act
24 as cancer drivers¹⁰. Our results support that G4s are promising targets for cancer intervention and
25 suggest novel candidates for further biological and mechanistic studies.

26

27

28

29

30

31

32 The formation of DNA and RNA secondary structures is of vital importance to fundamental
33 biological processes, such as replication, translation and splicing^{11,12}. While RNA structure-
34 mapping on a genomic-scale is established^{13,14}, extending these methodologies to interrogate
35 DNA secondary-structure formation remains a challenge. G4s are a particular class of DNA
36 secondary structures that is emerging as a regulatory element for key biological processes and an
37 important therapeutic target²⁻⁴. G4 structures can form in guanine-rich sequences from the
38 interaction of four guanine bases to generate a planar G-tetrad, which can subsequently self-
39 stack¹⁵. G4 formation is kinetically fast and they are thermodynamically very stable under
40 physiological conditions, particularly in the presence of K⁺¹⁵. Recently, G4 formation has been
41 visualised in human cells and tissues by means of immuno-fluorescence¹⁶⁻¹⁸. These and other
42 studies highlight the importance of G4 formation in specific genes, underpinning the value of
43 studying these structures at a larger scale. The formation of G4s can be assessed *in vitro* by
44 measuring the stalling of a polymerase along its template at G4 sites (polymerase stop assay)¹⁹.
45 Here, we adapt the polymerase stop assay together with Illumina® next-generation sequencing²⁰
46 to establish G4-Seq, the first method to detect and map DNA secondary structures on a genome-
47 wide scale. We altered sequencing conditions to either disfavour or promote G4 formation on the
48 sequencing array, comparing the respective sequencing readouts to elucidate the exact position of
49 the DNA structure (Fig. 1). We used two independent approaches to promote DNA G4
50 stabilisation: i) adding K⁺; ii) adding the G4 stabilising ligand pyridostatin (PDS, 1 μM)²¹. For
51 each condition, we compared sequencing quality and base calling before and after G4
52 stabilisation in a human genomic DNA library spiked with four known control sequences
53 (Methods, Fig. 1): two containing stable G4 structures (*c-myc* and *c-kit*), one mutated to prevent
54 G4 formation (*c-myc mut*) and the complementary C-rich strand of *c-myc* (*c-myc-opp*) that cannot
55 fold into a G4.

56 In our experiments, we supplemented standard Illumina sequencing buffers with either 50 mM
57 LiCl or NaCl, which do not cause strong G4 stabilization, or KCl that does stabilizes G4
58 structure²² (Methods), keeping the ionic strength of all buffers constant. The overall sequencing
59 quality, as quantified by Phred Quality scores²³ (Q, Methods), was not globally affected by any of
60 the added cations (Extended Data Fig. 1). However, quality was reduced *only* in the presence of
61 K⁺ for a subset of sequences, including the G4-positive controls *c-myc* and *c-kit* and sequences
62 computationally predicted to form a G4⁵. Conversely, the G4-negative controls *c-myc-opp* and *c-*
63 *myc-mut* showed no change in quality under any condition (Extended Data Fig. 2a). Sequencing
64 of the controls under Li⁺ and Na⁺ conditions revealed no alterations compared to the known input

65 sequences (i.e. base mismatches <2%), whereas under K⁺ conditions the G4-positive controls *c-*
66 *kit* and *c-myc* displayed 34% and 46% mismatches respectively (Extended Data Fig. 2b).
67 Therefore, we sequenced each genomic DNA template twice, with an initial sequencing run
68 (Read-1) in Na⁺, to ensure accurate sequencing and correct identification by alignment to the
69 human reference genome (*hg19*), and a second sequencing run (Read-2) under G4 stabilising
70 conditions (K⁺), to detect structure formation by mismatch quantification based on the sequence
71 obtained in Read-1.

72 We next explored whether specific stabilisation of G4s by the ligand PDS, previously shown to
73 induce polymerase stalling at G4 sites in cells²⁴, could also induce targeted sequencing errors.
74 We performed Read-1 in Na⁺ and Read-2 under the same cation conditions but with addition of
75 PDS (1 μM, *Methods*). Herein, we measured mismatches of 45% for *c-kit* and 66% for *c-myc* but
76 little effect (< 5% mismatches) for *c-myc-opp* and *c-myc-mut* that are unable to form G4s
77 (Extended Data Fig. 3). The inspection of mismatches along the *c-kit* control, which contains two
78 independent G4 motifs *c-kit1* and *c-kit2*,^{25,26} revealed that sequencing errors accumulated only
79 after the G4 start sites, suggesting that under both K⁺ and PDS conditions the formation of DNA
80 G4s cause polymerase stalling and mismatches in sequencing readout (Fig. 2a). In fact, when the
81 polymerase encounters a stable G4 in the DNA template a pausing is induced, which can
82 effectively truncate the reading of the template sequence. When this happens, the sequencer will
83 continue to generate what appears to be a scrambled sequence beyond this point, as illustrated by
84 Supplementary Figures 1 and 2. Ordinarily such reads are removed during the data analysis,
85 whereas we have retained them in our experiment to detect G4 sites. Our approach therefore
86 enables both the identification of G4-containing sequences and the exact location of the structure.
87 Interestingly, only PDS addition induced significant polymerase stalling at *c-kit1* in agreement
88 with the relative stability of the two G4s²⁵.

89 The analysis of 32 million reads, comprising a subset of ~110,000 Predicted Quadruplexes⁵
90 (PQs), showed higher mismatch-levels (median of 20% in K⁺ and 35% in PDS) in sequences
91 containing PQs as opposed to those without (non PQs; < 2%) (Fig. 2b). Mismatch levels were
92 generally high (> 38%) immediately after the PQ motif and negligible (< 1%) beforehand (Fig.
93 2c), confirming a G4-dependent effect, as observed for *c-kit*. Although, mismatch levels for non-
94 PQs were low on average (< 2%), a small fraction (~0.01) was found to have relatively high
95 mismatch levels (> 20%; ~149,000 sequences in K⁺ and ~216,000 in PDS), far greater than the
96 number of predicted PQs (~110,000; Fig 2b). Thus, suggesting that the number and nature of
97 human genomic G4s is substantially broader than previously predicted⁵.

98 This method, which we call G4-Seq, was applied to generate a high-resolution map of G4
99 structures in the human genome (*NA18507*, Methods), using the Illumina HiSeq platform, under
100 Na^+ conditions in Read-1 and either K^+ or PDS in Read-2. Each experiment was performed in
101 duplicate and yielded at least 285 million reads with an average coverage of 14x for the human
102 genome (Supplementary Table 1). We set thresholds of 25% and 18% mismatches for PDS and
103 K^+ , respectively, to ensure a similar false positive rate of $\sim 2\%$ (Methods). Thus, any read with
104 mismatches above these thresholds is considered a reliable indication of G4 formation and is
105 termed observed G4 sequence (OQ). By applying these criteria, we identified 716,310 OQs in
106 PDS and 525,890 OQs in K^+ within the human genome. Furthermore, 73% (in PDS) and 60% (in
107 K^+) of all 361,424 predicted canonical G4 forming sequences (PQs) were present in the
108 experimentally detected OQs (Extended Data Table 1). 90% of PQs found in K^+ were also
109 detected in PDS and 383,984 of the overall OQs were common to both conditions ($p < 10^{-16}$). The
110 high overlap between distinct G4 stabilising conditions provides independent validation of the
111 assignment of OQs. Our data indicates that the OQs detected exclusively with PDS do in fact also
112 display significantly high mismatch levels in K^+ (compared to random genomic intervals) and
113 accordingly for OQs detected exclusively in K^+ (Supplementary Figure 3), suggesting that it is the
114 extent of stabilisation under a given set of conditions that affects the likelihood of a G4 being
115 detected by G4-seq. The OQs detected in the presence of PDS could also reflect the binding
116 properties and specificity of the small-molecule for G4 stabilisation²⁷. The use of a different G4-
117 stabilising ligand, PhenDC3²⁸, showed a strong overlap (85%) with OQs detected in PDS
118 (Supplementary Figure 4), suggesting that no major differences in binding specificity were
119 observed with these two ligands.

120 Notably, the majority ($\sim 70\%$) of the OQs were actually *not* predicted from a classical description
121 of G4 structure⁵. Recent structural and biophysical studies have identified a small number of
122 cases of stable non-canonical G4 structures in which either the loops are exceptionally long (>7
123 bases)^{9,29}, or a discontinuity in the G-tracts leads to bulges⁸ (Extended Data Fig. 4). To elucidate
124 distinct structural features, the OQs were grouped as follows (Methods): 1) Canonical PQs: in
125 three categories according to loop length; 2) Long loops: sequences with any loop > 7 bases; 3)
126 Bulges: sequences with single-nucleotide interruptions in one or more of the G-runs or a longer
127 interruption in one G-run (e.g. GGH_{1-7}G); 4) Other: sequences not belonging to the previous
128 categories (Fig. 3a). Structural families are defined by a hierarchical assignment based on
129 sequence only (Methods). There is potential for multiple folding scenarios or polymorphism, that
130 is not accounted for in our assignment, but which could be assessed by dedicated structural

131 studies on a case-by-case basis. Long loops and Bulges accounted respectively for 21.5% and
132 21.6% of total OQs in K⁺ and 24% and 30% in PDS. The remaining OQs (category Other) may
133 have the potential to form G4s, such as structures containing multi-nucleotide bulges, two-tetrads
134 G4s, or topologies comprising both long loops and bulges (Extended Data Table 2). Collectively,
135 these findings have unraveled a dataset of stable G4 sequences that could not have been easily
136 identified *a priori* in genomic DNA by computational approaches.

137 We measured the fold enrichment of OQs compared to random genomic intervals to assess the
138 likelihood of each class to be detected by G4-Seq (Methods). Sequences with short loops have
139 high enrichment (>25 fold) under both PDS and K⁺ conditions, whereas sequences with longer
140 loops or bulges displayed lower enrichment (<15 fold; Fig. 3b) consistent with the relative
141 thermodynamic stability of the different G4 structures^{8,9,30}. Also, less stable G4s were more easily
142 detected by PDS (Extended Data Fig. 5).

143 To understand the potential functions of G4s we evaluated the existence of OQs in genomic
144 regions associated with promoters, 3' and 5'-UTRs, exons, introns and splicing junctions
145 (Extended Data Table 3). Notably, a large proportion of these regions (up to 49% in PDS and
146 46% in K⁺) comprise *exclusively* non-canonical G4s (i.e. Long loops or Bulges). The highest
147 density of G4s was found in 5' UTRs and splicing sites, consistent with a role in post-
148 transcriptional regulation, as supported by the recent finding in the 5' UTR of *eIF4A*².

149 Visual inspection of genes with biologically important G4s (*SRC*, *MYC*)^{24,31} or genes rich in PQs
150 (*MYL5*, *MYL9*; Fig. 4a, Extended Data Fig. 6) confirmed that G4-Seq is a powerful tool to
151 identify both predicted and uncharacterised G4s, and is highly specific for the G-rich strand
152 (Extended Data Fig. 7, Supplementary Table 2). We found non-canonical G4s within many genes
153 that have few or no PQs (Supplementary Table 3), including important cancer-related genes such
154 as *BRCA1*, *BRCA2* and *MAP3K8*. Genes with a high number of G4s may be particularly sensitive
155 to treatment with G4-stabilising ligands, as shown for the oncogene *SRC*²⁴. Our experimental map
156 also identified oncogenes and tumor suppressors with a notably high G4 density, such as *CUL7*,
157 *FOXAI*, *TUSC2* and *HOXB13* (Supplementary Table 4). This map further revealed significant
158 enrichment of G4s ($p = 4.5e^{-8}$) in somatic copy number alterations (SCNAs), which are signatures
159 of cancer¹⁰ (Fig. 4b). In particular, high G4 density is observed in regions containing oncogenes
160 such as *MYC*, *TERT*, *AKT1*, *FGFR3* and *BCL2L1* (Supplementary Table 5) that specifically relate
161 to SCN amplifications ($p = 2e^{-7}$) rather than deletions ($p = 0.01$). This is consistent with a
162 mechanistic link between G4s and the sites of genomic instability, a hallmark of cancer^{3,32}.

163 We have established a high-throughput, genome-wide method that profiles G4 DNA secondary
164 structure with high resolution. Our study reveals new insights into the nature of G4s that form in
165 the human genome, including non-canonical structural features. Our experimental dataset shows
166 enrichment of G4s in regulatory regions, in addition to oncogenes and SCNAs and provides a
167 resource of novel genomic targets for further biological and mechanistic studies and potential
168 future therapeutic intervention. We anticipate that our approach can be extended to study the
169 prevalence of G4s, and potentially other DNA secondary structures, in any genome. Furthermore,
170 G4-Seq can be exploited to detect DNA-small molecules interaction in a genomic context.

171

- 172 1 Rodriguez, R., Miller, K. M. Unravelling the genomic targets of small-molecules using
173 high-throughput sequencing *Nat. Rev. Genet.* **15**, 783-96, (2014).
- 174 2 Wolfe, A. L. *et al.* RNA G-quadruplexes cause eIF4A-dependent oncogene translation in
175 cancer. *Nature* **513**, 65-70, (2014).
- 176 3 Maizels, N. Genomic stability: FANCD1-dependent G4 DNA repair. *Curr. Biol.* **18**, R613-
177 614, (2008).
- 178 4 Haeusler, A. R. *et al.* C9orf72 nucleotide repeat structures initiate molecular cascades of
179 disease. *Nature* **507**, 195-200, (2014).
- 180 5 Huppert, J. L. & Balasubramanian, S. Prevalence of quadruplexes in the human genome.
181 *Nucleic Acids Res.* **33**, 2908-2916, (2005).
- 182 6 Eddy, J. & Maizels, N. Gene function correlates with potential for G4 DNA formation in
183 the human genome. *Nucleic Acids Res.* **34**, 3887-3896, (2006).
- 184 7 Kikin, O., D'Antonio, L. & Bagga, P. S. QGRS Mapper: a web-based server for
185 predicting G-quadruplexes in nucleotide sequences. *Nucleic Acids Res.* **34**, W676-682,
186 (2006).
- 187 8 Mukundan, V. T. & Phan, A. T. Bulges in G-quadruplexes: broadening the definition of
188 G-quadruplex-forming sequences. *J. Am. Chem. Soc.* **135**, 5017-5028, (2013).
- 189 9 Guedin, A., Gros, J., Alberti, P. & Mergny, J. L. How long is too long? Effects of loop
190 size on G-quadruplex stability. *Nucleic Acids Res.* **38**, 7858-7868, (2010).
- 191 10 Zack, T. I. *et al.* Pan-cancer patterns of somatic copy number alteration. *Nat. Genet.* **45**,
192 1134-1140, (2013).
- 193 11 Bochman, M. L., Paeschke, K. & Zakian, V. A. DNA secondary structures: stability and
194 function of G-quadruplex structures. *Nat. Rev. Genet.* **13**, 770-780, (2012).
- 195 12 Cruz, J. A. & Westhof, E. The dynamic landscapes of RNA architecture. *Cell* **136**, 604-
196 609, (2009).
- 197 13 Ding, Y. *et al.* In vivo genome-wide profiling of RNA secondary structure reveals novel
198 regulatory features. *Nature* **505**, 696-700, (2014).
- 199 14 Watts, J. M. *et al.* Architecture and secondary structure of an entire HIV-1 RNA genome.
200 *Nature* **460**, 711-716, (2009).
- 201 15 Davis, J. T. G-quartets 40 years later: from 5'-GMP to molecular biology and
202 supramolecular chemistry. *Angew. Chem. Int. Ed.* **43**, 668-698, (2004).
- 203 16 Biffi, G., Tannahill, D., McCafferty, J. & Balasubramanian, S. Quantitative visualization
204 of DNA G-quadruplex structures in human cells. *Nat. Chem.* **5**, 182-186, (2013).
- 205 17 Henderson, A. *et al.* Detection of G-quadruplex DNA in mammalian cells. *Nucleic Acids*
206 *Res.* **42**, 860-869, (2014).

- 207 18 Biffi, G., Tannahill, D., Miller, J., Howat, W. J. & Balasubramanian, S. Elevated levels of
208 G-quadruplex formation in human stomach and liver cancer tissues. *PloS one* **9**, e102711,
209 (2014).
- 210 19 Weitzmann, M. N., Woodford, K. J. & Usdin, K. The development and use of a DNA
211 polymerase arrest assay for the evaluation of parameters affecting intrastrand tetraplex
212 formation. *J. Biol. Chem.* **271**, 20958-20964, (1996).
- 213 20 Bentley, D. R. *et al.* Accurate whole human genome sequencing using reversible
214 terminator chemistry. *Nature* **456**, 53-59, (2008).
- 215 21 Rodriguez, R. *et al.* A novel small molecule that alters shelterin integrity and triggers a
216 DNA-damage response at telomeres. *J. Am. Chem. Soc.* **130**, 15758-15759, (2008).
- 217 22 Hud, N. V., Smith, F. W., Anet, F. A. L. & Feigon, J. The selectivity for K⁺ versus Na⁺
218 in DNA quadruplexes is dominated by relative free energies of hydration: A
219 thermodynamic analysis by H-1 NMR. *Biochemistry* **35**, 15383-15390, (1996).
- 220 23 Ewing, B., Hillier, L., Wendl, M. C., Green, P. Base-Calling of Automated Sequencer
221 Traces Using Phred. 1. Accuracy Assessment *Genome Research* **8**, 175-185, (1998).
- 222 24 Rodriguez, R. *et al.* Small-molecule-induced DNA damage identifies alternative DNA
223 structures in human genes. *Nat. Chem. Biol.* **8**, 301-310, (2012).
- 224 25 Fernando, H. *et al.* A conserved quadruplex motif located in a transcription activation site
225 of the human c-kit oncogene. *Biochemistry* **45**, 7854-7860, (2006).
- 226 26 Rankin, S. *et al.* Putative DNA quadruplex formation within the human c-kit oncogene. *J.*
227 *Am. Chem. Soc.* **127**, 10584-10589, (2005).
- 228 27 Marchand, A. *et al.* Ligand-Induced conformational changes with cation ejection upon
229 binding to human telomeric DNA G-quadruplexes. *J. Am. Chem. Soc.* **137**, 750-756,
230 (2015).
- 231 28 De Cian, A., DeLemos, E., Mergny, J-L., Teulade-Fichou, M-P., Monchaud, D. Highly
232 efficient G-quadruplex recognition by Bisquinolinium compounds. *J. Am. Chem. Soc.*
233 **129**, 1856-1857, (2007).
- 234 29 Palumbo, S. L., Ebbinghaus, S. W., Hurley, L. H. Formation of a Unique End-to-End
235 Stacked Pair of G-Quadruplexes in the hTERT Core Promoter with Implications for
236 Inhibition of Telomerase by G-Quadruplex-Interactive Ligands. *J. Am. Chem. Soc.* **131**,
237 10878-10891, (2009).
- 238 30 Bugaut, A. & Balasubramanian, S. A sequence-independent study of the influence of
239 short loop lengths on the stability and topology of intramolecular DNA G-quadruplexes.
240 *Biochemistry* **47**, 689-697, (2008).
- 241 31 Siddiqui-Jain, A., Grand, C. L., Bearss, D. J. & Hurley, L. H. Direct evidence for a G-
242 quadruplex in a promoter region and its targeting with a small molecule to repress c-
243 MYC transcription. *Proc. Natl. Acad. Sci. U S A* **99**, 11593-11598, (2002).
- 244 32 Paeschke, K. *et al.* Pif1 family helicases suppress genome instability at G-quadruplex
245 motifs. *Nature* **497**, 458-462, (2013).

246
247
248
249

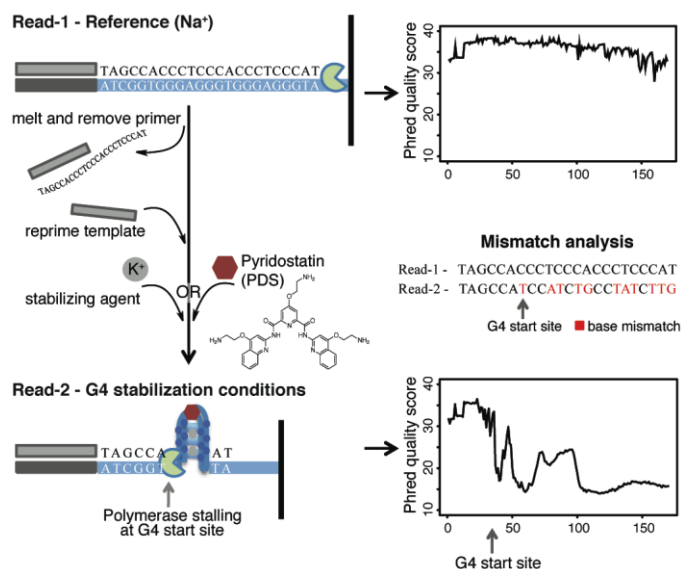
Supplementary Information is available in the online version of the paper.

250 **Acknowledgements** We thank Dr. Chris Lowe and Dr. David Tannahill for critical reading of the
251 manuscript and Dr. Dario Beraldi for technical support. We thank Patrick McCauley (Illumina)
252 who prepared the custom sequencing buffers. We are grateful to the Biotechnology and
253 Biological Sciences Research Council (BBSRC) and Illumina® for the studentship supporting

254 V.C (BB/I015477/1). The S.B. research group is supported by programme funding from Cancer
255 Research UK and from the European Research Council and project funding from BBSRC.

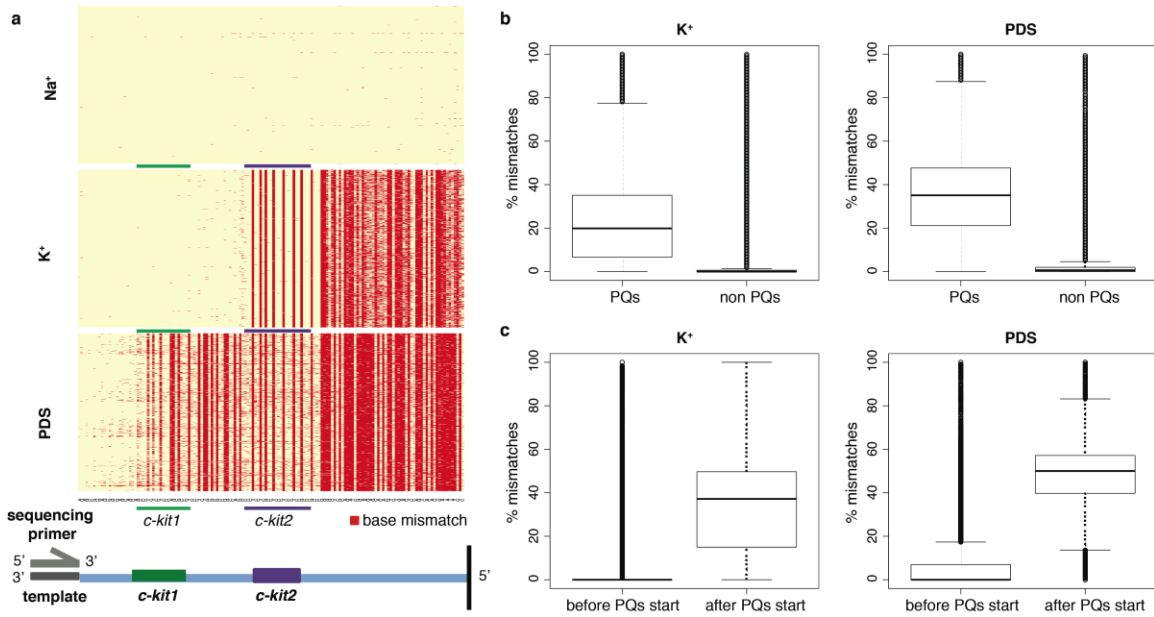
256 **Author Contributions** V.C. and J.B. carried out the experiments. G.M. designed, implemented
257 and performed the analysis. All authors designed the experiments. V.C., G.M., M.D.A. and S.B.
258 interpreted the results and co-wrote the manuscript with input from all authors.

259 **Author information** The data reported in this paper is available at the NCBI's GEO repository,
260 accession number GSE63874 (<http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE63874>).
261 Reprints and permissions information is available at www.nature.com/reprints. The authors
262 declare no competing financial interests. Correspondence and requests for materials should be
263 addressed to S.B. (sb10031@cam.ac.uk).



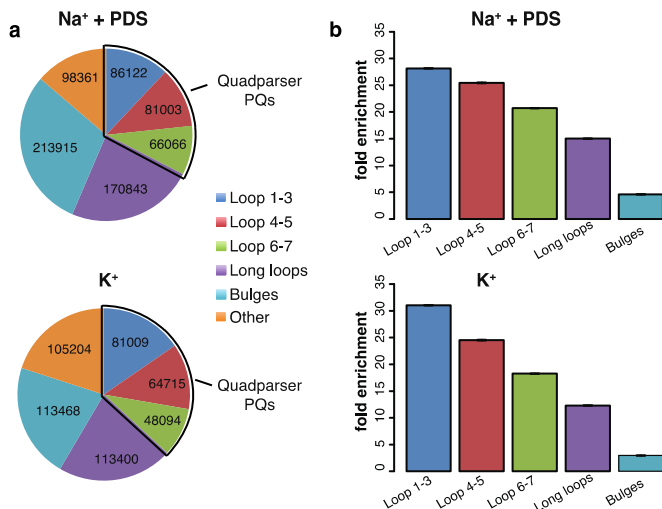
264

265 **Figure 1: A schematic of the G4-Seq method.** In a typical G4-Seq experiment sequencing is
266 performed twice. A first sequencing run under Na⁺ conditions (Read-1) enables accurate
267 sequencing and alignment of DNA fragments. Subsequently, the DNA synthesised during
268 sequencing is removed and the original template re-sequenced (Read-2) under conditions that
269 promote G-quadruplex (G4) stabilization: either by the addition of the G4-ligand PDS or by
270 supplementing sequencing buffers with K⁺. G4-induced polymerase stalling alters the sequencing
271 readout from the beginning of the G4 structure resulting in a drop in sequencing quality from that
272 point in Read-2 only. Differences in sequencing quality and mismatches between Read-1 and
273 Read-2 are analysed to provide a map of G4 structures in the human genome.



274

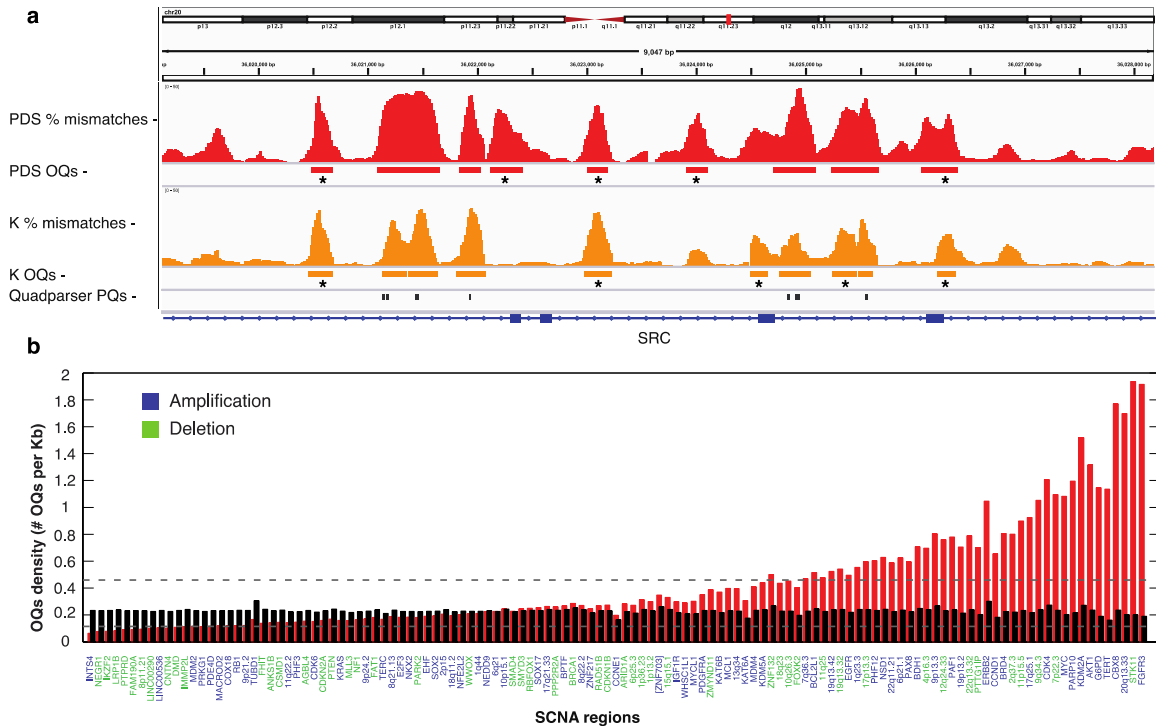
275 **Figure 2: Analysis of G4-seq for known G-quadruplex sequences.** a) Identification of base
 276 mismatches for the *c-kit* control sequence depicted in a heat-map plot under different sequencing
 277 conditions. Each row is an independent sequenced template, while each column corresponds to
 278 each the sequenced bases. Yellow background indicates no difference to the known input
 279 sequence and the sequence experimentally obtained, while red indicates mismatches. The two G-
 280 quadruplex motifs are indicated in green (*c-kit1*) and purple (*c-kit2*). Top: sequencing in Na^+ ,
 281 where negligible mismatches were observed. Middle: sequencing in K^+ showed mismatches
 282 accumulation starting at *c-kit2*, thus suggesting polymerase stalling. Bottom: sequencing in
 283 presence of PDS revealed stalling already at the first G4 motif (*c-kit1*) and significant mismatch
 284 accumulation. b) Boxplots showing the mismatch percentage between Read-1 and Read-2 for
 285 reads with *Quadparser*-predicted PQs (PQs; $N \sim 110,000$) and without (non-PQs; $N \sim 32$ million)
 286 for K^+ (left) and PDS (right). c) Boxplots representing the percentage mismatches for the reads
 287 containing a PQ, before or after the motif start site, for K^+ (left) and PDS (right).



288

289 **Figure 3: Structural analysis of Observed G-quadruplex sequences (OQs).** a) Number of
 290 OQs found in different G-quadruplex structural families, for Na⁺ + PDS or K⁺ sequencing
 291 conditions (Methods). The different families are defined as follows. Loop 1-3; Loop 4-5; Loop 6-
 292 7: OQs with at least one loop of the indicated length; Long loops: OQs with any loop of length >
 293 7; Bulges: OQs with a bulge of 1-7 bases in one G-run or multiple 1-base bulges; Other:
 294 sequences which do not fall into the categories above. b) Fold enrichment (ratio) of each
 295 structural family represented in OQs over random genomic sequences measured for Na⁺ + PDS
 296 (top) and K⁺ (bottom) conditions. Error bars are SEM of 3 independent randomizations. Fold
 297 enrichment values follow the relative thermodynamic stability of the different G4 families, with
 298 highest enrichment for G4 structures with short loops compared to longer loop counterparts.
 299 Treatment with PDS enables the detection of G4 structural variants with lower intrinsic stability.

300



301

302 **Figure 4: Genomic distribution of experimentally-determined OQs.** a) Genome browser view
 303 of PQs and OQs across the *SRC* oncogene. Red and orange show tracks for mismatches in reads
 304 aligning to the reverse strand (-) for PDS and K^+ , respectively. Bars of the same colors indicate
 305 OQ regions above threshold, while black bars indicate Quadparser PQs. OQs not predicted by
 306 Quadparser are indicated by * (sequences in Supplementary Table 2). G4-Seq enables detection
 307 of G4 sequences that were not previously predicted by computational algorithms in this
 308 oncogene. b) OQs density (red) in different SCNAs¹⁰ compared to random intervals (black),
 309 measured as number of OQs per kilobase. Blue gene labels: SCNAs representing amplifications.
 310 Green gene labels: SCNAs representing deletions. Dotted lines: values corresponding to 0.5 and 2
 311 times the average random density (0.22). Bars are sorted according to the fold enrichment of OQs
 312 density over random (Supplementary Table 5). A correlation was observed between the OQs
 313 density per kilobase and with SCNAs associated with amplifications, suggesting a potential role
 314 of G4 structure in carcinogenesis.

315

316

317

318

319 **Methods**

320 **Design of control sequences**

321 Full-length control sequences (sequence of interest underlined) are as follows:

322 *Control 1 (Positive): c-kit*

323 5'-Adapter 1-AGAGCCGCGAGCGGGCAGCAGCAGCCCTCTCCTCCCAGCGCCCTCCCTCTGCGCGCCGG
324 CCACGCCCTCCTCGCTCCCTCCCTCCGCCCGCCCGGGGCTCGCG-Adapter 2-3'.

325 *Control 2 (Negative): c-myc-opp*

326 5'-Adapter 1- ATTAGCGAGAGAGGATCTTTTTCTTTTCCCCACGCCCTCTGCTTTGGGAACCCGGGA
327 GGGCGCTTATGGGGAGGGTGGGGAGGGTGGGGAAGGGGGAGGAGAG-Adapter 2-3'.

328 *Control 3 (Positive): c-myc*

329 5'-Adapter 1- TCTCTCCCCACCTTCCCCACCTCCCCACCTCCCCATAAGCGCCCTCCCGGGTCCC
330 AAAGCAGAGGGCGTGGGGGAAAAGAAAAAAGATCCTCTTCGCTAATAG-Adapter 2-3'.

331 *Control 4 (Negative): c-myc-mut*

332 5'-Adapter 1- CTCCTCTTCACCTTCTTCACTCTTCACTCTCTTCATAAGCGCCCTCCCGGGTCCCAA
333 AGCAGAGGGCGTGGGGGAAAAAAAAAAGATCCTCTCTCGCTAATAG-Adapter 2-3'.

334

335 where:

336 Adapter 1- 5'-AATGATACGGCGACCACCGAGATCTACACTCTTCCCTACACGACGCTCTCCGATCT-3'

337 Adapter 2- 5'-AGATCGGAAGAGCACACGTCTGAACTCCAGTCACACTGATATATCTCGTATGCCGTCTT
338 CTGCTTG-3'

339

340 The *c-myc* and *c-kit* positive controls were designed based on the human genomic sequence of
341 two regions in the promoter of the oncogenes *MYC* and *KIT*, respectively, which are well-studied
342 examples of G-quadruplex (G4) forming motifs^{25,26,28}. Crucially, controls were designed
343 complementary to the G4 motif *i.e.* the C-rich sequence to ensure that during Illumina cluster
344 generation the G-rich sequence becomes immobilised to the flow cell surface and acts as the
345 template for sequencing. This protocol is necessary to allow the study of G4 structures on
346 polymerase procession. Two negative control sequences were also designed based on the *c-myc*
347 sequence: 1) *c-myc-opp*: the complementary G-rich strand of the *c-myc* G4, which becomes the
348 C-rich template sequence upon cluster generation; 2) *c-myc-mut*: a mutant of *c-myc* that can no
349 longer form a G4.

350 **Control sequence library preparation**

351 Synthetic oligonucleotides of the control sequences, and their complement sequences, with a 5'-
352 phosphate group and an A overhang (Biomers) were prepared using nuclease free water at the
353 final concentration of 1 µg/ml. The two complementary oligonucleotide sequences of each

354 control (100 ng/μl) were annealed in 10 mM Tris, 50 mM NaCl buffer by heating to 95 °C for 10
355 min and then cooled to 20 °C at 1 °C/min. The annealed DNA was prepared for Illumina
356 sequencing by ligation of Illumina adapters using a T4 DNA ligase at 30 °C for 10 min.
357 Following AMPure® bead clean-up, the adapted sequences were PCR amplified using standard
358 Illumina PCR primers and gel purified (Qiagen MinElute Gel Extraction kit). Purified fragments
359 were ligated into Life technologies PCR®-Blunt Vectors and transformed according to standard
360 methods. Plasmid DNA was purified from selected clones (Thermo scientific GeneJET plasmid
361 Miniprep Kit), followed by Sanger sequencing (GATC) to confirm the sequence identity and
362 directionality. DNA inserts of the chosen clones (C-rich variant of the insert in the case of *c-myc*,
363 *c-kit* and *c-myc-mut* and G-rich for *c-myc-opp*) were isolated by EcoRI-HF digestion and gel
364 purification to generate sequences ready for use in sequencing. Sequences were quantified using a
365 Qubit Fluorimeter (Life Technologies) and denatured according to standard Illumina protocols.
366 Control sequences were spiked into a human genomic library at a final concentration of 0.01 pM
367 for all sequencing experiments.

368 **Genomic library preparation**

369 Purified Human Genomic DNA isolated from primary human B-lymphocytes (NA18507) was
370 purchased from Coriell Institute for Medical Research and prepared for sequencing using TruSeq
371 DNA sample prep kit (Illumina) according to the manufacturers protocol. Human template DNA
372 was denatured as in standard Illumina protocols and used at 8 pM for sequencing on MiSeq
373 instruments (Illumina) and 12 pM for all sequencing on an Illumina HiSeq 2500 in Rapid Run
374 mode (with the addition of 0.01 pM of each control sequence).

375 **Modified sequencing buffer preparation**

376 In collaboration with Illumina, the standard sequencing buffers (incorporation, wash and cleavage
377 buffers) were supplemented with K⁺, Na⁺ or Li⁺ at a final concentration of 50 mM for the
378 incorporation and wash buffers and 1 M for the cleavage buffer. In addition, for small-molecule
379 experiments with PDS, all buffers were prepared using Na⁺ at 50 mM final concentration, and
380 PDS⁴ (1 μM) was added to the incorporation buffer on the instrument. All other reagents used
381 were from standard proprietary Illumina sequencing kits.

382 **G4-Seq Protocol**

383 Illumina sequencing was performed using either MiSeq or HiSeq 2500 Rapid Run
384 instrumentation, using the same basic protocol. A human genomic library containing synthetic
385 control sequences (prepared as above) was used as template. Cluster generation and amplification

386 were carried out according to standard procedures. The template DNA was then sequenced using
387 buffer conditions containing Na⁺ (Read-1) for 250 cycles (MiSeq) or 150 cycles (HiSeq 2500).
388 The newly synthesised DNA strand was removed by denaturation to leave the original template
389 DNA strand. The Read-1 sequencing primer (HP10) was then added to the flow-cell and
390 hybridised as per standard sequencing protocols. Annealing buffer (10mM Tris and 100mM KCl,
391 pH 7.4) was added to the flow cell and the temperature increased to 65°C for 5 min, followed by
392 cooling to 20 °C at 1 °C/min, in order to promote G4 formation in immobilised template DNA.
393 For sequencing experiments with PDS or PhenDC3, the small-molecule was added to the flow
394 cell (1 µM in annealing buffer) and equilibrated for 30 min at room temperature. Sequencing was
395 then performed on the template DNA (Read-2) in G4-stabilisation conditions, i.e. either K⁺
396 sequencing buffers or with PDS addition in Na⁺ buffer. The sequencing read length was 250 and
397 150 base pairs (bp) for the MiSeq and HiSeq 2500 respectively. Base-calling log (bcl) files from
398 the sequencing run were processed to generate FASTQ files for further analysis.

399 **FASTQ files**

400 The FASTQ format³³ consists of: 1) a read identifier to allow identification of sequences from the
401 same cluster when performing different sequencing reads, hence Read-1 and Read-2; 2) a
402 measure of base-calling quality- the Phred quality score, Q, which is inversely related to the
403 probability that the corresponding base-call is incorrect (i.e. a high Q score indicates a low
404 probability of erroneously calling the given base, while a lower Q score indicates greater
405 probability that the given base is incorrectly called); 3) the actual base-call, where the nucleotide
406 with highest confidence is assigned to each sequencing position. Read quality was calculated as
407 the average Phred quality of all bases; the quality difference was calculated as Read-1 quality
408 minus Read-2 quality; the percentage of mismatches was calculated comparing base calling at
409 Read-1 and Read-2 and counting the fraction of different calls across the whole read.

410 **Different cation analysis**

411 Sequencing was performed in Li⁺, Na⁺ and K⁺ as described above. Two replicates were
412 performed for K⁺ and Li⁺ conditions and three replicates for Na⁺ conditions. FASTQ files were
413 obtained from MiSeq 250 bp single-end reads. Files were aligned to the human genome (*hg19*) by
414 using the bwa mem aligner with default parameters (<http://bio-bwa.sourceforge.net/>).

415 **K⁺ and PDS genomic analysis**

416 Sequencing was performed as described above. Two technical replicates were performed for each
417 G4-stabilisation condition on HiSeq instrumentation. FASTQ files were obtained from HiSeq

418 2500 150 bp single-end reads. FASTQ files from Read-1 were aligned to the human genome
419 (*hg19*) using the bwa mem aligner with default parameters (<http://bio-bwa.sourceforge.net/>). Bam
420 alignment files were processed using bedtools (<https://code.google.com/p/bedtools/>): 1) bam files
421 were converted to bed files (command bamToBed); 2) bed files were expanded 30 bases
422 downstream (command slopBed -s -r 30); 3) expanded bed files were grouped to keep only the
423 best alignments for each read (command groupBy -g 4 -c 5 -o max); 4) FASTA sequence files
424 were extracted from the bed intervals (command bedtools getfasta -s); 5) FASTA sequence files
425 and the FASTQ files from both Read-1 and Read-2 were loaded in R (<http://www.r-project.org/>)
426 for analysis. Sequence tails beyond poly-A tails (≥ 9 bases) were trimmed as they represent the
427 end of the DNA fragment attached to the flow cell. The difference in the quality score and
428 percentage of mismatches (% mismatches) between Read-1 and Read-2 for each individual base
429 was calculated and stored for each read, together with coverage count of +1. All single-base
430 values calculated from the processed reads were then pooled to generate genomic tracks of
431 mismatch percentage (average of values) and total coverage (sum of values). To ease data
432 handling, genomic tracks were finally binned in intervals of length 15 bases and smoothed with a
433 moving average of order 15 (i.e. window size around the point value to be smoothed).

434 **Control sequences analysis**

435 FASTQ files were generated from the MiSeq (cations experiments) or the HiSeq 2500 (K⁺ and
436 PDS experiments) sequencing platforms. FASTQ were aligned to a FASTA file containing only
437 the control sequences by using the bwa mem aligner with default parameters ([http://bio-](http://bio-bwa.sourceforge.net/)
438 [bwa.sourceforge.net/](http://bio-bwa.sourceforge.net/)). The Phred quality score (Q) and the base-calling extracted from reads
439 were successfully aligned to each control sequence then were analysed.

440 **PQs identification and positional analysis**

441 For each sequencing read, the aligned sequence information was extracted as above and PQs were
442 identified according to the *Quadparser* algorithm by searching for the regular expression
443 '(G{3,}[ATGC]{1,7}){3,}G{3,}'. For positional analysis, “before PQs start” is defined as the
444 sequence up to 12 bases upstream of the PQ start site (12 bases is the approximate footprint of
445 DNA polymerase). “After PQs start” is defined as the remaining sequence, from 12 bases
446 upstream the PQ start site until the end of the sequence (excluding any sequencing beyond the
447 poly-A tail).

448 **OQ detection**

449 *Quadparser*-predicted PQs were considered as a positive set (PQs) and reads without PQs as a
450 negative set (non PQs). For all reads, % mismatches were calculated (range 0-100 %). For each
451 threshold t_i , the following numbers were calculated: TP_i - true positives i.e. reads with PQs above
452 the threshold t_i , FP_i -false positives, i.e. reads without PQs above the threshold t_i , FN_i - false
453 negatives i.e. reads with PQs below the threshold t_i and TN_i - true negatives, i.e. reads without
454 PQs below the threshold t_i . The false positive rate, $FPR_i = (FP_i / (FP_i + TN_i))$ was calculated for each
455 threshold t_i and the thresholds for OQ detection were set in order to have $FPR \approx 0.02$ (high
456 specificity), i.e. 2% of the non PQs would be detected as OQs. This yielded thresholds of 18%
457 and 25 for K⁺ and PDS sequencing respectively. A sequence with a % mismatch value above
458 these thresholds was defined as an Observed G-quadruplex Sequence (OQs). For the genomic
459 analysis, continuous regions with a maximal peak summit above the threshold (18% for K⁺ and
460 25% for PDS) were considered as OQ regions. OQ regions displaying multiple peak were split
461 into separated OQs using PeakSplitter (<http://www.ebi.ac.uk/research/bertone/software>). Regions
462 from two replicates were analysed independently, keeping strand information separated. We only
463 considered high confidence OQ regions in genomic intervals common to both replicates for
464 further analyses (command intersectBed -s of the bedtools).

465 **Structural analysis of OQ categories**

466 OQ sequences were stratified into different OQ categories by searching for different regular
467 expressions (Fig. 3). To assign univocally an OQ region to a specified category and avoid
468 considering the same region multiple times, we followed priority rules based on the predicted
469 stability from high to low (Loop 1-3 > Loop 4-5 > Loop 6-7 > Long loops > Bulges > Other). The
470 different categories were defined as follows: Loop 1-3: (G₃N_{1,3})₃G₃, with N =
471 [ATCG]; Loop 4-5: (G₃N_{1,5})₃G₃ and not in previous category; Loop 6-7:
472 (G₃N_{1,7})₃G₃ and not in a previous category; Long loops: (G₃N_{1,12})₃G₃,
473 or G₃N_{1,7}G₃N_{13,21}G₃N_{1,7}G₃ and not in a previous category; Bulges: OQ
474 sequences with any G-run being GH₁₋₇GG or GHGGN_{1,7}GGHG, with H = [ATC] and not in a
475 previous category; Other: not in any other category. The other category was further stratified into
476 sub-categories containing OQs having either multiple bulges with more than one nucleotide (e.g.,
477 GH_{2,5}GGN_{1,7}GGH_{2,5}G) or two-tetrads motifs (GGN_{1,7}GGN_{1,7}GGN_{1,7}GG)
478 (Extended Data Table 2). Finally, the ratio of the numbers of each category in PDS and K⁺ was
479 calculated (Extended Data Fig. 5).

480 **Fold-enrichment analysis of OQ structural categories**

481 The 525,890 K⁺ OQ intervals were randomly shuffled three times across the genome (command
482 shuffleBed in bedtools) to generate random sequences of the same size distribution as the OQs.
483 This was also done for the 716,310 PDS OQ intervals. The different OQ categories were
484 identified and counted in both the experimental OQs and the three randomized intervals. For each
485 category, the ratio of real OQ over the average of three random cases was calculated and plotted
486 as fold-enrichment for PDS and K⁺ (Fig. 4b). Error bars were calculated for each category as the
487 standard error of the mean (SEM) of three random replicates, and each SEM was then divided by
488 the average of random counts in the category to adapt it to the fold enrichment plot.

489 **Genomic regions analysis**

490 Gene annotation files were downloaded from the UCSC genome browser website
491 (<https://genome.ucsc.edu/>), genome version *hg19*, and different genomic regions (5'UTRs,
492 3'UTRs, exons, introns, promoters, TSSs and splice regions) were extracted and stored as
493 genomic intervals (bed file format). For each region, the total number of regions, the total region
494 size and the number of PDS or K⁺ OQs overlapping to the region intervals (command
495 intersectBed of the bedtools) were calculated. The number of regions overlapping exclusively
496 with *Quadparser* PQs and with non-canonical PQs (i.e., Long loops and Bulges) were calculated
497 (Extended Data Table 3). Any intervals overlapping sequences from both categories were
498 excluded from analysis to avoid ambiguity.

499 **Genes and oncogenes analysis**

500 For each gene annotated in the version *hg19* of the human genome, the number of *Quadparser*
501 predicted PQs, of OQs in PDS and OQs in K⁺ were counted. The density of PQs or OQs was
502 calculated by dividing the respective counts by the gene body length and multiplying by 1000
503 (i.e. density is the number of structures per kilobase). For oncogene analyses, we considered 498
504 oncogenes and 766 tumour suppressors²⁴. Genes with a PQs density less than half of *SRC* PQs
505 density but with a OQs density higher than *SRC* OQs density were extracted (Supplementary
506 Table 4).

507 **Somatic copy number alteration (SCNA) analysis**

508 140 SCNAs previously identified as being associated with cancer were considered¹⁰, of which 70
509 were amplifications and 70 were deletions. Only SCNA less than 10 Mb in size were analysed,
510 leaving a total of 123 regions (50 deletions and 73 amplifications). For each region t

511 he number of OQs was counted. OQ genomic intervals were then randomly reshuffled three times
512 (random-OQs) and the number of random-OQs in each SCNA was calculated and averaged. The
513 OQs and random-OQs counts were divided by each region size and multiplied by 1000, to give a
514 density per kilobase. The OQs and random-OQs densities were then compared and their ratio
515 calculated such that SCNA regions with ratio > 1 are enriched in OQs compared to random,
516 whereas SCNAs with ratio < 1 are depleted (Supplementary Table 5; Fig. 4b). The difference
517 between OQs and random densities was statistically assessed for the 123 regions using the two-
518 tailed t-test; SNCA amplifications (n=73) and deletions (n=50) were also tested in the same way
519 against their counterpart (random-OQs for amplification and deletion regions only, respectively).

520

521 31. Cock, P. J., Fields, C. J., Goto, N., Heuer, M. L. & Rice, P. M. The Sanger FASTQ file
522 format for sequences with quality scores, and the Solexa/Illumina FASTQ variants.
523 *Nucleic Acids Res.* **38**, 1767-1771, (2010).

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547 **Extended Data Figure 1: Overall sequencing quality in sequencing experiments with the**
548 **different cations Li⁺, Na⁺ and K⁺.** Each plot visually shows base calling quality (Phred quality
549 score, Q; y-axes) for the 250 sequenced bases (x-axes), in two independent experiments, with
550 sequencing buffers containing Li⁺ (top), Na⁺ (middle) and K⁺ (bottom), as generated by the
551 program FastQC (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>). Yellow bars and
552 black whiskers are box plots for the respective base positions; red lines are median values; blue
553 lines are mean values.

554

555 **Extended Data Figure 2: Sequencing quality and sequencing errors (% mismatches) for**
556 **control sequences.** Bar plots showing the average Phred quality score (Q) (a) and % mismatches
557 (b) for the 4 control sequences when sequencing with different cations Li⁺ (left), Na⁺ (middle)
558 and K⁺ (right). *c-kit*, *c-myc*: positive controls; *c-myc-opp*, *c-myc-mut*: negative controls (see
559 Methods). Data is taken from a number of independent sequencing experiments: 3 for Na⁺, 2 for
560 Li⁺ and K⁺. The numbers of different control sequences (i.e. independent sequencing clusters on
561 the flow cell) in the combined experiments are (order; *c-kit*, *c-myc-opp*, *c-myc*, *c-myc-mut*): 2741,
562 1139, 1040, 10945 for Li⁺; 8235, 3076, 2787, 26974 for Na⁺; 2935, 1315, 1, 12809 for K⁺. Bars
563 are standard deviations. No error bar present for *c-myc* in K⁺ (n=1).

564 **Extended Data Figure 3: Sequencing errors for controls in PDS conditions.** % mismatches
565 for the control sequences in the same sequencing experiment with Na⁺ sequencing buffers during
566 the first read (Read-1; left) followed by the addition of the small-molecule PDS in Na⁺ throughout
567 the second read (Read-2; right). Error bars are SEMs (respectively: 0.16, 0.02, 0.18 and 0.07 for
568 left plot; 0.12, 0.08, 0.15 and 0.09 for right plot). N = 948, 367, 367 and 3990 for *c-kit*, *c-myc-*
569 *opp*, *c-myc*, *c-myc-mut*.

570 **Extended Data Figure 4: Different families of G-quadruplex structures:** Left: canonical PQs
571 predicted by Quadparser (L1-3=N1-7, with N=A|C|T|G). Middle: PQs with longer loops (L1-
572 3=N8-12 or L2=N8-21). Right: PQs with a single bulge B1=H1-7 or multiple bulges B2=H1-5
573 (H=A|T|C).

574

575 **Extended Data Figure 5: Detection of OQs representing different G-quadruplex structural**
576 **families in PDS versus K⁺ conditions.** Fold enrichment (ratio) between the numbers of OQs in
577 PDS over K⁺ for each category (see B). Values > 1 indicate higher numbers in PDS. G-
578 quadruplex structural families: Loop 1-3; Loop 4-5; Loop 6-7: OQs with at least one loop of the

579 indicated length; Long loops: OQs with any loop of length 8 to 12 for L1-3 or 8 to 21 for L2;
580 Bulges: OQs with one bulge of 1 to 7 bases (A, T, C) or multiple bulges of 1 base.

581 **Extended Data Figure 6: Comparison of genomic regions in PDS and K⁺ sequencing**
582 **conditions.** a) Genome browser view of a genomic region within *MYC* oncogene. Red and orange
583 tracks: % mismatches in reads aligning to the reverse strand (-) for PDS and K⁺, respectively. OQ
584 intervals are shown as red and orange bars below the corresponding peaks.. b) Genome browser
585 view of a genomic region within the *MYL5-MFSD7* gene. Black and blue tracks: % mismatches
586 in reads aligning to the forward strand (+) for PDS and K⁺, respectively. OQ intervals are shown
587 as black and blue bars below the corresponding peaks. c) Genome browser view of a genomic
588 region within the *MYL9* gene. All colours and features as in a). See Supplementary Table 2 for
589 sequence details. For all panels, OQs not predicted by Quadparser are indicated by * and
590 Quadparser PQs are shown as black bars.

591 **Extended Data Figure 7: Comparison of forward versus reverse strands in PDS sequencing**
592 **conditions.** A) Genomic region within the *MYL9* gene. Red and black tracks: % mismatches in
593 reads aligning to the reverse strand (-) and forward strand (+), respectively. OQs intervals are
594 shown as red and black bars below corresponding peaks. Quadparser PQs are shown below in
595 black. OQs not predicted by Quadparser are indicated by asterisks (*). See Supplementary Table
596 2 for sequence details.

597 **Extended Data Table 1: Quadparser PQs detected by G4-Seq.** The number and percentage of
598 Quadparser PQs detected by G4-Seq under PDS or K⁺ conditions or common to both. Two
599 replicate Illumina HiSeq sequencing runs were performed for each condition. These data show
600 the high degree of reproducibility and overlap between two different G-quadruplex stabilisation
601 conditions.

602 **Extended Data Table 2: Number of OQs in the category “Other”.** Multiple bulges =
603 G[ATC]₂₋₅GGLGG[ATC]₂₋₅G; two-tetrads = GGLGGLGGLGG, with L = N₁₋₇; % G content =
604 percentage of G nucleotides in the OQ sequence.

605

606

607 **Extended Data Table 3: Distribution of OQs in different genomic regions.** Several
608 measurements reporting the genomic distribution of OQs in PDS (top half) or K⁺ (bottom half)
609 are listed in the table; columns are as follows. Region: genomic features- UTR: untranslated
610 region; TSS: transcription start site; promoters 1000 up: 1000 bases upstream the TSS; TSS 1000
611 up down: 1000 bases up- and down-stream the TSS; splice 50: 50 bases up- and down-stream
612 splice sites (i.e. exon-intron junctions); # regions: number of disjoint genomic regions; total
613 region size: sum of all disjoint genomic regions; OQs density: $1000 * (\# \text{ OQs}) / (\text{total region}$
614 $\text{size})$; # regions with OQs: number of genomic regions overlapping with at least one OQ. #
615 regions with non-canonical OQs: number of genomic regions overlapping exclusively with OQs
616 having a long loop or a bulge; # regions with PQs: number of genomic regions overlapping
617 exclusively with *Quadparser*-predicted PQs (loop size 1-7). Ratio non-canonical OQs / PQs: ratio
618 of the number of regions with non-canonical and canonical PQs.