

1 "This is an un-copied authored manuscript copyrighted by the  
2 American Association for Clinical Chemistry (AACC). This may not be  
3 duplicated or reproduced, other than for personal use or within the rule of  
4 'Fair Use of Copyrighted Materials' (section 107, Title 17, U.S. Code)  
5 without permission of the copyright owner, AACC. The AACC disclaims  
6 any responsibility or liability for errors or omissions in this version of the  
7 manuscript or in any version derived from it by the National Institutes of  
8 Health or other parties. The final publisher-authenticated version of the  
9 article is available at <http://www.clinchem.org>."

10

11 **Ultra-sensitive mutation detection and genome-wide DNA copy**  
12 **number reconstruction by error corrected circulating tumor DNA**  
13 **sequencing**

14

15 Sonia Mansukhani<sup>1\*</sup>, Louise J. Barber<sup>1\*</sup>, Dimitrios Kleftogiannis<sup>1</sup>, Sing Yu  
16 Moorcraft<sup>2</sup>, Michael Davidson<sup>2</sup>, Andrew Woolston<sup>1</sup>, Paula Zuzanna Proszek<sup>3</sup>,  
17 Beatrice Griffiths<sup>1</sup>, Kerry Fenwick<sup>4</sup>, Bram Herman<sup>5</sup>, Nik Matthews<sup>4</sup>, Ben O'Leary<sup>6</sup>,  
18 Sanna Hulkki<sup>3</sup>, David Gonzalez De Castro<sup>7</sup>, Anisha Patel<sup>8</sup>, Andrew Wotherspoon<sup>9</sup>,  
19 Aleruchi Okachi<sup>2</sup>, Isma Rana<sup>2</sup>, Ruwaida Begum<sup>2</sup>, Matthew N. Davies<sup>1,10</sup>, Thomas  
20 Powles<sup>11</sup>, Katharina von Loga<sup>1</sup>, Michael Hubank<sup>3</sup>, Nick Turner<sup>6,12</sup>, David Watkins<sup>2</sup>,  
21 Ian Chau<sup>2</sup>, David Cunningham<sup>2</sup>, Stefano Lise<sup>1</sup>, Naureen Starling<sup>2</sup> and Marco  
22 Gerlinger<sup>1,2</sup>

23

24 **Running title: Error corrected circulating tumor DNA sequencing**

25

26 **Affiliations:**

27 <sup>1</sup>Centre for Evolution and Cancer, Division of Molecular Pathology, The Institute of  
28 Cancer Research, London, United Kingdom.

29 <sup>2</sup>Gastrointestinal Cancer Unit, The Royal Marsden NHS Foundation Trust, London  
30 and Sutton, United Kingdom.

31 <sup>3</sup>Centre for Molecular Pathology, The Royal Marsden NHS Foundation Trust, Sutton,  
32 United Kingdom.

33 <sup>4</sup>Tumour Profiling Unit, The Institute of Cancer Research, London, United Kingdom.

34 <sup>5</sup>Diagnostics and Genomics Group, Agilent Technologies Inc., Santa Clara, USA

35 <sup>6</sup>Breast Cancer Now Research Centre, The Institute of Cancer Research, London,  
36 United Kingdom.<sup>7</sup>Centre for Cancer Research and Cell Biology, Belfast, United  
37 Kingdom.

38 <sup>8</sup>Department for Radiology, The Royal Marsden NHS Foundation Trust, London and  
39 Sutton, United Kingdom.

40 <sup>9</sup>Department of Histopathology, The Royal Marsden NHS Foundation Trust, London  
41 and Sutton, United Kingdom.

42 <sup>10</sup>Current address: Achilles Therapeutics, Francis Crick Institute, London, United  
43 Kingdom.

44 <sup>11</sup>Barts Cancer Institute, Queen Mary University of London, London, United  
45 Kingdom.

46 <sup>12</sup>Breast Cancer Unit, The Royal Marsden NHS Foundation Trust

47 \*These authors contributed equally to this work

48

49 **Corresponding Author:**

50 Dr Marco Gerlinger

51 Centre for Evolution and Cancer, The Institute of Cancer Research

52 237 Fulham Road, London SW3 6JB, United Kingdom

53 Tel: +44 207 153 5234

54 email: marco.gerlinger@icr.ac.uk

55

56 **Keywords:** cancer genomics, circulating tumor DNA, liquid biopsy, molecular  
57 barcodes, sequencing error correction.

58

### 59 **Abstract**

60 **Background:** Circulating free DNA sequencing (cfDNA-Seq) can portray cancer  
61 genome landscapes but highly sensitive and specific technologies are necessary to  
62 accurately detect mutations with often low variant frequencies.

63 **Methods:** We developed a customizable hybrid-capture cfDNA-Seq technology  
64 using off-the-shelf molecular barcodes and a novel duplex DNA-molecule  
65 identification tool for enhanced error correction.

66 **Results:** Modelling based on cfDNA-yields from 58 patients showed this technology,  
67 requiring 25 ng cfDNA, could be applied to >95% of patients with metastatic  
68 colorectal cancer (mCRC). cfDNA-Seq of a 32-gene/163.3kbp target region detected  
69 100% of single nucleotide variants with 0.15% variant frequency in spike-in  
70 experiments. Molecular barcode error correction reduced false positive mutation  
71 calls by 97.5%. In 28 consecutively analyzed patients with mCRC, 80 out of 91  
72 mutations previously detected by tumor tissue sequencing were called in the cfDNA.  
73 Call rates were similar for point mutations and indels. cfDNA-Seq identified typical  
74 mCRC driver mutations in patients where biopsy sequencing had failed or did not  
75 include key mCRC driver genes. Mutations only called in cfDNA but undetectable in  
76 matched biopsies included a subclonal resistance driver mutation to anti-EGFR

77 antibodies in *KRAS*, parallel evolution of multiple *PIK3CA* mutations in two cases,  
78 and *TP53* mutations originating from clonal hematopoiesis. Furthermore, cfDNA-Seq  
79 off-target read analysis allowed simultaneous genome-wide copy number profile  
80 reconstruction in 20 of 28 cases. Copy number profiles were validated by low-  
81 coverage whole genome sequencing.

82 **Conclusions:** This error-corrected ultra-deep cfDNA-Seq technology with a  
83 customizable target region and publicly available bioinformatics tools enables broad  
84 insights into cancer genomes and evolution.

## 85            **Introduction**

86            Many tumors release cell free DNA (cfDNA) into the circulation, allowing the  
87 analysis of cancer genetic aberrations from blood samples [1-6]. Such 'liquid  
88 biopsies' can inform tailored therapies [7] or predict recurrences after surgery [8, 9].  
89 cfDNA analysis also permits subclonal mutation detection that is often missed by  
90 biopsies due to spatial intratumor heterogeneity [10, 11]. Genetic techniques with  
91 high analytical sensitivity and low false positive error rates are crucial for accurate  
92 cfDNA-Seq due to low tumor-derived cfDNA fractions and low abundances of  
93 subclonal mutations. Digital droplet PCR (ddPCR) and BEAMing assays can  
94 accurately detect point mutations present at frequencies  $\leq 0.1\%$  but are restricted to  
95 the analysis of a small number of genomic loci [8, 12]. Targeted next generation  
96 sequencing (NGS) can interrogate larger regions such as gene panels but the error  
97 rate of NGS complicates the calling of mutations with variant allele frequencies  
98 (VAFs)  $< 5\%$  [13]. Error correction through random molecular barcodes (MBC) has  
99 been incorporated into NGS cfDNA assays to reduce this error rate [14, 15] and has  
100 enabled mutation calling with VAFs  $\leq 0.1\%$ . However, these methods have often  
101 used amplicon sequencing, which can hamper coverage of entire genes due to  
102 primer design restrictions. Some methods have employed solution hybrid-capture,  
103 which is ideal to target entire genes, but used bespoke or proprietary rather than off-  
104 the-shelf reagents and publicly available bioinformatics tools, limiting their broad  
105 application for clinical or research purposes.

106            Here we assessed how novel, commercially available off-the-shelf MBC  
107 reagents combined with customized capture-based target enrichment technology  
108 could be optimized for ultra-deep error-corrected cfDNA-Seq. We developed a

109 duplex-DNA molecule calling tool to improve the calling accuracy and assessed  
110 concordance of mutation calls from cfDNA with clinical grade tumor tissue  
111 sequencing in patients with metastatic colorectal cancer (mCRC).

## 112 **METHODS**

### 113 **Patients and samples**

114 Plasma samples and clinical data were available from the FOrMAT trial  
115 (Feasibility Of Molecular characterization Approach to Treatment [16], Chief  
116 Investigator: N Starling ClinicalTrials.gov NCT02112357). Healthy donor (HD) cfDNA  
117 was obtained through the Tissue Collection Framework to Improve Outcomes in  
118 Solid Tumours (Chief Investigator: T Powles). Both trials were approved by UK  
119 ethics committees and all patients provided written informed consent. Details of  
120 clinical trials, patients, samples, sample processing and experimental techniques are  
121 provided in the online Supplemental Methods file.

122

### 123 **cfDNA sequencing**

124 SureSelect<sup>XT-HS</sup> (Agilent) was used to prepare sequencing libraries using our  
125 optimized protocol (online Supplemental Methods file) and a custom designed  
126 SureSelect bait-library (online Supplemental Table 1). Sequencing libraries were  
127 clustered using the cBot and sequenced with paired-end 75 reads on an Illumina  
128 HiSeq2500 in rapid mode.

129 SureCall software (version 4.0.1.45, Agilent) was used to trim and align fastq  
130 reads to the hg19 reference genome with default parameters and for MBC de-  
131 duplication, permitting one base mismatch within each MBC. Consensus families

132 comprising of single reads were removed, on-target depths were assessed and  
133 variants were called with SureCall SNPPEP.

134 To identify variants supported by duplexes we developed the freely available  
135 duplexCaller bioinformatics tool [17].

136 All variant positions identified in patient cfDNA were assessed in six HD  
137 samples using bam-readcount [18]. Most called variants were absent in HD samples  
138 (online Supplemental Table 2) but mutations with VAF less than double that of an  
139 identical variant in HD were removed as false positives.

140 BAM files resulting from MBC de-duplication before removal of single-read  
141 consensus families were used to generate genome-wide DNA copy number profiles  
142 with CNVkit [19], with Antitarget average size set to 30 kb. HD samples were used  
143 as the normal reference pooled dataset.

144

#### 145 **Low coverage whole genome sequencing (lcWGS)**

146 Genomic libraries were constructed from 10 ng cfDNA with the NEBNext Ultra  
147 II kit and sequenced with 100bp single-end reads on HiSeq2500 in rapid mode  
148 (0.42x median coverage). Data was aligned (hg19 reference genome) with Bowtie  
149 (v0.12.9), and processed as described [20]. logRatios were normalised against a  
150 gender-matched pooled dataset from HD cfDNA (9 male, 8 female) before  
151 segmentation and median centering.

152

#### 153 **ddPCR**

154 ddPCR was performed in case 8 (BRAF V600E) and to validate discordant  
155 variants between cfDNA and tumor tissue. 4 of 11 such cases had sufficient  
156 remaining cfDNA to validate subclonal variants (online Supplemental Methods file).

157

### 158 **cfDNA sequencing with a commercial kit**

159 25 ng, 17 ng and 25 ng cfDNA (cases 3, 15, and 23, respectively) were  
160 processed using the Roche AVENIO Expanded kit as per the manufacturer's  
161 protocol. Libraries were sequenced with 151 bp paired-end reads on Illumina  
162 NextSeq500 to 2,689-6,420x depth after de-duplication. Data was analyzed using  
163 the Roche AVENIO ctDNA Analysis Software v1.0.0 with default parameters.

164

## 165 **RESULTS**

### 166 **cfDNA sequencing optimization**

167 Modelling based on cfDNA yields from 58 patients with mCRC showed that 25  
168 ng of cfDNA could be extracted from 20-30 ml blood from >95% of cases (online  
169 Supplemental Figure 2C). 25 ng was therefore chosen as our standard cfDNA input  
170 quantity. We designed a solution hybrid-capture panel targeting 32 genes including  
171 all major CRC driver genes, (163.3 kb, online Supplemental Table 1) and used  
172 Agilent SureSelect<sup>XT-HS</sup> kit, which tags each DNA strand with a random 10-base  
173 MBC, for sequencing library preparation. The SureSelect<sup>XT-HS</sup> protocol was optimized  
174 to perform reliably with 25 ng cfDNA input (online Supplemental Methods file). The  
175 fraction of on-target reads is usually low when using small targeted sequencing  
176 panels and low input DNA, so we first assessed how the on-target fraction could be  
177 optimized by varying the stringency of the post-capture wash. Two library



178 preparations were started in parallel from each of four cfDNA samples, using the 1.5  
179 h fast-hybridization protocol. Then, post-capture washes were performed at 65°C in  
180 one library and at 70°C in the other. Sequencing generated similar read numbers  
181 (65°C: 92,820,887; 70°C: 102,582,694 median reads/sample) and the on-target  
182 fraction significantly ( $p=0.0011$ ) increased from 30-35% to 71-74% with the 70°C  
183 protocol (Figure 1A). Hence, the more stringent conditions were chosen for our  
184 standard protocol. Target exon coverage was even with this solution hybrid-capture  
185 technique and was not subject to the gaps commonly seen with commercial  
186 amplicon sequencing designs (online Supplemental Figure 3). This would be  
187 particularly advantageous for the analysis of tumor suppressor genes where driver  
188 mutations often spread across large parts of the gene.

189 We next used MBCs to de-duplicate sequencing data and perform error  
190 correction. SureCall creates families of reads with matched MBC that also align to  
191 the same genomic position and then identifies the most likely consensus sequence  
192 for each family (Figure 1B). This reduces random errors arising during PCR and  
193 sequencing, as these are not common to all reads of a family. Consensus families  
194 contained a median of 8 to 15 supporting reads in samples sequenced with the  
195 optimized protocol (online Supplemental Figure 4), which was within the optimal  
196 range for barcode error correction [21]. After MBC de-duplication, the median on-  
197 target depth with the 70°C protocol was 1,782x. This was theoretically sufficient to  
198 achieve a detection limit as low as 1 mutated DNA fragment in 1,782 molecules  
199 (0.056%). However, the analytical sensitivity for *de novo* mutation detection is lower  
200 in practice since more than one read is required to support robust bioinformatics  
201 calling. Thus, we designed a mixing experiment to test the ability to detect and  
202 bioinformatically call mutations with low VAFs.

203 **Assay sensitivity and specificity**

204 cfDNA from two donors that differed in 16 homozygous single nucleotide  
205 polymorphisms (SNPs) within the targeted region were used to prepare a dilution  
206 series with 0.15%, 0.075% and 0.0375% cfDNA from donor A spiked into cfDNA  
207 from donor B. Sequencing a median of 74,030,118 reads/sample generated a  
208 median on-target depth of 21,651x before de-duplication. Data from each sample  
209 was then processed in two ways: first, we used MBCs for de-duplication and calling  
210 of consensus sequences; second, we performed standard de-duplication using only  
211 the genomic position of each read pair. The median on-target depth was higher after  
212 MBC de-duplication (MBC 2,420x versus 1,587x with standard de-duplication; Figure  
213 1C). This was anticipated as different MBCs tag distinct DNA fragments that would  
214 otherwise be counted as duplicates. For example, the forward and reverse strands of  
215 each original 'duplex' dsDNA molecule were separately tagged by MBC and so were  
216 retained as independent consensus families (Figure 1D). Standard de-duplication  
217 cannot distinguish these reads from PCR duplicates.

218 We first investigated whether the spiked-in SNPs could be re-identified in the  
219 MBC de-duplicated BAM files using the Integrative Genomics Viewer (IGV) [22] and  
220 tried to understand patterns associated with true positive variants. All 16 SNPs were  
221 detected in the 0.15% mix, 14/16 at 0.075% and 11/16 at 0.0375% mixing ratios  
222 (Figure 1E). Thus, our ultra-deep cfDNA-Seq assay allowed robust detection of  
223 variants at 0.15% and retained a high detection capability at 0.075%. We then  
224 assessed if MBC error correction improved the bioinformatics calling accuracy of  
225 ultra-low frequency variants, which is more challenging than re-identification of  
226 known variants. While interrogating sequencing data manually in IGV, we had  
227 observed that all true variants were at least supported by two consensus families

228 mapping to the same genomic position but differing in whether the variant was seen  
229 in read 1 or read 2 in paired-end sequencing (Figure 1D). These were highly likely to  
230 represent the forward and reverse strand of the double-stranded input cfDNA  
231 molecule as observed previously [15]. Based on this observation, we developed the  
232 duplexCaller bioinformatics tool that identified variants supported by duplex reads  
233 (online Supplemental Methods file) and added the requirement for such a 'duplex-  
234 configuration' to be present to accept a mutation as genuine. The presence of a  
235 variant in at least one additional family with a different alignment position was also  
236 added to the post-call filters to assure high specificity. Thus, a variant had to be  
237 present in  $\geq 3$  consensus DNA families in order to be accepted as a mutation call in  
238 the MBC de-duplicated data. For a meaningful comparison, mutations in the  
239 standard de-duplicated data were also required to be present in  $\geq 3$  reads.

240 We then compared SureCall calls for the mixing experiment on standard-  
241 versus MBC-deduplicated data and quantified how many of the homozygous SNPs  
242 from sample A that were present at 0.15% in the cfDNA mixture were called.  
243 Although samples A and B differed at 16 homozygous SNP positions, only the 9  
244 variant SNPs in spiked-in sample A could be assessed for capability to call at low  
245 frequency against the reference genome. The other 7 SNPs were reference wild-  
246 type in spiked-in sample A and so could not be called. Mutation calling after standard  
247 de-duplication with low stringency caller settings (variant call quality threshold  
248 [VCQT]=40) detected 5/9 homozygous SNPs (Figure 1F) but also generated 156  
249 additional calls. These additional variants were likely false positives, since they had  
250 not been identified by deep sequencing of the individual cfDNA samples used in the  
251 mixing experiment. Stepwise increase of the VCQT reduced false positives but this  
252 was accompanied by a loss of analytical sensitivity. When the same data were called

253 using MBCs and a low stringency VCQT=40 (Figure 1F), 4 of the spiked-in SNPs  
254 were called with only 2 likely false positive variants. We assessed why calling with  
255 MBC error correction failed to identify the 5 other SNPs. Each of these had VAFs  
256 <0.1% when visualized in IGV [23], which was below the minimum VAF of 0.1% that  
257 can be called by SureCall. We also assessed the number of false positive calls in  
258 standard de-duplicated data at the maximum VCQT that identified the same four true  
259 positive variants detected with MBC: 81 likely false positives were called compared  
260 to just 2 using MBC. Hence at the same analytical sensitivity, de-duplication using  
261 the MBCs dramatically decreased false positives by 97.5%. Mutation calling in 6 HD  
262 samples subjected to cfDNA-Seq only identified heterozygous and homozygous  
263 SNPs but no mutations with lower frequency (online Supplemental Table 3), further  
264 supporting the high analytical specificity of this MBC technology.

### 265 **Concordance of cfDNA- and tumor-sequencing in mCRC patients**

266 cfDNA from 28 patients with mCRC were consecutively analyzed. Seven were  
267 sequenced with the 65°C protocol and 21 with the 70°C protocol. The median  
268 sequencing depth was higher with 70°C (2,087x) than 65°C (1,205x) (Figure 2A).

269 We then analyzed the concordance and discordance of mutation calls within  
270 the target regions common to the tumor biopsy sequencing assay and our cfDNA-  
271 Seq panel. Biopsies of 23 cases had been sequenced with the FOrMAT NGS panel  
272 (online Supplemental Table 4) and four biopsies had been subjected to routine  
273 clinical amplicon sequencing of 5 genes (*BRAF*, *KRAS*, *NRAS*, *PIK3CA* and *TP53*).  
274 One case had failed tissue sequencing.

275 88% (80/91) of all mutations that had been found by tumor sequencing were  
276 called in the cfDNA (Figure 2A). All 11 mutations not called in cfDNA were from only

277 3 cases. Inspection of the sequencing data on IGV revealed that 5/11 mutations  
278 were present in cfDNA at VAFs below the SureCall detection limit (Figure 2B).  
279 Sufficient cfDNA remained from case 8 for orthogonal analysis by ddPCR. Using  
280 manufacturer-validated ddPCR-probes for the *BRAF* V600E mutation we identified  
281 2,830 wild type DNA fragments but no mutated fragments (data not shown). This  
282 confirmed that the absence of sufficiently abundant tumor-derived cfDNA molecules,  
283 rather than technical failure, explained the inability to detect mutations.

284 We next assessed mutations called by cfDNA-Seq in genes that had not been  
285 sequenced in corresponding tumor tissue. *APC* mutations were detected in each of 4  
286 cases whose tumors had only been analyzed with the 5-gene amplicon panel (Figure  
287 2A). Furthermore, one mutation was found in each of *FBXW7*, *CTNNB1*, *TCF7L2*,  
288 *ATM* and *SMAD4*. We also detected mutations in *APC*, *TP53* and *KRAS* in case 28  
289 that had failed prior tumor tissue sequencing attempts. In total, 11 of these 13  
290 mutations (85%) encoded protein changes previously reported in the COSMIC  
291 cancer mutation database [24] and all variants in the tumour suppressor genes *APC*  
292 and *FBXW7* were truncating and hence likely driver mutations. This demonstrated  
293 that our assay could detect biologically and clinically important cancer mutations  
294 directly from cfDNA.

295 We then investigated mutations that had been called in cfDNA but were  
296 absent when the same gene had been analyzed in tumor tissue: 7 in *TP53*, 7 in  
297 *ATM*, 3 in *PIK3CA*, 2 in *SMAD4* and one each in *KRAS*, *FBXW7* and *TCF7L2*. All  
298 four mutations called in the oncogenes *KRAS* and *PIK3CA* were canonical cancer  
299 driver mutations. 8/18 mutations (44%) located in tumor suppressor genes were  
300 nonsense mutations or encoded for amino acid changes found recurrently in cancer  
301 [24], suggesting that these were also driver mutations. Together, 54.5% (12/22) of

302 variants detected only in cfDNA were likely cancer driver mutations. The VAFs of  
303 mutations that were only detected in cfDNA but not in tumor tissue were a mean  
304 105-fold lower than the VAF of the most abundant mutation detected in the same  
305 cfDNA sample (online Supplemental Figure 1); these variants likely originated from  
306 small cancer subclones. However, two *TP53* mutations present in cfDNA but not in  
307 matched tumor tissue (Cases 9, 13) were also detected with similar VAF in DNA  
308 from blood cells (online Supplemental Table 5). These *TP53* mutations hence  
309 originated from a clonal expansion of blood cells [9], termed clonal hematopoiesis  
310 [25, 26].

311 An activating mutation in *KRAS* (Q61H) was detected with a VAF of 0.37% in  
312 cfDNA but not in the matched tumor (case 10). This was the only patient that had  
313 received treatment with the anti-EGFR antibody cetuximab prior to blood collection  
314 and the *KRAS* mutation was likely a driver of acquired resistance that evolved during  
315 therapy [27]. ddPCR testing of cfDNA provided orthogonal validation (Figure 2C),  
316 showing that our technology is suitable for the detection of subclonal resistance  
317 driver mutations. Suspected driver mutations in *PIK3CA* were frequently discordant  
318 with 3/7 mutations only detectable in cfDNA (E545K, H1046R, R1023\*). Two cases  
319 (17,26) harbored parallel evolution events, as further activating *PIK3CA* mutations  
320 were present in the tumors and the cfDNA. These results are consistent with studies  
321 showing that intratumor heterogeneity of *PIK3CA* mutations is common in mCRCs  
322 whereas heterogeneity is rare for mutations in *APC* and, in tumors not previously  
323 treated with anti-EGFR antibodies, for *KRAS*, *NRAS* and *BRAF* mutations [28].

324 Mutations in *ATM* tumor suppressor gene were called in 8/28 cfDNA samples.  
325 Sequencing of matched tumor showed wild-type sequence in seven of these and  
326 one tumor had only been sequenced with the 5-gene panel. All *ATM* mutations had

327 low VAFs (median: 0.17%) and only 2/8 encode protein changes previously  
328 catalogued in cancer [24], making it difficult to interpret their functional relevance. No  
329 *ATM* mutations were called in 6 healthy donors, indicating that the mutation calls in  
330 cfDNA from mCRC patients are unlikely the result of a high false positive call rate in  
331 this gene.

332         Next, we used ddPCR to validate further subclonal mutations called in cfDNA  
333 but not in tumor tissue. All subclonal variants with VAF <2% from samples where  
334 sufficient cfDNA material was available and where a custom ddPCR-assay could be  
335 designed were assessed (online Supplemental Methods file). ddPCR validated all 6  
336 tested mutations and VAFs were similar to those found by our error-corrected cfDNA  
337 technology (Figure 2D, online Supplemental Table 6).

338         Additionally, we re-sequenced three cfDNA samples containing low VAF  
339 (<2%) mutations (cases 3, 15, 23) with the commercially available AVENIO ctDNA  
340 kit. 9/10 point mutations in genes targeted by both panels were concordant (online  
341 Supplemental Table 7). The low frequency *TP53* R175H variant in case 3 was not  
342 called by AVENIO software but was seen to be present upon manual review of the  
343 BAM file. Three indels in *APC* (cases 3,23) were not called by AVENIO analysis.  
344 This comparison further confirmed the reliable performance of our customizable  
345 cfDNA assay.

#### 346         **Genome wide DNA copy number aberration analysis**

347         We finally assessed if we could maximise the information gain from a targeted  
348 cfDNA assay through simultaneous reconstruction of genome-wide copy number  
349 aberration (CNA) profiles. Applying the CNVkit-package [19] that uses off-target  
350 reads to infer copy number changes, we generated genome-wide CNA profiles for

351 20/28 cases (71%) (Figure 3A-B). Chromosome arm losses (Chr17p and 18q) and  
352 gains (Chr1q, 7, 8q, 13 and 20), which are typical for mCRC, were observed [29]. All  
353 8 samples with a flat CNA profile had very low maximum VAFs  $\leq 5.6\%$ . A high-level  
354 targetable amplification involving the *ERBB2* oncogene was detected despite a low  
355 tumor-derived cfDNA fraction (8.6% VAF) in case 11 (Figure 3C). This amplification  
356 had also been detected in the matched tumor, validating the ability to profile CNAs  
357 with our cfDNA-Seq technology. No other amplifications had been detected in tumor  
358 biopsies with the FOrMAT NGS panel. Low-coverage whole genome sequencing is  
359 an established approach for genome wide copy number profiling and we applied this  
360 to 18 samples with sufficient cfDNA. This independent validation showed a median  
361 weighted Spearman correlation of 0.886 with the profiles generated from cfDNA-Seq  
362 using CNVkit (online Supplemental Figure 5).

## 363 **DISCUSSION**

364 Our ultra-deep and error-corrected cfDNA-Seq protocol that uses off-the-shelf  
365 MBCs in combination with a custom-designed solution hybrid capture panel detected  
366 100% of the known variants with VAFs of 0.15% in a mixing experiment. The use of  
367 MBC error correction and the requirement for variants to be supported by a duplex-  
368 pair of consensus families reduced false positive mutation calls by 97.5% while  
369 maintaining true positives. We developed the DuplexCaller bioinformatics tool, which  
370 can be run directly after MBC de-duplication to facilitate mutation calling; all  
371 bioinformatics tools for the analysis of data generated with this technology are hence  
372 freely available. Our approach did not rely on background error correction models  
373 that are constructed from large numbers of healthy donor samples and are therefore  
374 impractical for applications requiring frequently changing custom gene panels,  
375 including clinical assay development.



376           Importantly, the 1.5 h fast-hybridization step (standard protocol: 16h) used in  
377 our assay dramatically reduces library preparation time which is advantageous when  
378 fast turnaround is critical. Increasing the wash temperature after capture dramatically  
379 reduced off target reads. The higher temperature likely relaxes the target/bait-bond  
380 in hybridised molecules with a higher number of mismatches, reducing the non-  
381 specific carry over of DNA fragments into the library.

382           cfDNA-Seq of 28 mCRC patients demonstrated that 88% of mutations  
383 detected by clinical grade tumor tissue sequencing were also called in cfDNA. This  
384 detection capability is similar to that reported for MBC-error corrected cfDNA-Seq  
385 with a 5-gene assay using amplicons (87.2%) [1] and a 54-gene assay using target-  
386 capture (85%) [14, 30]. Furthermore, indels are more difficult to call than point  
387 mutations. Yet, our cfDNA assay called 23/26 indels (88.5%) that were known based  
388 on tumor sequencing, showing a similar performance to point mutation detection  
389 (87.7% called).

390           cfDNA-Seq detected several additional driver mutations not reported by tumor  
391 sequencing. Seven were in *TP53*. Two were also observed in the matched blood  
392 cells, indicating that they originated from clonal hematopoiesis. The discovery of  
393 clonal hematopoiesis in 7% of our cohort demonstrates the importance of  
394 sequencing DNA extracted from blood cells to avoid misinterpreting such variants as  
395 cancer-associated mutations. In one patient who received cetuximab therapy, we  
396 detected a *KRAS* Q61H variant that was absent from the matched tumor and likely  
397 represents the evolution of a drug resistant subclone. Multiple *PIK3CA* activating  
398 mutations detected in two anti-EGFR therapy naive patients represent parallel  
399 evolution events. These examples show that our cfDNA assay can provide insights  
400 into cancer evolution. Because the minimally invasive nature of cfDNA-Seq allows

401 application at multiple time-points, this could be used to monitor the evolution of  
402 subclonal drug resistance driver mutations without prior knowledge of specific loci  
403 where resistance mutations will occur. We finally demonstrate that cfDNA-Seq allows  
404 genome-wide CNA reconstruction and validate this against low-coverage genome  
405 sequencing. As the number of targeted therapies increases, custom target  
406 enrichment panels that can be readily adapted and scaled for the tumor type and  
407 therapeutic agent in question could be used to investigate the full tumor genomic  
408 landscape of point mutations, indels and CNAs. This would facilitate the identification  
409 of novel resistance mechanisms. Importantly, this ultra-sensitive cfDNA-Seq  
410 technology can also address the subset of 20% of patients with mCRC who cannot  
411 be molecularly profiled due to unobtainable or inadequate biopsy tissues [16, 31].

412 In conclusion, this cfDNA-Seq approach with customizable and off-the-shelf  
413 reagents showed a similar performance to published techniques that use bespoke  
414 reagents and more complex analyses.

415

416

417

418

419

420

421

422 **Data access**

423 Sequencing fastq files have been deposited into the NCBI Sequence Read  
424 Archive (SRA submission code SUB3510375).

425

426 **Acknowledgements**

427 We would like to thank all patients participating in the FORMAT clinical trial  
428 and the clinical research team members at the Royal Marsden Hospital who  
429 supported the sample collection. The study was supported by charitable donations  
430 from Tim Morgan to the Institute of Cancer Research, from Philip Moodie to The  
431 Royal Marsden Cancer Charity and by a Clive and Ann Smith Fellowship. The study  
432 received funding by Cancer Research UK, a Wellcome Trust Strategic Grant  
433 (105104/Z/14/Z), the Royal Marsden Hospital/Institute of Cancer Research National  
434 Institute for Health Research Biomedical Research Centre for Cancer and by a  
435 Cancer Research UK Clinical PhD Studentship.

436

437 **Disclosure Declaration**

438 The authors had pre-marketing access to Agilent SureSelect<sup>XT-HS</sup> reagents. BH is an  
439 employee of Agilent. The other authors received no financial support or  
440 compensation from Agilent.

441

442 **Statement of Author Contributions**

443 SM, LJB and MG conceived the study and wrote the manuscript; SM, LJB and BG  
444 processed samples; SM, LJB and BH developed the cfDNA-Seq assay; DK and SL  
445 developed the DuplexCaller tool; SM, LJB, DK, AW, MND and MG analyzed the  
446 data; SYM, MD, AP, AO, IR, RB, DW, AWoth, KvL, IC, DC, NS and TP provided

447 clinical data and samples. KF and NM sequenced the cfDNA libraries, PZP, DGDC,  
448 SH and MH provided tumor biopsy sequencing data from the FORMAT panel and  
449 ran the Avenio analysis, NT and BOL provided support for ddPCR.

450

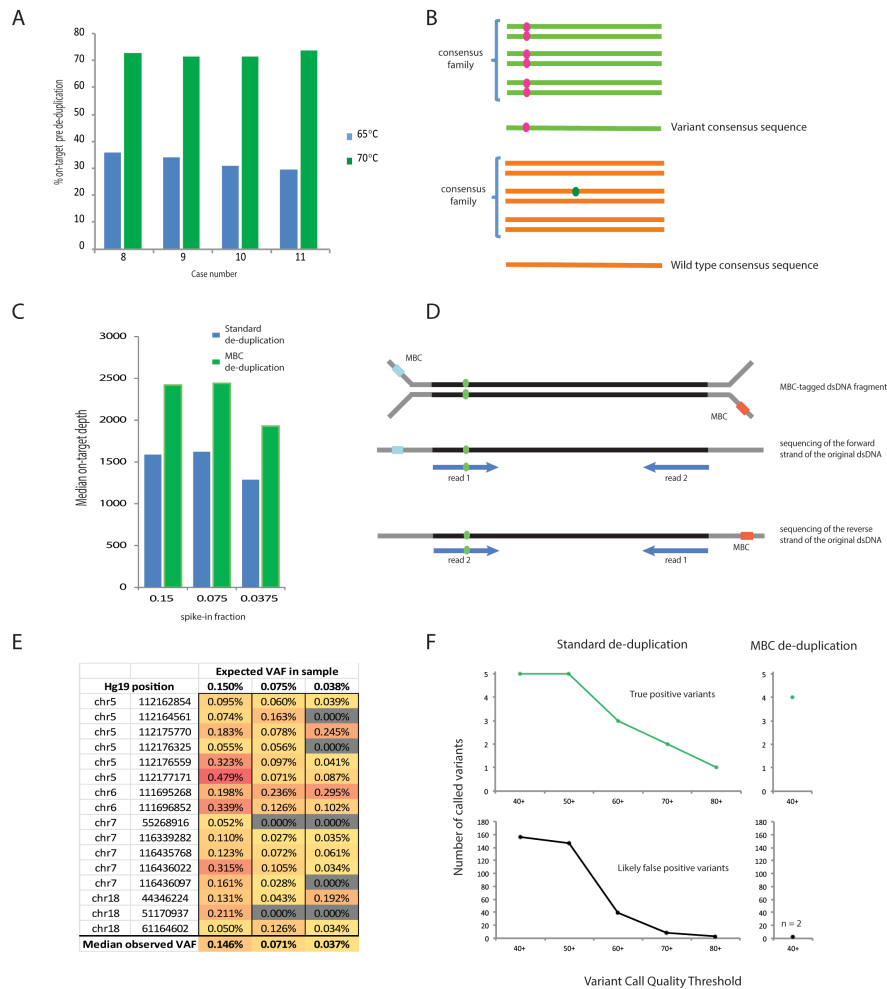
451 Sonia Mansukhani<sup>1\*</sup>, Louise J. Barber<sup>1\*</sup>, Dimitrios Kleftogiannis<sup>1</sup>, Sing Yu  
452 Moorcraft<sup>2</sup>, Michael Davidson<sup>2</sup>, Andrew Woolston<sup>1</sup>, Paula Zuzanna Proszek<sup>3</sup>,  
453 Beatrice Griffiths<sup>1</sup>, Kerry Fenwick<sup>4</sup>, Bram Herman<sup>5</sup>, Nik Matthews<sup>4</sup>, Ben O'Leary<sup>6</sup>,  
454 Sanna Hulkki<sup>3</sup>, David Gonzalez De Castro<sup>7</sup>, Anisha Patel<sup>8</sup>, Andrew Wotherspoon<sup>9</sup>,  
455 Aleruchi Okachi<sup>2</sup>, Isma Rana<sup>2</sup>, Ruwaida Begum<sup>2</sup>, Matthew N. Davies<sup>1,10</sup>, Thomas  
456 Powles<sup>11</sup>, Katharina von Loga<sup>1</sup>, Michael Hubank<sup>3</sup>, Nick Turner<sup>6,12</sup>, David Watkins<sup>2</sup>,  
457 Ian Chau<sup>2</sup>, David Cunningham<sup>2</sup>, Stefano Lise<sup>1</sup>, Naureen Starling<sup>2</sup> and Marco  
458 Gerlinger<sup>1,2</sup>

## References

1. Bettgowda C, Sausen M, Leary RJ, Kinde I, Wang Y, Agrawal N, et al. Detection of circulating tumor DNA in early- and late-stage human malignancies. *Sci Transl Med* 2014, 6(224):224ra224.
2. Siravegna G, Marsoni S, Siena S, Bardelli A: Integrating liquid biopsies into the management of cancer. *Nat Rev Clin Oncol* 2017, 14(9):531-548.
3. Heitzer E, Ulz P, Belic J, Gutsch S, Quehenberger F, Fischereder K, et al. Tumor-associated copy number changes in the circulation of patients with prostate cancer identified through whole-genome sequencing. *Genome Med* 2013, 5(4):30.
4. Haber DA, Velculescu VE: Blood-Based Analyses of Cancer: Circulating Tumor Cells and Circulating Tumor DNA. *Cancer Discov* 2014, 4(6):650-661.
5. Leary RJ, Sausen M, Kinde I, Papadopoulos N, Carpten JD, Craig D, et al. Detection of chromosomal alterations in the circulation of cancer patients with whole-genome sequencing. *Sci Transl Med* 2012, 4(162):162ra154.
6. Dawson SJ, Tsui DW, Murtaza M, Biggs H, Rueda OM, Chin SF, et al. Analysis of circulating tumor DNA to monitor metastatic breast cancer. *N Engl J Med* 2013, 368(13):1199-1209.
7. Mok T, Wu YL, Lee JS, Yu CJ, Sriuranpong V, Sandoval-Tan J, et al. Detection and Dynamic Changes of EGFR Mutations from Circulating Tumor DNA as a Predictor of Survival Outcomes in NSCLC Patients Treated with First-line Intercalated Erlotinib and Chemotherapy. *Clin Cancer Res* 2015, 21(14):3196-3203.
8. Garcia-Murillas I, Schiavon G, Weigelt B, Ng C, Hrebien S, Cutts RJ, et al. Mutation tracking in circulating tumor DNA predicts relapse in early breast cancer. *Sci Transl Med* 2015, 7(302):302ra133.
9. Phallen J, Sausen M, Adleff V, Leal A, Hruban C, White J, et al. Direct detection of early-stage cancers using circulating tumor DNA. *Sci Transl Med* 2017, 9(403).
10. Abbosh C, Birkbak NJ, Wilson GA, Jamal-Hanjani M, Constantin T, Salari R, et al. Phylogenetic ctDNA analysis depicts early-stage lung cancer evolution. *Nature* 2017, 545(7655):446-451.
11. Gerlinger M, Horswell S, Larkin J, Rowan AJ, Salm MP, Varela I, et al. Genomic architecture and evolution of clear cell renal cell carcinomas defined by multiregion sequencing. *Nat Genet* 2014, 46(3):225-233.
12. Diehl F, Li M, He Y, Kinzler KW, Vogelstein B, Dressman D: BEAMing: single-molecule PCR on microparticles in water-in-oil emulsions. *Nat Methods* 2006, 3(7):551-559.
13. Perakis S, Speicher MR: Emerging concepts in liquid biopsies. *BMC Med* 2017, 15(1):75.
14. Lanman RB, Mortimer SA, Zill OA, Sebisano D, Lopez R, Blau S, et al. Analytical and Clinical Validation of a Digital Sequencing Panel for Quantitative, Highly Accurate Evaluation of Cell-Free Circulating Tumor DNA. *PLoS One* 2015, 10(10):e0140712.
15. Newman AM, Lovejoy AF, Klass DM, Kurtz DM, Chabon JJ, Scherer F, et al. Integrated digital error suppression for improved detection of circulating tumor DNA. *Nat Biotechnol* 2016, 34(5):547-555.
16. Moorcraft SY, Gonzalez de Castro D, Cunningham D, Jones T, Walker BA, Peckitt C, et al. Investigating the feasibility of tumour molecular profiling in gastrointestinal malignancies in routine clinical practice. *Ann Oncol* 2018, 29(1):230-236.
17. GitHub. <https://github.com/dkleftogi/duplexFiltering/blob/master/duplexCaller.py>. duplexCaller. (Accessed January 2018).
18. GitHub. bam-readcount. <https://github.com/genome/bam-readcount> (Accessed October 2017)
19. Talevich E, Shain AH, Botton T, Bastian BC: CNVkit: Genome-Wide Copy Number Detection and Visualization from Targeted DNA Sequencing. *PLoS Comput Biol* 2016, 12(4):e1004873.

20. Baslan T, Kendall J, Rodgers L, Cox H, Riggs M, Stepansky A, et al. Genome-wide copy number analysis of single cells. *Nat Protoc* 2012, 7(6):1024-1041.
21. Kennedy SR, Schmitt MW, Fox EJ, Kohn BF, Salk JJ, Ahn EH, et al. Detecting ultralow-frequency mutations by Duplex Sequencing. *Nat Protoc* 2014, 9(11):2586-2606.
22. Robinson JT, Thorvaldsdottir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov JP: Integrative genomics viewer. *Nat Biotechnol* 2011, 29(1):24-26.
23. James T. Robinson HT, Wendy Winckler, Mitchell Guttman, Eric S. Lander, Gad Getz, Jill P. Mesirov: Integrative Genomics Viewer. *Nat Biotechnol* 2011, 29:24-26.
24. Forbes SA, Beare D, Bindal N, Bamford S, Ward S, Cole CG, et al: COSMIC: High-Resolution Cancer Genetics Using the Catalogue of Somatic Mutations in Cancer. *Curr Protoc Hum Genet* 2016, 91:10 11 11-10 11 37.
25. Xie M, Lu C, Wang J, McLellan MD, Johnson KJ, Wendl MC, et al: Age-related mutations associated with clonal hematopoietic expansion and malignancies. *Nat Med* 2014, 20(12):1472-1478.
26. Coombs CC, Zehir A, Devlin SM, Kishtagari A, Syed A, Jonsson P, et al. Therapy-Related Clonal Hematopoiesis in Patients with Non-hematologic Cancers Is Common and Associated with Adverse Clinical Outcomes. *Cell Stem Cell* 2017, 21(3):374-382.e374.
27. Misale S, Yaeger R, Hobor S, Scala E, Janakiraman M, Liska D, et al: Emergence of KRAS mutations and acquired resistance to anti EGFR therapy in colorectal cancer. *Nature* 2012, 486(7404):532-536.
28. Brannon AR, Vakiani E, Sylvester BE, Scott SN, McDermott G, Shah RH, et al. Comparative sequencing analysis reveals high genomic concordance between matched primary and metastatic colorectal cancer lesions. *Genome Biol* 2014, 15(8):454.
29. TCGA: Comprehensive molecular characterization of human colon and rectal cancer. *Nature* 2012, 487(7407):330-337.
30. Kim ST, Lee WS, Lanman RB, Mortimer S, Zill OA, Kim K-M, et al. Prospective blinded study of somatic mutation detection in cell-free DNA utilizing a targeted 54-gene next generation sequencing panel in metastatic solid tumor patients. *Oncotarget* 2015, 6(37):40360-40369.
31. Khakoo S, Georgiou A, Gerlinger M, Cunningham D, Starling N: Circulating tumour DNA, a promising biomarker for the management of colorectal cancer. *Crit Rev Oncol Hematol* 2018, 122:72-82.

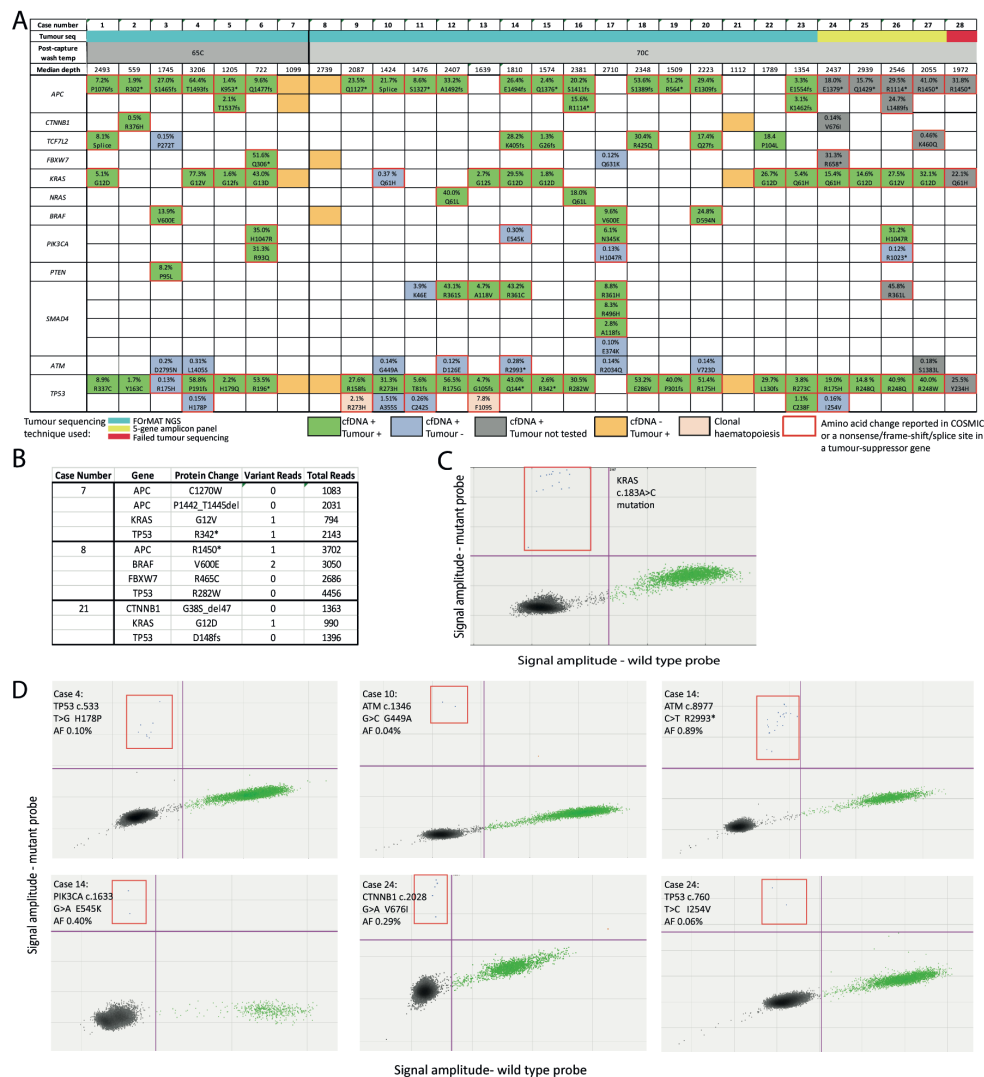
## Figures



**Figure 1. (A)** Percentage of reads on-target before de-duplication in samples prepared with 65°C vs 70°C post-capture washes. **(B)** Graphic depicting the principles of MBC error correction. Reads with the same MBC that map to the identical genomic location are grouped into a consensus family. If a variant (pink) occurs in all reads then the consensus read sequence will be variant for that base (top). However if a variant (green) is only detected in a small fraction of the reads in the family, it will be disregarded and the consensus read sequence will be wild-type (bottom). **(C)** cfDNA mixing experiment: 25 ng mixes of donor A spiked into donor B at 0.15%, 0.075% and 0.0375%. **(D)** Illustration of duplex read pair detection. A double stranded cfDNA fragment (black) containing a variant (green) is depicted, ligated to Y-shaped MBC-tagged adapters (grey). **(E)** Expected and observed variant allele frequencies (VAF) and genomic positions for the 16 SNPs in the cfDNA mixing

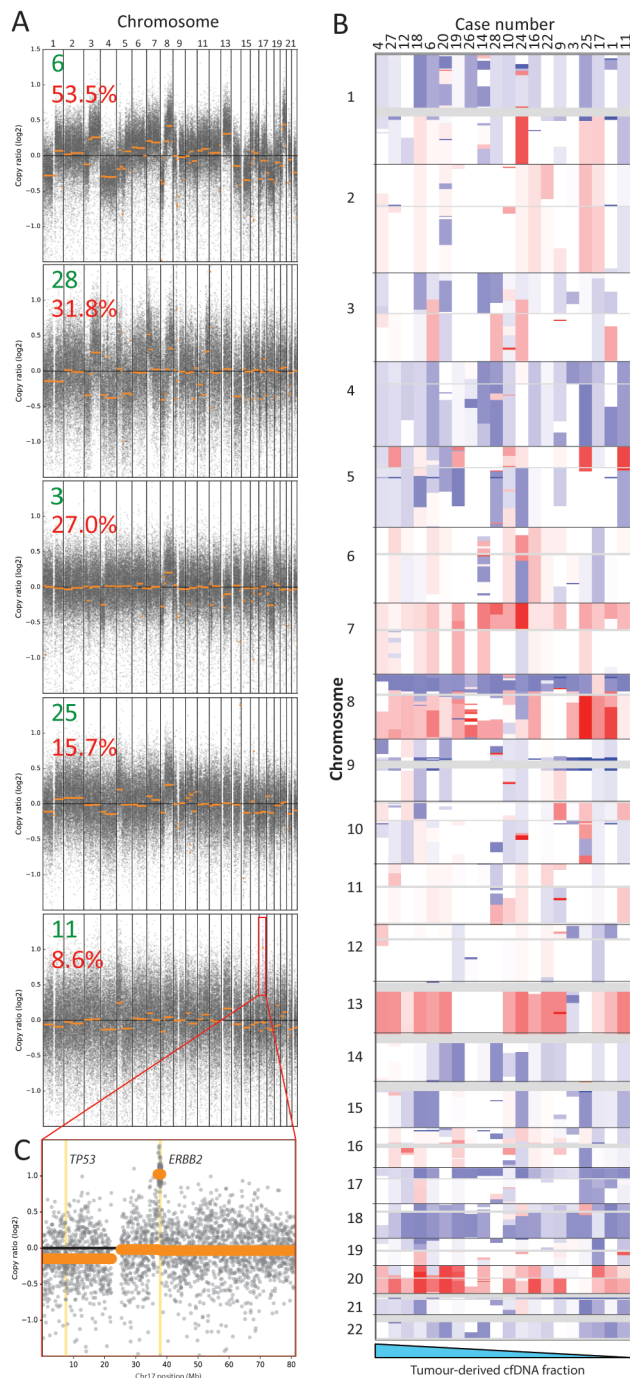
experiment. **(F)** Impact of MBC error correction on true positive and false positive calls. The top panels show the number of true positive variants (expected SNPs) that were bioinformatically called in the mixing experiment with standard de-duplication (left) and MBC de-duplication (right) using different variant call quality thresholds. The lower panel shows the number of likely false positive variant calls (not observed in the deep sequencing of either cfDNA sample used in the mix) for standard de-duplication (left) and MBC de-duplication (right).





**Figure 2 (A)** Concordance of mutations identified by cfDNA-Seq and by sequencing of tumor material. Mutations identified in both cfDNA-Seq and tumor sequencing are colored green. Novel variants called by cfDNA-Seq and not by tumor sequencing are colored blue. Variants not detected by cfDNA-Seq that were detected in tumor sequencing are colored orange. Pink indicates clonal hematopoiesis. Red outlines indicate mutations reported as tumorigenic in COSMIC. Variants in grey have been identified in the cfDNA of patients that either had been sequenced using the limited 5-gene amplicon panel or failed FORMAT sequencing. Percentages indicate VAF in cfDNA. **(B)** Read depth and number of consensus family reads supporting each of the 11 variants in cases 7, 8, and 21 that had not been called in cfDNA but had previously been detected in tumor tissue. Median VAF 0.066%. **(C)** ddPCR validation

of the *KRAS* c.183A>C mutation that results in the amino acid change Q61H in case 10. Green dots: droplets with wild-type DNA, blue dots (outlined by the red quadrant): droplets with mutant DNA, black dots: droplets that have no incorporated DNA. **(D)** ddPCR validation of 6 subclonal mutations called in cfDNA but not in tumor tissue.



**Figure 3 (A)** Genome wide copy number aberrations can be detected from targeted cfDNA-Seq, even where tumor content is low. Representative log copy ratio plots for five cases (green number) in our cohort with tumor content ranging from 53.5% to 8.6% (red number indicates max VAF) are shown. **(B)** Genome wide heat map of segmented copy number raw log ratio data after amplitude normalization. Gains are red and losses are blue. Profiles are ordered (left to right) from highest to lowest tumor content (based on maximum VAF) for all 20 cases that had a visible CNA

profile. **(C)** Focused log copy ratio plot of chromosome 17 for case 11 which had a high level amplification of *ERBB2*.