**AmpliconReconstructor: Integrated analysis of NGS and optical mapping resolves the complex structures of focal amplifications in cancer**

**Authors:** Jens Luebeck[1,2], Ceyda Coruh[3], Siavash R. Dehkordi[2], Joshua T. Lange[4,5], Kristen M. Turner[5], Viraj Deshpande[2], Dave A. Pai[6], Chao Zhang[1], Utkrisht Rajkumar[2], Julie A. Law[3], Paul S. Mischel[5,7,8], Vineet Bafna[2]*

**Affiliations:**

[1]Bioinformatics and Systems Biology Graduate Program, University of California at San Diego, La Jolla, CA 92093, USA

[2]Department of Computer Science and Engineering, University of California at San Diego, La Jolla, CA 92093, USA

[3]Plant Molecular and Cellular Biology Laboratory, Salk Institute for Biological Studies, La Jolla, CA 92037, USA

[4]Biomedical Sciences Graduate Program, University of California at San Diego, La Jolla, CA 92093, USA

[5]Ludwig Institute for Cancer Research, University of California at San Diego, La Jolla, CA 92093, USA

[6]Bionano Genomics, Inc., San Diego, CA 92121, USA

[7]Moores Cancer Center, University of California at San Diego, La Jolla, CA 92093, USA

[8]Department of Pathology, University of California at San Diego, La Jolla, CA 92093 USA

*To whom correspondence should be addressed (vbafna [at] cs.ucsd.edu)

**Abstract:**

Oncogene amplification, a major driver of cancer pathogenicity, is often mediated through focal amplification of genomic segments. Recent results implicate extrachromosomal DNA (ecDNA) as the primary mechanism driving focal copy number amplification (fCNA) - enabling gene amplification, rapid tumor evolution, and the rewiring of regulatory circuitry. Resolving an fCNA's structure is a first step in deciphering the mechanisms of its genesis and the subsequent biological consequences. Here, we introduce a powerful new computational method, AmpliconReconstructor (AR), for integrating optical mapping (OM) of long DNA fragments (>150kb) with next-generation sequencing (NGS) to resolve fCNAs at single-nucleotide resolution. AR uses an NGS-derived breakpoint graph alongside OM scaffolds to produce high-fidelity reconstructions. After validating performance by extensive simulations, we used AR to reconstruct fCNAs in seven cancer cell lines to reveal the complex architecture of ecDNA, breakage-fusion-bridge cycles, and other complex rearrangements. By distinguishing between chromosomal and extrachromosomal origins, and by reconstructing the rearrangement signatures associated with a given fCNA's generative mechanism, AR enables a more thorough understanding of the origins of fCNAs, and their functional consequences.

**Main:**

Oncogene amplification is a major driver of cancer pathogenicity[1–5]. Genomic signatures of oncogene amplification include somatic focal Copy Number Amplifications (fCNAs) of small (typically < 10Mbp) genomic regions[5,6]. Multiple mechanisms cause fCNAs including, but not limited to, extrachromosomal DNA (ecDNA) formation[5,7,8], chromothripsis[9], tandem duplications[10,11] and breakage-fusion-bridge (BFB) cycles[12–14]. ecDNA, in particular, enables tumors to achieve far higher oncogene genomic copy numbers and maintain far greater levels of intratumor genetic heterogeneity than previously anticipated, due to their non-chromosomal mechanism of inheritance - enabling tumors to evolve rapidly[5,15,16]. In addition, the very high DNA template level generated by ecDNA-based amplification, coupled to its highly accessible chromatin architecture, permits massive oncogene transcription[17–19].

63    While ecDNA elements are a common form of fCNA[5], other mechanisms can also result

64    in amplification with very different functional consequences[6]. Thus, accurate identification

65    and reconstruction of the fCNA structure not only describes the rearranged genomic

66    landscape, but also represents a first step in identifying the generative mechanism.

67    Reconstruction of fCNA architecture involves determining the order and orientation of the

68    genomic segments that constitute the amplicon. There are many methods to detect single

69    genomic breakpoints from sequencing data using a variety of different sequencing

70    technologies[20–23]. However, fewer methods are available to handle the more difficult

71    problem of ordering and orienting multiple genomic segments joined by breakpoints into

72    high confidence scaffolds which enable complete reconstructions of complex

73    rearrangements[6,24].

74

75    A previous method for characterizing the identity of focally amplified genomic regions,

76    AmpliconArchitect (AA), generates an accurate breakpoint graph from next-generation

77    sequencing (NGS) data[6]. The graph encodes the genomic segments involved in fCNAs,

78    their copy numbers, and breakpoint edges connecting the segments. Unambiguous

79    reconstruction of fCNA architecture requires extracting paths and cycles from the

80    breakpoint graph, to reveal the true structure of the underlying rearranged genome.

81    However, in practice, path/cycle extraction is often confounded by duplications of large

82    genomic regions inside an amplicon (Supplemental Fig. S1a), imperfections in the graph

83    arising from errors in estimation of segment copy numbers, or erroneous and/or missing

84    breakpoints.

85

86    We hypothesized that an approach combining the strengths of NGS with long-range

87    mapping data would enable larger and more unambiguous reconstructions of fCNA

88    architectures. To produce the highest-quality reconstructions of large, complex focal

89    amplifications, we utilized both optical mapping (OM) data as well as the breakpoint

90    graphs generated from AA with NGS data. OM provides single-molecule information

91    about the approximate locations of fluorescently-labeled sequence motifs on long

92    fragments of DNA[25]. The median molecule ("map") length used in assembly across all

93    samples used in this study is 244 kbp (molecule N50 340 kbp), while the median segment

length in breakpoint graphs used in this study is 100 kbp, highlighting that OM data can span multiple junctions in breakpoint graphs derived from focal amplifications (Supplemental Table 1). The integrated NGS data and OM data provide an orthogonal pairing of short- and long-range information about genomic structural variation. We utilized Bionano (Bionano Genomics, Inc., San Diego, CA) whole-genome imaging to generate single-molecule optical maps, which were subsequently *de novo* assembled into OM contigs (contig N50 72.8 Mbp) to improve confidence and reduce noise.

Here we present a novel computational method for reconstructing large complex fCNAs, AmpliconReconstructor (AR). AR takes a breakpoint graph and long-range Bionano OM data as inputs. AR produces an ordering and orientations of graph segments, with fine-structure information from the breakpoint graph embedded into the large-scale reconstructions. AR outputs megabase-scale reconstructions of fCNA amplicons. We demonstrated the fine-scale accuracy of AR using simulated OM data derived from previously analyzed cancer cell lines[6,21]. Furthermore, we reconstructed fCNAs at unprecedented resolution and size in seven cancer cell lines (CAKI2, GBM39, NCI-H460, HCC827, HK301, K562, T47D). Finally, we validated the reconstructions using cytogenetics.

**Results:**

AR separates the computational tasks involved in reconstruction of fCNAs into three primary modules (Fig. 1a,b). First, an OM alignment module, SegAligner, for aligning reference segments to assembled OM contigs generated by either the Bionano Irys or Bionano Saphyr instruments (Supplemental Fig. S1b-d, Methods - "Optical map contig alignments with SegAligner"). SegAligner is critical to the imputation process as it can score placements of short genomic segments onto an OM contig, which wasn't possible with other aligners. Second, a scaffolding module, which takes a collection of breakpoint graph segments aligned to OM contigs as input and creates scaffolds represented by directed acyclic graphs (DAGs) (Fig. 1c-e, Methods – "Reconstructing amplicon paths with AmpliconReconstructor"). Third, AR also relies on a novel scaffold-path imputation technique (Fig. 1f-h, Methods – "Imputing paths in the scaffold with

125    AmpliconReconstructor") to connect and chain together breakpoint graph segments that

126    may individually be too small to be informatively labeled and aligned with optical mapping

127    (Fig. 1f). Finally, a pathfinding module, which links scaffolds and searches for paths in a

128    copy number (CN)-aware manner, to identify possible reconstructions of the amplicon.

129    AR outputs a collection of sequence resolved paths supported by the linked scaffolds. To

130    visualize output from AR, we implemented a visualization utility, CycleViz, to show the

131    integrated OM- and NGS-derived breakpoint graph data (Supplemental Fig. S2).

132    AmpliconReconstructor is implemented in Python 2.7, and SegAligner is implemented in

133    C++.        Both        tools        are        available        publicly        at

134    https://github.com/jluebeck/AmpliconReconstructor.

135

**AR accurately reconstructs simulated amplicons**

137    We utilized multiple simulation strategies to measure the performance of AR. For a

138    ground-truth set of amplicon structures, we used 85 non-trivial amplicon breakpoint graph

139    paths previously reported by AmpliconArchitect from 25 cancer cell lines[6]. The breakpoint

140    graph paths included both cyclic and non-cyclic paths with lengths varying from 260 kbp

141    to 2.8 Mbp (median 1.1 Mbp) and the number of graph segments varying from 3 to 47

142    (mean 17.5 segments; Supplemental Table 2). These paths were used as a reference to

143    simulate OM molecules. (Methods – "Simulation of amplicons to measure AR

144    performance"). Simulated molecules were assembled into contigs using the Bionano

145    Assembler[26,27].

146

147    For each of the 85 simulation cases, we then ran AR on the corresponding breakpoint

148    graph and the *de novo* assembled contigs, and examined four different variables that

149    could affect the performance of AR. First, we tested AR performance using SegAligner

150    for OM alignment, versus AR using other OM alignment tools to replace SegAligner.

151    Second, we evaluated the performance of AR across a range of amplicon copy numbers.

152    Third, we measured performance with false edges present in the breakpoint graph.

153    Finally, we generated and tested mixtures of three similar amplicons from the same

154    samples, simulated with different amplicon copy numbers, to measure the effects of

155    potential amplicon heterogeneity on AR performance.

156

157 We measured the accuracy of AR by computing precision and recall across the four

158 simulation conditions. As precision and recall could be quantified in multiple ways when

159 comparing ground-truth and reconstructed simulation paths, leading to different

160 understandings of performance, we described three ways of measuring the similarity of

161 the paths ("Length (bp)", "Nseg", "Breakpoint"; Methods – "Measuring AR simulation

162 performance"), based on the longest common substring (LCS) between ground-truth and

163 reconstructed path sequences. We report the "Length (bp)" measurement in the analysis

164 described here, while results with other measurements are presented in Supplemental

165 Table 2 and Supplemental Figure S3.

166

167 AR using SegAligner achieved a mean F1 score (harmonic mean of the precision and

168 recall) of 0.88 for the highest copy number level (CN 20) and 0.68 for the lowest copy

169 number level (CN 2) (Fig. 1i, Supplemental Fig. S3, Supplemental Table 2). In contrast,

170 when OMBlast[28] or Bionano RefAligner[26,29] were used in place of SegAligner, we noticed

171 a decrease in both precision and recall. For RefAligner and OMBlast, respectively, we

172 report mean F1 scores of 0.52, 0.43 for CN 20, and 0.42, 0.41 for CN 2. When imputation

173 was omitted from AR, the mean F1 score for CN 20 decreased from 0.88 to 0.70. We

174 observed similarly consistent trends using other methods of measuring precision and

175 recall – "Nseg" and "Breakpoint" (Supplemental Fig. S3). We saw a few cases of

176 'assembly failure,' where no paths differing from the reference genome involving the

177 amplicon segments were assembled.  Figure 1i shows cumulative precision and recall

178 values for AR using SegAligner (with and without imputation), and with assembly failures

179 filtered. We additionally reported simulation F1 scores with and without filtering for

180 possible OM assembly failure (Supplemental Table 2).

181

182 False edges in the breakpoint graph increase the possible number of path imputations

183 that AR must consider, potentially leading to erroneous scaffolds. We designed another

184 simulation study where after simulating CN 20 amplicon OM data, additional false edges

185 were added between existing graph segments. We tested three scenarios with the

186 proportion of additional false edges ranging from 0%, 50% and 100% of the number of

187  true graph edges. The three scenarios resulted in nearly identical mean F1 scores of

188  0.881, 0.880, 0.881 across the 85 amplicon simulations (Supplemental Table 2,

189  Supplemental Fig. S4a), highlighting the robustness of the path imputation method.

190

191  To understand how AR performed when faced with structural heterogeneity, we designed

192  a simulation study involving 123 combinations of breakpoint graph paths where each

193  combination was derived from a single sample at varying copy number mixtures. We

194  simulated amplicons from heterogeneous mixtures with (1) a single dominant amplicon

195  (CNs 20-2-2); (2) a linear mixture of CNs (CNs 20-15-10); (3) equally abundant amplicons

196  (CNs 20-20-20). We report mean F1 scores of 0.92, 0.89, and 0.91, respectively for the

197  three cases (Supplemental Table 2). To explain the increase in performance of the

198  mixture simulations as compared to the single amplicon simulations, we hypothesize that

199  the greater total number of molecules improved the assembly process. Regardless, the

200  high similarity between the precision and recall in each mixture case (Supplemental Fig.

201  S4b) indicates AR can reconstruct an accurate amplicon path even in the context of

202  heterogeneity. Based on these metrics, we found AR to be robust, and to outperform

203  other methods. To further demonstrate its ability to reconstruct a variety of complex

204  fCNAs, we ran AR on seven cancer cell lines with evidence of fCNA.

205

206  **AR reconstructs ecDNA in multiple forms**

207  Three cell lines in our data set were previously reported to contain ecDNA[5] - GBM39,

208  NCI-H460, and HK301. In a previous study[17], we analyzed the glioblastoma multiforme

209  (GBM) cell line GBM39 using a preliminary version of AR that used RefAligner and

210  manual merging of graph segments, but without path imputation or scaffold linking

211  capabilities. Re-analysis reproduced an unambiguous 1.26 Mbp EGFRvIII-containing

212  circular ecDNA that was identical to the previously published structure[17] (Supplemental

213  Fig. S5). The entire structure was captured by a single non-circular OM contig, with

214  circularity confirmed by an overlapping graph segment aligned to both ends of the contig.

215

216  Previous studies of ecDNA have documented their integration into chromosomes over

217  time, linearizing and appearing as homogeneously staining regions (HSRs), often in non-

7

218 native locations[5,7,15]. In a previous study[5], The GBM cell line, HK301, had been

219 cytogenetically determined to have circular ecDNA; however, we observed from FISH

220 (fluorescence in situ hybridization) data that the sample's ecDNA had become HSR-like

221 at the time of this study (Fig. 2a). AA generated a breakpoint graph supporting

222 amplification of both EGFRvIII and EGFR wild-type (Fig. 2c), however an unambiguous

223 reconstruction from the graph alone was not possible. The AR reconstruction of the

224 HK301 fCNA indicated a complex and cyclic structure supported by three contigs (Fig.

225 2d), which explained 98.1% of the amplified genomic regions. The graph segments came

226 predominantly from chr7, but also included two small regions (2890 bp, 4591 bp) from

227 chr6 (Fig. 2c,d). We noted a ~20 kbp deletion inside EGFR, showing a lower CN than

228 the surrounding region, but which was still amplified over the baseline, non-amplicon

229 regions of chr7. This indicates heterogeneity of EGFR wild-type/vIII mutation status.

230 Despite the heterogenous status of this allele, AR reconstructed the EGFRvIII version –

231 which is the dominant form of the amplicon (Fig. 2d).

232

233 The lung cancer cell line NCI-H460 has previously been documented to bear MYC

234 amplification[30], and our cytogenetic analysis showed evidence for both its HSR-like and

235 ecDNA amplification (Fig. 2e,f). Despite the heterogeneous nature of the amplicon's

236 integration status, AA generated a breakpoint graph for a contiguous 2.15 Mbp region of

237 chr8 (Fig. 2g). AR reconstructed a single 4.10 Mbp structure supported by five OM contigs

238 (Fig. 2h). This structure contained all amplified segments from the breakpoint graph and

239 explained the breakpoint graph segment copy number ratios of the duplicated segments.

240 For example, segment chr8:129,404,278-129,591,422 appeared 4 times,

241 chr8:128,690,200-129,404,277 (carrying MYC & PVT1) appeared twice,

242 chr8:129,591,423-129,911,811 appeared twice, and chr8:129,911,812-130,640,594

243 appeared once, making the ratios consistent with the estimated graph segment copy

244 numbers (46, 25, 25, 12, respectively; Fig. 2g). The status of the long non-coding RNA

245 PVT1 (a known regulator of MYC)[31] on this amplicon is heterogeneous, as one copy of

246 PVT1 does not contain breakpoints, while the other shows a disrupted copy of PVT1. AR

247 also identified a self-inversion at the end of the amplicon (black arrows in Fig. 2h),

248     suggestive of an alternating forward-backward orientation (segmental tandem

249     aggregation with inversion) of the amplicon in the agglomerated ecDNA.

250

251     In summary, AR reconstructed paths that were consistent with the expected copy number

252     ratios and graph structures in GBM39, HK301, and NCI-H460, explaining 99.9%, 98.1%,

253     and 100% of the amplified genomic content in the breakpoint graphs for each cell line,

254     respectively. Furthermore, the AR reconstructions of ecDNA in HSR-like form lend

255     additional evidence to the agglomerative model of ecDNA integration (Fig. 2b)[8,32,33].

256

257     **AR reconstructs a rearranged Philadelphia chromosome in K562**

258     The classical model of the BCR-ABL1 (Philadelphia chromosome) fusion involves a

259     reciprocal translocation of the q arms of chromosomes 9 and 22[34]. However, this

260     mechanism alone does not explain the copy number amplification of BCR-ABL1 fusion

261     commonly observed in chronic myeloid leukemia (CML), highlighting a need for methods

262     to better understand the genesis of the BCR-ABL1 amplification[35,36]. To reconstruct the

263     fine structure of a Philadelphia chromosome, we used the CML cell line K562 where a

264     BCR-ABL1 fusion had previously been reported[37].

265

266     The AA reconstructed breakpoint graph of the BCR-ABL1 fCNA in K562 (Fig. 3a) contains

267     8.5 Mbp of amplified genomic segments. The graph shows signatures of complex

268     rearrangements alongside the BCR-ABL1 fusion, which AA predicted to have a copy

269     number of 17 (Fig. 3a). We generated both Bionano Irys and Bionano Saphyr OM data

270     for K562 cells and observed consistent results in the independent reconstructions of

271     amplicons from both sources (Supplemental Fig. S6a,b). Using the breakpoint graph and

272     OM contigs, AR reconstructed a complex linear structure that chained together 1.7 Mbp

273     from chr22 (containing BCR), 548 kbp of chr9 (containing ABL1), and multiple regions

274     from chr13 (732 kbp; including a disrupted copy of GPC5) (Fig. 3b). In Figure 3b, we show

275     one possible scaffolding of the given regions, whose structure was reproduced in both

276     Saphyr and Irys datasets. AR also reported junctions between segments in the breakpoint

277     graph where NGS-derived breakpoint edges were not reported, as indicated by the

278  missing half-height grey bars between adjacent genomic segments in the genome tracks
279  of Figure 3b.

280

281  We performed multiple FISH experiments using combinations of probes for BCR, ABL1,
282  GPC5, and chr22 centromere probe CEP22. The FISH images confirmed the co-
283  localization of the BCR-ABL1 fusion and GPC5 on a common HSR-like structure (Fig.
284  3c). Furthermore, it validated the status of the K562 BCR-ABL1 fusion as being located
285  on chr22 (Supplemental Fig. S7).

286

287  In addition to the reconstruction reported in Figure 3b, AR additionally identified other
288  scaffolds, indicating that the genomic structure surrounding the BCR-ABL1 translocation
289  may be varied across the multiple copies (Supplemental Fig. S6c,d; Supplemental Fig.
290  S8a-f). In particular, the genomic segment bearing CLTCL appears in both forward and
291  reverse directions (Supplemental Fig. S8b,c). Other amplified regions of chr13 include a
292  self-inversion at the 3' end of GPC5 (Supplemental Fig. S6c,d, Supplemental Fig. S8e).
293  A scaffold from the Irys-based reconstruction indicated a secondary reconstruction could
294  be joined with the BCR-ABL1 reconstruction (Supplemental Fig. S6d; overlap of segment
295  20). From the AR reconstructions of the BCR-ABL1 amplicon and the co-existence of
296  BCR, ABL1 and GPC5 in overlapping locations, as shown by FISH (Fig. 3c 'Zoom'), AR
297  enabled us to hypothesize a potential sequence of events by which the fCNA formed. The
298  AR reconstructions support the formation of the BCR-ABL1 translocation (Supplemental
299  Fig. S8g;i-ii) followed by incorporation of chr13 regions (Supplemental Fig. S8g;iii-iv),
300  which subsequently undergo rearrangement (Supplemental Fig. S8g;v), and ultimately a
301  series of inverted repeats, possibly mediated through dicentrism (Supplemental Fig.
302  S8g;vi).

303

304  These results are consistent with previous reports that used cytogenetic approaches in
305  BCR-ABL1-positive samples to identify the presence of additional chromosomal
306  segments besides chr9 and chr22 involved in the Philadelphia chromosome[30,31]. AR
307  reconstructed the first base-pair resolved structures of the surrounding complex
308  rearrangement. The rearrangement of BCR-ABL1 and chr13 segments was followed by

10

309   additional duplications leading to a focal amplification. This example demonstrates the

310   utility of AR in resolving complex fCNAs, enhancing our understanding of the fundamental

311   mechanisms of cancer pathogenesis.

312

313   **AR enabled the first sequence-based reconstruction of a breakage-fusion-bridge**

314   The BFB mechanism of genomic amplification involves the loss of telomeres and

315   subsequent fusion of two sister chromatids[12,13]. In subsequent cellular division, the

316   asymmetric breaking of the fused dicentric chromosome structure results in one daughter

317   cell having an increased copy number of pieces of the previously fused chromosome. The

318   structure of various BFBs have been analyzed using cytogenetic techniques[14] and also

319   by computational models that predict a BFB mechanism based on copy number

320   counts[38,39]. Both methods are imprecise, to a degree, and may fail to capture the fine

321   structure of the BFB or handle imprecise copy number counts and/or additional structural

322   variants (SVs) inside the BFB. We deployed AR on the HCC827 lung cancer cell line

323   where we AA and cytogenetics previously suggested a BFB containing EGFR, though an

324   unambiguous structure was not identifiable[5,6].

325

326   We observed a banded pattern of EGFR and CEP7 (a chr7 centromeric D7Z1 repeat) in

327   a DNA FISH experiment on HCC827 cells, suggestive of a BFB mechanism (Fig. 4a). AA

328   generated a breakpoint graph of a 4.2 Mbp amplified region of chr7 containing EGFR

329   (Fig. 4b). The amplified BFB segments in the AA output ranged in size from 217 kbp to

330   1176 kbp. AR enabled the reconstruction of 16 unique OM scaffolds which, when

331   combined, covered the entirety of a BFB structure (Fig. 4c,d). The five most informative

332   single scaffolds ranged in size from 750 kbp to 2.3 Mbp, containing multiple junctions

333   which validate the order and orientation of the BFB breakpoint graph segments, resulting

334   in a 9.4 Mbp BFB structure, hereafter referred to as a BFB repeat unit. The BFB repeat

335   unit was amplified across the chromosome (Fig. 4a, e-f). AR also revealed a region

336   outside the AA amplicon, near the centromere of chr7, which explained the observed

337   EGFR and CEP7 repeat ("F"). In segment "B", we observed both a 600 bp deletion across

338   the entire BFB repeat unit and an 11 kbp inversion. The latter is labeled throughout Figure

339    4 with a black asterisk and only appears when segment "B" is duplicated and inverted,
340    suggesting that the SV arose during the formation of the BFB.

341

342    When the AR scaffolds were combined with the copy number data present in the
343    breakpoint graph, we identified a single BFB structure, that was consistent with the
344    theoretical BFB model of BFB formation[40]. A putative sequence of BFB cycles and
345    additional structural variation that results in the final BFB structure is shown in Fig. 4f (also
346    Supplemental Fig. S9a,b). Note that the copy number information and the theoretical
347    model together could not have reconstructed this BFB, as it contains heterogeneous
348    interior structural variants. We further validated the BFB patterning in HCC827 cells with
349    multi-FISH for segments "A", "C", and "D" from the BFB, using FISH (Fig. 4e,
350    Supplemental Fig. S9c). Together, these results on HCC827 show the power of AR as a
351    method to elucidate a complex mechanism of BFB-driven fCNA, even in the presence of
352    additional structural variant heterogeneity.

353

354    In addition to the EGFR-bearing amplicon, AA detected 5 other amplicons containing
355    MYC and NCOA2, among other oncogenes, in HCC827. The graphs were complex
356    (Supplemental Fig. S10a) and in many cases AA did not identify discordant edges
357    between distinctly amplified regions. Given the dearth of breakpoint edges, we combined
358    the amplicon breakpoint graphs for all six HCC827 amplicons and ran AR on the
359    combined graph, containing 555 segments. AR identified 206 contigs having alignments
360    to one or more graph segments. AR reconstructed multiple possible scaffolds and
361    captured overlapping subsets of amplicon regions from different graphs, suggestive of
362    possible heterogeneity. One scaffold showed NCOA2 located on a native region of chr8,
363    while another showed NCOA2 joined to MYC through a segment of chr21 (Supplemental
364    Fig. S10b,c).

365

366    **Other focal amplifications reconstructed by AR**
367    In breast cancer cell line T47D, where the AA breakpoint graph suggested amplification
368    of a 634 kbp region, AR reconstructed a 430 kbp segmental tandem duplication supported
369    by both AR and the AA breakpoint graph, containing oncogene GSE1 (Supplemental Fig.

370    S11a,b). This highlighted the ability of AR to also reconstruct classes of ultra-large, albeit

371    less-complex SVs.

372

373    In the renal cancer cell line, CAKI-2, AA generated a breakpoint graph spanning 12.0

374    Mbp, joining regions from chr3 and chr12 (Supplemental Fig. S11c,d). Despite the lower

375    overall copy number of this amplicon (~5), AR still reconstructed a 13.1 Mbp amplicon

376    explaining 99.9% of the amplified genomic content in the AA-detected fCNA. Both

377    amplicons for CAKI-2 and T47D appear to be intrachromosomal events given the AR

378    results.

379

380    Across the focal amplifications we studied in seven cancer cell lines, we reported 64

381    individual amplified breakpoints detected by both AA and validated by AR (Supplemental

382    Table S3). We also reported a summary of reconstruction findings for each sample and

383    provided a list of reconstructed paths in Supplemental Table S4. Taken together, our data

384    demonstrate the power of AR to combine NGS and OM data to elucidate a variety of

385    complex fCNAs commonly found in cancer - enabling a deeper understanding of the

386    fundamental mechanisms that give rise to fCNAs and promote cancer pathogenesis.

387

388    **Discussion:**

389    Revealing the architecture of fCNAs, particularly at a large scale, is critical to

390    understanding their functional implications. For instance, rearrangements present in

391    fCNAs can directly increase oncogene copy number, disrupt gene structure[41], and lead

392    to dysregulation of chromatin[17–19]. Thus, understanding the organization and content of

393    fCNAs is essential in predicting the behavior of the underlying sequences. Accurate

394    reconstruction of fCNA architecture can provide insights into the mechanisms of their

395    formation, leading to an improved understanding of the biological consequences of fCNA

396    that would not be available solely from methods characterizing individual breakpoints.

397

398    While previous methods have characterized complex structural variation using both OM

399    and NGS data[21,42], these methods have typically focused on the identification of individual

400    variants and breakpoints[38]. AR represents a more robust and comprehensive algorithmic

13

401  approach to reconstructing the fine architecture of a target fCNA. Indeed, while some of

402  the individual junctions reported by AR in these cell lines were already known[21], by

403  focusing on reconstructing entire amplicons through the propagation of breakpoint

404  information into larger scaffolds, AR provides a deeper insight into the complex

405  mechanisms that generate fCNA.

406

407  Genomic structural heterogeneity is problematic for any genome reconstruction, including

408  focal amplifications and the structure of fCNA in stable cell lines may evolve over time.

409  Despite the change in topology between linear HSR-like and circular ecDNA fCNAs, the

410  breakpoint graphs between both circular and linear forms of the same samples are highly

411  similar[6], suggesting ecDNA genomic structure is often not altered during reintegration.

412  We further note that assembled OM contigs may fail to capture rare instances of structural

413  heterogeneity in the genome. However, previous results suggest that focal amplifications

414  conferring a fitness advantage to cancer cells are clonally amplified[5,43], allowing for an

415  accurate reconstruction of the dominant structure.

416

417  AR produced a high-confidence reconstruction of the K562 BCR-ABL1 focal amplification,

418  which is thought to be derived at least in part from a balanced translocation. Despite the

419  presence of the AR-supported and FISH-validated HSR-like status of the BCR-ABL1

420  translocation in K562, there does not exist a completely validated model that explains the

421  increased copy number of BCR-ABL1 in one single location. We cannot rule out the

422  possibility that the BCR-ABL1 amplification in K562 is mediated through an ecDNA

423  stage[44], given the transient nature of the emergence and retreat of ecDNA[15] and the

424  highly rearranged genomic landscape surrounding BCR-ABL1.

425

426  The collection of paths reconstructed by AR represent possible reconstructions of the

427  fCNA, and the collection of paths may contain multiple similar explanations for the fCNA

428  architecture. This may be in part due to genomic heterogeneity, limitations of the optical

429  map assembly process, or errors in linking scaffolds across overlapping graph segments.

430  Furthermore, technological limitations related to the quality of OM assembly may affect

431  the ability to reconstruct high-fidelity amplicons. Thus, identifying a single best path or

432  collection of scaffolds which represent a reconstruction best explaining the breakpoint

433  graph and OM data still requires some manual identification and interpretation. From the

434  collection of output structures, AR does not automatically produce a prediction of the

435  mechanism of amplification.

436

437  We have not yet adapted AR to accept breakpoint graphs generated by other tools or to

438  accept breakpoint graphs derived from more balanced rearrangements - though the AR

439  algorithm is designed to handle them if an accurate breakpoint graph was provided as

440  input. Furthermore, recent advances in other long-range sequencing technologies[45]

441  highlight the need to adapt the AR algorithm to work with more general long-read

442  technologies – an aspect we plan to address in future development.

443

444  The accurate, multi-megabase scale, complex fCNAs reconstructed by AR not only

445  describe fine structural features of fCNA architecture, but also reveal mechanistic

446  signatures of fCNA formation, allowing for future interrogation of the relationship between

447  fCNA architecture and the biological consequences of fCNA structure. In particular,

448  methods to accurately characterize fCNAs will enable better classifications of cancer

449  subtypes and their associated prognoses.

450

451  **Methods**

452  **Cell culture**

453  NCI-H460, K562, and HCC827 cells were obtained from ATCC and cultured in RPMI-

454  1640 media supplemented with 10% FBS. HK301 cells were cultured as neural spheres

455  in DMEM/F12 media supplemented with B27, EGF (20 ng/ml), FGF (20 ng/ml), and

456  heparin (1 ug/ml). All cells were incubated under standard conditions.

457

458  **Metaphase chromosome spreads**

459  Metaphase cells were enriched by treating cells with Karyomax (Gibco) at a final

460  concentration of 0.1µg ml$^{-1}$. Cells were collected, washed in PBS, and resuspended in

461  75mM KCl for approximately 15 minutes at 37°C. Cells were fixed by addition of an equal

462  volume of Carnoy's fixative (3:1 methanol:glacial acetic acid). Cells were washed three

463  additional times in Carnoy's fixative and dropped onto humidified glass slides.

464

465  **FISH**

466  Metaphase spreads were equilibrated in 2x SSC (30mM sodium citrate, 300mM NaCl, pH

467  7) for approximately 5 minutes. They were dehydrated using successive washes of 75%,

468  85%, and 100% ethanol for two minutes each and allowed to dry. FISH probes were

469  diluted in hybridization buffer (Empire Genomics) and added to metaphase spreads on

470  slides, along with 22mm$^2$ coverslips. Samples were denatured at 70-75°C for 30 seconds

471  – 2 minutes. Probe hybridization was performed at 37°C for around 3 hours or overnight

472  in a humid and dark chamber. Samples were washed successively in 0.4x SSC and 2x

473  SSC with 0.1% Tween-20. Samples were incubated with DAPI (0.1µg ml$^{-1}$ in 2x SSC) for

474  10 minutes, then washed with 2x SSC and briefly rinsed with H$_2$O. Samples were

475  mounted with Prolong Gold, #1.5 coverslips, and sealed with nail polish.

476

477  **Microscopy**

478  Confocal microscopy was performed on a Leica SP8 Confocal microscope with white light

479  laser and Lightning deconvolution. Fluorescent microscope images were acquired using

480  an Olympus BX43 microscope with a QiClick cooled camera. Images were subsequently

481  analyzed in ImageJ[46] (using the Bio-Formats plugin[47]), to perform cropping, add scale

482  bars and perform global adjustments to image brightness.

483

484  **Acquisition of WGS data**

485  We previously published[5,6] WGS data to SRA for six of the seven cancer cell lines

486  (GBM39, NCI-H460, HCC827, HK301, K562, T47D) analyzed here. For CAKI-2, we used

487  WGS data published by the Cancer Cell Line Encyclopedia on SRA. A list of SRA

488  accession numbers used is available in Supplemental Table 1.

489

490  **Breakpoint graph generation**

491  WGS data was aligned to hg19 with BWA-MEM[48] (version 0.7.17-r1188, default

492  parameters) and the resulting alignments along with SNV calls produced by Freebayes[49]

16

493    (version v1.3.1-17-gaa2ace8) were supplied as input to the Canvas[50] CNV caller (version

494    1.39.0.1598). The alignments and CNV seeds were filtered using AmpliconArchitect's

495    amplified_intervals.py module. Seeds exceeding 40 kbp with copy number 5 were

496    subsequently analyzed with AmpliconArchitect. AmpliconArchitect outputs a breakpoint

497    graph encoding segmented CN calls and the discordant reads connecting the segments.

498    We note that in most cases identical amplicon regions are identified when CNV caller

499    ReadDepth[51] is used for seeding instead.

500

501    We standardized the breakpoint graph generation process into a workflow called

502    PrepareAA, available on Github: https://github.com/jluebeck/PrepareAA. We used the

503    default parameters specified by PrepareAA in this analysis. To produce *in silico* digestions

504    of breakpoint graph segments into reference optical maps, we used the

505    generate_cmap.py utility in AmpliconReconstructor. This method for *in silico* digestion

506    can produce labeling patterns for the Bionano Saphyr DLE-1 labeling pattern, while many

507    previous methods for *in silico* digestion do not.

508

509    **OM data generation**

510    High molecular weight (HMW) DNA was extracted from GBM39, HCC827, HK301, and

511    K562 cells using the Bionano Prep Blood and Cell Culture DNA Isolation Kit (Bionano

512    Genomics #80004), with minor modifications to recover good quality HMW gDNA. As

513    detailed below, the Nick, Label, Repair, and Stain (NLRS) and Direct Label and Stain

514    (DLS) reactions were carried out for the Bionano Irys and Saphyr platforms, respectively.

515    To generate the Irys data, DNA was nicked using Nt.BspQI nicking endonuclease (NEB),

516    followed by labeling, repairing, and staining, using the Bionano Prep NLRS DNA Labeling

517    Kit (Bionano Genomics #80001) along with recommended NEB reagents. To generate

518    the Saphyr data, DNA was labeled with DLE-1 enzyme, followed by proteinase digestion

519    and a membrane clean-up step, using the Bionano Prep DLS DNA Labeling Kit (#80005).

520    BspQI-labeled DNA was loaded onto the Irys Chip (Bionano Genomics #20249) and the

521    run conditions were manually optimized on the Irys system (Bionano Genomics #30047)

522    to ensure efficient DNA loading into the nanochannels. DLS-labeled DNA was loaded

523    onto a Saphyr Chip (Bionano Genomics #20319), and run conditions were automatically

524  optimized on the Saphyr system (Bionano Genomics #60239) using the Saphyr

525  Instrument Control Software to maximize DNA loading. Raw images generated by Irys

526  were processed into digital "Molecules" files using the Bionano software AutoDetect[25].

527  Images from the Saphyr system were processed into digital "Molecules" files via the

528  Saphyr Instrument Control Software. For Irys data, molecules ≥150 kilobase pairs (kbp)

529  were assembled into consensus genome maps using the Bionano Assembler[26,27] (version

530  5122), using default parameters; for Saphyr data, molecules ≥150 kbp were assembled

531  into maps using Bionano Access (version 1.2.1)[26]. Bionano Genomics separately

532  provided Saphyr OM data for cell lines K562, T47D, NCI-H460, and CAKI-2. The methods

533  by which OM data was generated for those four cell lines were previously published[21].

534

535  **Optical map contig alignments with SegAligner**

536  SegAligner uses a dynamic programming (DP) approach to optical map alignment, with

537  a recursion similar to previously proposed DP algorithms for OM alignment[52,53].

538  SegAligner scores OM alignments in a novel way which accounts for collapsed pairs of

539  labels in the assembled OM contig and uses an E-value approach to compute alignment

540  significance as method of controlling false alignments. We define label collapse as the

541  phenomenon where two nearby labels on an OM contig or map are measured as a single

542  label due to limitations of imaging[54].

543

544  SegAligner supports alignment of *in silico* digested segments of the reference genome

545  (including entire chromosomes of the reference genome) and assembled optical map

546  contigs. SegAligner supports models of error for data from both the Bionano Irys and

547  Bionano Saphyr instruments, and we parameterize our methods for them separately

548  (Supplemental Table 5). SegAligner also supports multiple modes of alignment including

549  semi-global, fitting, and overlap alignment.

550

551  To motivate the notion of an OM alignment, we first define the concept of an OM matching

552  region. Similarly to Valouev et al.,[53] a matching region is defined as the region between

553  and including two labels on a map. For example, $j$ and $i$ in Supplemental Fig. S1b

554  constitute a matching region with size $j - i$ and one unmatched label in-between. The

555    alignment score for two matching regions depends on the size discrepancy of the

556    matching regions and the number of unmatched labels in each matching region.

557

558    We define the following variables:

559    -    *b* is a sorted list of real numbers corresponding to the positions of labels on the

560         optical map contig in base pair units.

561    -    *x* is a sorted list of real numbers corresponding to the positions of labels on a single

562         *in silico* reference segment in base pair units.

563    -    *P* is a matrix storing backtracking references

564    -    *U* is a set storing reference segment label to contig label pairings which have

565         already been used in previous iterations of the alignment process.

566    -    *d* is the width of the band to consider for a banded alignment (default 6).

567    -    *M* is a map which relates each label, *j* on a genomic segment, *x*, to the estimated

568         probabilities for the left neighbor and right neighbor of *j,* that *j* and a neighbor would

569         be observed as a single label (i.e. "collapse").

570

571    Next, define $S[j][q]$ as the best score of aligning a subsequence of the first *j* labels on *b*

572    with a subsequence of the first *q* labels on segment *x*, where *j* and *q* are included in the

573    subsequences. Given two labels on the assembled contig *i, j*, and two labels on the

574    reference genome segment *p, q* where *i < j,* and *p < q,* The DP recurrence used by

575    SegAligner is (Algorithm 1)

576

$$S[j][q] = max_{\begin{pmatrix} max(0,j-d) \leq i < q \\ max(0,q-d) \leq p < d \end{pmatrix}} \{S[i][p] + Score(i,j,p,q)\}$$

578

579    Where *Score* is the SegAligner scoring function for two OM matching regions. *Score*

580    includes a function which computes the number of expected reference labels between p

581    and q after accounting for label collapse. A backtracking matrix P is used to record the

582    decision made in filling each cell S[j][q].  The DP Algorithm has complexity $O(mnd^2)$

583    where $m = |b|,\ n = |x|$ and $d$ is the width of the band. Backtracking is performed in $O(m)$

584    steps by backtracking through the coordinates stored in *P*. We find a most-likely path by

585  initializing the backtracking at $\mathrm{argmax}_{j,q} S[j][q]$ or $S[|b| - 1][|x| - 1]$ for fitting alignment.

586  Values used to parameterize the scoring function and label collapse map generation

587  function given below are provided in Supplemental Table S1.

588

589  Algorithm 2: SegAligner scoring function

```
590  function Score(b,x,i,j,p,q,M):
591      fₙ = c*(j - (i + 1))
592      e_ref = M(p,q)
593      f_p = c*e_ref
594      Δ = (abs((b[j] - b[i]) - (x[q] - x[p])))ᵏ
595      return 2c - (fₙ + f_p + Δ)
```

596

597  As multiple regions of a long OM query might match similar regions of the reference, we

598  extend the DP by masking out the best alignment path from the DP scoring matrix and

599  recomputing the next best alignment.

600

601  Labels within approximately 2000 bp on an OM molecule may be read as a single label

602  due to limitations of imaging, with increasing probability for smaller label-to-label intervals

603  (Supplemental Fig. S1c). SegAligner captures that behavior in its scoring method, by

604  precomputing the number of expected labels appearing in a collapsed label-set, given the

605  reference.

606

607  To compute probabilities of label collapse, we assume a model in which the probability

608  that a label at position $r$ has merged with its right neighbor at position $s$ is given by

609  $P(r \to s) = \min\left(1, \left(\frac{(s-r)^t}{w^t}\right)\right)$. The map $M$, encoding the expected number of uncollapsed

610  labels between two points on an *in silico* reference segment, is generated iteratively, by

611  evaluating the following sum. $M(p,q)$ represents the sum of probabilities for each label

612  between, but not including $p$ and $q$ that the label has collapsed with a neighbor. The sum

613  of probabilities for [0,1] binary random variables to be 1 naturally gives the expected value

614  of the sum of the binary random variables.

615

616

$$M(p,q) = \begin{cases} \sum_{p<k<q} \left(1 - \min\left(1, \frac{(x[k] - x[k-1])^t}{w^t}\right)\right)\left(1 - \min\left(1, \frac{(x[k+1] - x[k])^t}{w^t}\right)\right) & \text{if } x[q] - x[p] \geq \eta \\ 0 \text{ if } x[q] - x[p] < \eta \end{cases}$$

617

618  A genomic segment may appear multiple times in an optical map contig.

619  Parameterizations of $w, t$ and $\eta$ are parameterized separately depending on the Bionano

620  instrument used (Supplemental Table 5). SegAligner uses a set ($U$) to keep track of the

621  pairings of segment labels ($q$) and reference labels ($j$) which form each significant high-

622  scoring alignment. After a best-scoring alignment is found, the label pairings ($j,q$) are

623  added to $U$. Subsequent alignments of that segment cannot re-use any pairings in $U$. This

624  limits the creation of many nearly identical local alignments which differ by small indels,

625  only one of which (the best scoring) is useful from a practical standpoint. We also placed

626  a threshold on the number of times a single segment can be aligned to a single contig,

627  so that low-complexity segments do not cause the aligner to stall (default 12).

628

629  **Identifying significant high-scoring alignments**

630  To compute statistically significant alignments, SegAligner uses a strategy similar to

631  BLAST[55]. For each reference segment, $r$, SegAligner constructs a distribution of

632  alignment scores representing the best scoring alignments of $r$ to all contigs

633  (Supplemental Fig. S1b). As this distribution may contain true alignments between $r$ and

634  one or more contigs, violating the random pairing assumption of the E-value model,

635  SegAligner removes the highest 25 values from the distribution. From the remaining

636  distribution of scores, we define a set of high scoring segment pairs (HSPs) which are the

637  distribution of scores from the 85th percentile and up, from which SegAligner estimates

638  parameters in the E-value model. We note that this region of the HSP scoring distribution

639  tends to behave linearly (Supplemental Fig. S1c), allowing for a linear regression

640  approach to parameter estimation.

641

642  SegAligner assigns an empirical E-value for each element in the sorted distribution of

643  HSP alignment scores based on its rank (highest scoring having E-value 1). SegAligner

644  then performs a local linear regression to estimate unknown variables in the E-value

645  model. Generally, the E-value model is given by

21

646
$$E = Kmn_r e^{-\lambda S}$$

647 which implies

648
$$\log(E) = \log(Kmn_r) - \lambda S$$

649

650 where $m$ is the size of the combined collection of contig labels, $n_r$ is the number of labels

651 on the reference segment, and $S$ is the alignment score. As $K$ and $\lambda$ are unknown and

652 represent the intercept and slope, respectively, SegAligner determines them from the

653 empirical distribution of scores and E values using linear regression.

654

655 With all parameters known, the number of random high-scoring alignments, $a$, with score

656 $\geq S$ is given by a Poisson distribution

657
$$P(a) - \frac{e^{-a} E^a}{a!}$$

658 This implies that finding at least one HSP for a given value of $E$ is

659
$$P = 1 - e^{-E}$$

660 Thus, the score-cutoff $S_r^*$ corresponding to a given probability, $P$, for segment $r$, is

661
$$S_r^* = \frac{-\log\left(-\frac{\log(1-P)}{Kmn_r}\right)}{\lambda}$$

662 SegAligner assigns to each reference segment a score which corresponds to the p-value

663 cutoff for alignment significance. Default p-values are; $10^{-4}$ for semi-global alignment,

664 $10^{-6}$ for overlapping alignment, and $10^{-9}$ for detection of new genomic reference segments

665 aligning to contigs where the reference segment is not specified in the provided

666 breakpoint graph segments (detection mode). SegAligner also computes the mean and

667 median of segment-contig label pair alignment scores for each alignment exceeding the

668 significance thresholds. Statistically significant scoring alignments failing mean and

669 median thresholds (Supplemental Table 2) are filtered out. By default, AR attempts to

670 align graph segments with at least 10 (Irys) or 12 (Saphyr) labels in the segment.

671 However, the fitting mode of alignment only requires two endpoint labels, and so it is used

672 in the path imputation step in AR.

673

674 **AmpliconReconstructor – ARAlignDetect module**

675 AmpliconReconstructor coordinates the alignment of in-silico digested breakpoint graph

676 segments to optical map contigs using SegAligner (Fig. 1b). Alternately, AR can take as

677 input XMAP-formatted alignments produced by other alignment tools. If OM contigs with

678 alignments to graph segments contain unaligned regions with between 20 and 500

679 unmatched labels, and 200 kbp to 5 Mbp in length, those regions are extracted and

680 searched against the reference genome. ARAlignDetect calls SegAligner in the

681 "detection" mode, which then aligns the extracted unaligned region of the contig(s) to the

682 specified reference genome. If significant alignments are found between unaligned

683 regions of the contig and chromosomal segments in the reference, those segments are

684 extracted, and their identity is added to the collection breakpoint graph segments. Finally,

685 a new breakpoint graph is output containing the newly detected segments.

686

687 **Reconstructing amplicon paths with AmpliconReconstructor**

688 Optical map alignments of segments with contigs are converted into a scaffold, which we

689 define as a collection of alignments where the genomic distance between each pair of

690 alignment endpoints is known. AR represents the scaffolded alignments as a directed

691 acyclic graph (DAG), where the nodes are an abstract representation of each OM

692 alignment. Directed edges connect adjacent alignment endpoints. Overlapping

693 alignments are connected with special directed edges referred to as "forbidden" edges

694 (Fig. 1h). Two nodes are only connected by a non-forbidden edge if the right endpoint of

695 the source node has one or fewer labels of overlap with the left endpoint of the destination

696 node. Each contig with at least one alignment to a graph segment will comprise an

697 individual scaffold.

698

699 **Imputing paths in the scaffold with AmpliconReconstructor**

700 Some segments in the breakpoint graph may be too short to be uniquely aligned to an

701 OM contig. AR attempts to impute corrected paths in the scaffold using the structure of

702 the breakpoint graph. For every non-forbidden edge in the scaffold graph with a gap size

703 less than 400 kbp, AR identifies breakpoint graph nodes corresponding to the source and

704 destination endpoints, which we will denote as $s$, and $t$. AR then uses a constrained depth-

705 first search (DFS) strategy to identify paths in the breakpoint graph between $s$ and $t$.

23

706 Finding all possible paths between two nodes may produce infinitely many solutions
707 should a cycle exist between the two nodes, so the recursion is constrained to terminate
708 if certain conditions are reached. The constraints used in the search procedure are:

709

710 1) The multiplicity of the segments in the candidate path must always remain less
711 than or equal to the copy number of the segment as specified in the breakpoint
712 graph.
713 2) If a candidate path reaches the destination vertex, its length in base-pair units must
714 not be more than $\min(25000,10000L_p)$ shorter than the distance between the
715 source and destination vertices as expected given the scaffold backbone, where
716 $L_p$ is the length of the path in number of segments.
717 3) During path construction, the length of a candidate path must not exceed
718 $\min(25000,10000L_p)$ beyond the of the expected distance given the scaffold
719 backbone.
720 4) The number of valid candidate paths connecting source to destination must not
721 exceed $2^{10}$.
722 5) The path may not form a trivial cycle from ultra-short breakpoint graph segments
723 less than 100 bp long. Such cycles appearing in an NGS-derived breakpoint graph
724 we assumed to be erroneous or artifactual.

725

726 As constraint #4 may cause failure of the DFS whereby a tractable number of paths is not
727 found, AR implements a constrained BFS search as a fallback option, which is used when
728 the DFS fails for that reason. By parsimony, shorter paths between two nodes are more
729 likely to be correct, thus AR applies the same set of criteria for the BFS search, with the
730 threshold in constraint #4 increased to $2^{16}$.

731

732 All valid candidate imputation paths discovered by AR are scored by a fitting alignment
733 procedure using SegAligner. To score a candidate path, the ordered path segments, as
734 well as the first and last labels on the source and destination endpoints, are converted to
735 a compound CMAP composed of the concatenated CMAPs of the individual segments.
736 A fitting alignment is performed between the compound CMAP and the region of the

24

737  contig between the alignment endpoints, using SegAligner. The path with the alignment

738  score which most improves the junction score is kept. If no valid candidate path improves

739  the score of the junction, it remains unimputed. The scaffold is then updated to contain

740  the imputed breakpoint graph path.

741

742  **Identifying linked scaffold paths with AR**

743  Given the collection of scaffold DAGs, AR first searches for paths in the individual DAGs

744  which represent "heaviest" paths in the scaffold DAG, where the weight of a path is the

745  sum of the lengths of its segments in base pairs. AR stores the heaviest path(s) for each

746  scaffold prior to performing scaffold linking.

747

748  AR leverages the two orthogonal sources of information encoded in the breakpoint graph

749  and OM contigs to link individual scaffolds. As the breakpoint graph segments are not

750  detected to contain interior breakpoints, two endpoint alignments of the same breakpoint

751  graph segment may be linked across two contigs. AR searches for prefix paths and suffix

752  paths in each DAG. From the collection of prefixes and suffixes, AR searches for overlap

753  between scaffolds generated from different contigs. Given that a contig can be assembled

754  in either direction, overlapping reverse oriented suffixes or prefixes can also be matched.

755  AR exhaustively finds sub-paths hitting either end of a scaffold DAG, which have overlap

756  with other endpoint sub-paths, where the endpoint sequence of the scaffold may be

757  assembled in either direction.

758

759  **Finding reconstructions in the linked scaffold graph**

760  Given the graph of linked scaffolds, AR searches for paths in the graph which conform to

761  the copy number ratios in the breakpoint graph. AR starts by searching for all paths in the

762  graph which begin at endpoint nodes in the individual scaffolds. AR then uses a greedy

763  approach to identify the longest unique paths which conform to the copy number

764  restrictions. From the candidate paths, AR checks each path segment's multiplicity

765  against the copy numbers encoded in the breakpoint graph in a ratio-dependent manner.

766

767  AR iterates over all the segment multiplicities in the reconstructed path, and at each
768  multiplicity level determines the maximum estimated genomic copy number of path
769  segments with that multiplicity. If a path segment has a multiplicity that is greater than the
770  genomic copy number of that segment divided by the maximum copy number of all
771  segments with multiplicities less than the given segment, then the path violates the copy
772  number ratio check. AR allows each segment in the reconstructed path to exceed by 1
773  copy the copy number expected given the ratio between breakpoint graph copy numbers
774  and segment multiplicity. If $n_p$ is the multiplicity of segment $n$ in the candidate path, $P$, and
775  $n_g$ is the copy number of graph segment $n$ in the breakpoint graph, then $n_p$ must satisfy

776

777
$$n_p \leq \max\left(c, \frac{n_g}{m_g}\right) + 1, \qquad \forall\, n \in P$$

778  where

779
$$m_g = \max\left(i_g, \forall\, i \in P, i_p == c\right)$$

780
$$c \in \mathbb{Z}$$

781
$$n_p > c > 0$$

782

783  If a candidate path passes the copy number ratio check, it undergoes a pairwise
784  comparison with other paths passing this criterion, to check for path uniqueness. A path
785  is unique if it does not represent a subsequence of a previously identified unique path.
786  Furthermore, no rotation of the path sequence may be a subsequence of a previously
787  identified unique path. AR assess subsequence paths by computing a longest common
788  substring between a candidate path and a previously identified unique path (Algorithm 3).
789  As the paths are first sorted by total alignment score prior to the iterative approach, this
790  method is a greedy algorithm which prioritizes long, heavy paths as being more likely to
791  be identified as unique non-subsequence paths. AR categorizes paths as being cyclic if
792  the first and last scaffold graph node in the path are the same, and the path length is
793  greater than two, as this distinguishes cyclic paths from paths which appear cyclic such
794  as singleton paths or paths which represent segmental tandem duplications. Paths
795  reported by AR are output in the AmpliconArchitect "cycles" file format.

796

797     Algorithm 3: Greedy filtering of subsequence paths

798

```
799   Function FilterSubsequencePaths(sorted_paths):
800      kept = empty array
801      for P in sorted_paths do:
802         isSubsequence = False
803         for J in kept do:
804            for R in the set of all rotations of path P:
805               if R is a subsequence of J then:
806                  isSubsequence = True
807
808         if not isSubsequence then:
809            append P to kept
810      return kept
```

811

## Simulation of amplicons to measure AR performance

813     We used OMSim[56] (version 1.0) to simulate Bionano Irys OM data from the hg19

814     reference as well as from 85 non-trivial paths (i.e. not directly consistent with the reference

815     genome) in AA-generated breakpoint graphs from 25 cancer samples, including both

816     cyclic and non-cyclic breakpoint graph paths. OM molecules were simulated at 40x

817     baseline coverage for each chromosome arm in hg19. The combined hg19 maps from all

818     arms were assembled into a set of OM contigs using Bionano Assembler (version 5122).

819     A similar process was performed using high-confidence breakpoint graph paths, which

820     were converted to FASTA format and used for map simulation. For each simulated path,

821     molecules were simulated at a range of copy numbers, and simulated molecules from the

822     chromosome arm(s) (downsampled to the appropriate CN) from which the path segments

823     came were combined and *de* novo assembled into OM contigs with BioNano Assembler.

824     The resulting contigs from each amplicon simulation were combined with the previously

825     simulated reference contigs and used as input to AR. For combination sets of three

826     amplicons from the same sample, a similar downsampling and combination strategy was

827     used, where molecules from each of the three amplicon simulations was separately

828     downsampled based on the copy number settings of the mixture then combined. As

829    heterogeneous combinations of amplicons may occur at different ratios, we selected three

830    sets of copy numbers for this combination simulation cases: 20-20-20, 20-15-10, and 20-

831    2-2.

832

833    **Measuring AR simulation performance**

834    We computed the longest common substring (LCS) between the AR paths and the

835    ground-truth path and considered only the path having the LCS between AR and AA paths

836    when computing precision and recall. We define the LCS here using the identities of the

837    breakpoint graph segments and their orientations. We pre-filtered some possible

838    assembly error reflected in the paths by removing ends of reconstructed paths which were

839    trivial reconstructions of the reference genome and which were not supported by the AA

840    path. To measure the accuracy of AR-reconstructed paths against the ground truth

841    simulated paths, we developed a set of three measurements which were used in

842    calculating performance and recall.

843

844    1) Length (bp): Reports the length of a breakpoint graph path in base pair units.

845    2) Nsegs: Reports the length of a breakpoint graph path in terms of the number of

846        graph segments (unbiased towards genomic length)

847    3) Breakpoint: Reports the length of a breakpoint graph in terms of the number of

848        breakpoint graph segment junctions in the path.

849

850

851    We define precision and recall as follows, where $M$ is the path measurement function

852    (Length (bp), Nsegs, or Breakpoint), $LCS$ is the longest common substring function, $P_{AA}$

853    is the sequence of segments in the AA path, and $P_{AR}$ is the sequence of segments in the

854    reconstructed AR path:

855
$$Precision: \frac{M\big(LCS(P_{AA}, P_{AR})\big)}{M(P_{AR})}$$

856
$$Recall: \frac{M\big(LCS(P_{AA}, P_{AR})\big)}{M(P_{AA})}$$

857

858 To summarize the precision and recall metrics in a single value, we computed a mean F1

859 score across all the simulated amplicons for a given set of simulation conditions as

860
$$mean\ F1 = \frac{\sum_i \left(2\ \frac{precision_i * recall_i}{precision_i + recall_i}\right)}{n}$$

861

## Reconstructed path visualizations

863 We developed a visualization utility, CycleViz (https://github.com/jluebeck/CycleViz),

864 which produces circular and linear visualizations of AR or AA reconstructed amplicons

865 (Supplemental Fig. S2a,b), to create topologically correct visualizations of AR

866 reconstructions. CycleViz accepts inputs including the path files reported by AR (in the

867 AA "cycles" format) as well as the path OM alignment files (optional) and produces

868 visualizations which show the reconstructed path, *in silico* digestion of the path segments

869 and the alignments of the digested segments with assembled OM contigs. For circular

870 and linear visualizations, CycleViz places path segments in the visualization based on the

871 length of the segments and their position in the path. For circular visualization layouts,

872 the relative positions are converted to polar coordinates and a circular layout is formed.

873 We also developed a visualization utility for visualizing JSON-encoded scaffold graphs

874 formed by AR using CytoscapeJS (Supplemental Fig. S2c).

875

## Contributions:

877 J.L., V.B., and P.S.M conceived the work and designed the study. J.L. and V.B. developed

878 the AmpliconReconstructor algorithm and software. C.C. and D.A.P. generated Bionano

879 OM data and provided technical advice. J.L., S.R.D., and V.B. conceived and conducted

880 the simulation study. J.T.L and K.M.T. conducted FISH and microscopy experiments and

881 provided technical advice. V.B., P.S.M., and J.A.L. supervised all experiments. V.D., C.Z.,

882 and U.R. performed computational analysis and provided technical advice. J.L., C.C.,

883 J.T.L., K.M.T., V.D., D.A.P., J.A.L., P.S.M., and V.B. wrote the paper.

884

## Competing Interests:

886 P.S.M. and V.B. are co-founders of Boundless Bio, Inc. (BB), and serve as consultants.

887 K.T. is currently employed by and receives income from BB. V.B. is a co-founder, and

888    has equity interest in Digital Proteomics, LLC, and receives income from Digital

889    Proteomics (DP). D.A.P. is employed by and receives income from Bionano Genomics,

890    Inc. The terms of this arrangement have been reviewed and approved by the University

891    of California, San Diego in accordance with its conflict of interest policies. BB and DP

892    were not involved in the research presented here.

893

894    **Acknowledgements:**

895    The authors thank members of the Bafna and Mischel labs, as well as Dr. Marcy Erb

896    (UCSD SOM Microscopy Core) for advice on FISH experiments, the Ecker Lab at the

897    Salk Institute for use of the Bionano Irys optical mapping instrument, and Bionano

898    Genomics, Inc., (Alex Hastie, Jian Wang, Ernest Lam, Andy Pang, and others) for

899    supplying data and providing input on this project.

900

901    **Code Accessibility:**

902    The following tools are available online.

903      - AmpliconArchitect: https://github.com/virajbdeshpande/AmpliconArchitect

904      - PrepareAA: https://github.com/jluebeck/PrepareAA

905      - AmpliconReconstructor (& SegAligner):

906      https://github.com/jluebeck/AmpliconReconstructor

907      - CycleViz: https://github.com/jluebeck/CycleViz

908      - ScaffoldGraphViewer: https://github.com/jluebeck/ScaffoldGraphViewer

909

910

911

912

913

914

915

916

917

918

**References**

1.  Hanahan, D. & Weinberg, R. A. Hallmarks of cancer: The next generation. *Cell* **144**, 646–674 (2011).

2.  Bignell, G. R. *et al.* Signatures of mutation and selection in the cancer genome. *Nature* **463**, 893–898 (2010).

3.  Stuart, D. & Sellers, W. R. Linking somatic genetic alterations in cancer to therapeutics. *Current Opinion in Cell Biology* **21**, 304–310 (2009).

4.  Zack, T. I. *et al.* Pan-cancer patterns of somatic copy number alteration. *Nat. Genet.* **45**, 1134–1140 (2013).

5.  Turner, K. M. *et al.* Extrachromosomal oncogene amplification drives tumour evolution and genetic heterogeneity. *Nature* **543**, 122–125 (2017).

6.  Deshpande, V. *et al.* Exploring the landscape of focal amplifications in cancer using AmpliconArchitect. *Nat. Commun.* **10**, (2019).

7.  Carroll, S. M. *et al.* Double minute chromosomes can be produced from precursors derived from a chromosomal deletion. *Mol. Cell. Biol.* **8**, 1525–1533 (1988).

8.  Oobatake, Y. & Shimizu, N. Double-strand breakage in the extrachromosomal double minutes triggers their aggregation in the nucleus, micronucleation, and morphological transformation. *Genes, Chromosom. Cancer* gcc.22810 (2019). doi:10.1002/gcc.22810

9.  Stephens, P. J. *et al.* Massive genomic rearrangement acquired in a single catastrophic event during cancer development. *Cell* **144**, 27–40 (2011).

10. Jones, D. T. W. *et al.* Tandem duplication producing a novel oncogenic BRAF fusion gene defines the majority of pilocytic astrocytomas. *Cancer Res.* **68**, 8673–8677 (2008).

11. Menghi, F. *et al.* The Tandem Duplicator Phenotype Is a Prevalent Genome-Wide Cancer Configuration Driven by Distinct Gene Mutations. *Cancer Cell* **34**, 197-210.e5 (2018).

12. McClintock, B. The Stability of Broken Ends of Chromosomes in Zea Mays. *Genetics* **26**, 234–82 (1941).

13. Soler, D., Genescà, A., Arnedo, G., Egozcue, J. & Tusell, L. Telomere dysfunction

950   drives chromosomal instability in human mammary epithelial cells. *Genes*

951   *Chromosom. Cancer* **44**, 339–350 (2005).

952   14.   Kitada, K. & Yamasaki, T. The complicated copy number alterations in

953         chromosome 7 of a lung cancer cell line is explained by a model based on

954         repeated breakage-fusion-bridge cycles. *Cancer Genet. Cytogenet.* **185**, 11–9

955         (2008).

956   15.   Nathanson, D. A. *et al.* Targeted Therapy Resistance Mediated by Dynamic

957         Regulation of Extrachromosomal Mutant EGFR DNA. *Science (80-. ).* **343**, 72–76

958         (2014).

959   16.   Verhaak, R. G. W., Bafna, V. & Mischel, P. S. Extrachromosomal oncogene

960         amplification in tumour pathogenesis and evolution. *Nature Reviews Cancer* **19**,

961         283–288 (2019).

962   17.   Wu, S. *et al.* Circular ecDNA promotes accessible chromatin and high oncogene

963         expression. *Nature* (2019). doi:10.1038/s41586-019-1763-5

964   18.   Morton, A. R. *et al.* Functional Enhancers Shape Extrachromosomal Oncogene

965         Amplifications. *Cell* (2019). doi:10.1016/j.cell.2019.10.039

966   19.   Mitsuda, S. H. & Shimizu, N. Epigenetic Repeat-Induced Gene Silencing in the

967         Chromosomal and Extrachromosomal Contexts in Human Cells. *PLoS One* **11**,

968         (2016).

969   20.   Kosugi, S. *et al.* Comprehensive evaluation of structural variation detection

970         algorithms for whole genome sequencing. *Genome Biol.* **20**, (2019).

971   21.   Dixon, J. R. *et al.* Integrative detection and analysis of structural variation in

972         cancer genomes. *Nat. Genet.* **50**, 1388–1398 (2018).

973   22.   Layer, R. M., Chiang, C., Quinlan, A. R. & Hall, I. M. LUMPY: A probabilistic

974         framework for structural variant discovery. *Genome Biol.* **15**, (2014).

975   23.   Sedlazeck, F. J. *et al.* Accurate detection of complex structural variations using

976         single-molecule sequencing. *Nat. Methods* **15**, 461–468 (2018).

977   24.   Dzamba, M. *et al.* Identification of complex genomic rearrangements in cancers

978         using CouGaR. *Genome Res.* **27**, 107–117 (2017).

979   25.   Cao, H. *et al.* Rapid detection of structural variation in a human genome using

980         nanochannel-based genome mapping technology. *Gigascience* **3**, 34 (2014).

981    26.   Software Downloads - Bionano Genomics. Available at:

982          https://bionanogenomics.com/support/software-downloads/. (Accessed: 14th

983          November 2019)

984    27.   Anantharaman, T., Mishra, B. & Schwartz, D. Genomics via optical mapping. III:

985          Contiging genomic DNA. *Proceedings. Int. Conf. Intell. Syst. Mol. Biol.* 18–27

986          (1999).

987    28.   Leung, A. K.-Y. *et al.* OMBlast: alignment tool for optical mapping using a seed-

988          and-extend approach. *Bioinformatics* btw620 (2016).

989          doi:10.1093/bioinformatics/btw620

990    29.   Anantharaman, T. S., Mishra, B. & Schwartz, D. C. Genomics via Optical Mapping

991          II: Ordered Restriction Maps. *J. Comput. Biol.* **4**, 91–118 (1997).

992    30.   Barr, L. F. *et al. c-Myc Suppresses the Tumorigenicity of Lung Cancer Cells and

993          Down-Regulates Vascular Endothelial Growth Factor Expression 1. CANCER

994          RESEARCH* **60**, (2000).

995    31.   Cho, S. W. *et al.* Promoter of lncRNA Gene PVT1 Is a Tumor-Suppressor DNA

996          Boundary Element. *Cell* **173**, 1398-1412.e22 (2018).

997    32.   Vogt, N. *et al.* Amplicon rearrangements during the extrachromosomal and

998          intrachromosomal amplification process in a glioma. *Nucleic Acids Res.* **42**,

999          13194–13205 (2014).

1000   33.   Storlazzi, C. T. *et al.* Gene amplification as doubleminutes or homogeneously

1001          staining regions in solid tumors: Origin and structure. *Genome Res.* **20**, 1198–

1002          1206 (2010).

1003   34.   Rowley, J. D. A new consistent chromosomal abnormality in chronic myelogenous

1004          leukaemia identified by quinacrine fluorescence and Giemsa staining. *Nature* **243**,

1005          290–293 (1973).

1006   35.   Virgili, A. & Nacheva, E. P. Genomic amplification of BCR/ABL1 and a region

1007          downstream of ABL1 in chronic myeloid leukaemia: A FISH mapping study of

1008          CML patients and cell lines. *Mol. Cytogenet.* **3**, (2010).

1009   36.   Chandran, R. K. *et al.* Genomic amplification of BCR-ABL1 fusion gene and its

1010          impact on the disease progression mechanism in patients with chronic

1011          myelogenous leukemia. *Gene* **686**, 85–91 (2019).

1012   37.   Grosveld, G. *et al.* The chronic myelocytic cell line K562 contains a breakpoint in bcr and produces a chimeric bcr/c-abl transcript. *Mol. Cell. Biol.* **6**, 607–616 (1986).

1015   38.   Zakov, S., Kinsella, M. & Bafna, V. An algorithmic approach for breakage-fusion-bridge detection in tumor genomes. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 5546–51 (2013).

1018   39.   Zakov, S. & Bafna, V. Reconstructing Breakage Fusion Bridge Architectures Using Noisy Copy Numbers. *J. Comput. Biol.* **22**, 577–594 (2015).

1020   40.   Kinsella, M. & Bafna, V. Combinatorics of the breakage-fusion-bridge mechanism. *J. Comput. Biol.* **19**, 662–678 (2012).

1022   41.   Koche, R. P. *et al.* Extrachromosomal circular DNA drives oncogenic genome remodeling in neuroblastoma. *Nat. Genet.* (2019). doi:10.1038/s41588-019-0547-z

1025   42.   Chan, E. K. F. *et al.* Optical mapping reveals a higher level of genomic architecture of chained fusions in cancer. *Genome Res.* **28**, 726–738 (2018).

1027   43.   Decarvalho, A. C. *et al.* Discordant inheritance of chromosomal and extrachromosomal DNA elements contributes to dynamic disease evolution in glioblastoma. *Nat. Genet.* **50**, 708–717 (2018).

1030   44.   Morel, F. *et al.* Double minutes containing amplified bcr-abl fusion gene in a case of chronic myeloid leukemia treated by imatinib. *Eur. J. Haematol.* **70**, 235–9 (2003).

1033   45.   Jain, M. *et al.* Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nat. Biotechnol.* **36**, 338–345 (2018).

1035   46.   Downloads - ImageJ. Available at: https://imagej.net/Downloads. (Accessed: 9th December 2019)

1037   47.   Linkert, M. *et al.* Metadata matters: Access to image data in the real world. *Journal of Cell Biology* **189**, 777–782 (2010).

1039   48.   Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. (2013).

1041   49.   Garrison, E. & Marth, G. Haplotype-based variant detection from short-read sequencing. (2012).

1043  50.  Roller, E., Ivakhno, S., Lee, S., Royce, T. & Tanner, S. Canvas: Versatile and
1044       scalable detection of copy number variants. *Bioinformatics* **32**, 2375–2377 (2016).
1045  51.  Miller, C. A., Hampton, O., Coarfa, C. & Milosavljevic, A. ReadDepth: A parallel R
1046       package for detecting copy number alterations from short sequencing reads.
1047       *PLoS One* **6**, (2011).
1048  52.  Huang, X. & Waterman ', M. S. Dynamic programming algorithms for restriction
1049       map comparison. **8**, 1–520 (1992).
1050  53.  Valouev, A. *et al.* Alignment of Optical Maps. *J. Comput. Biol.* **13**, 442–462
1051       (2006).
1052  54.  Das, S. K. *et al.* Single molecule linear analysis of DNA in nano-channel labeled
1053       with sequence specific fluorescent probes. *Nucleic Acids Res.* **38**, (2010).
1054  55.  Karlin, S. & Altschul, S. F. Methods for assessing the statistical significance of
1055       molecular sequence features by using general scoring schemes. *Proc. Natl. Acad.*
1056       *Sci. U. S. A.* **87**, 2264–8 (1990).
1057  56.  Miclotte, G. *et al.* OMSim: a simulator for optical map data. *Bioinformatics* **33**,
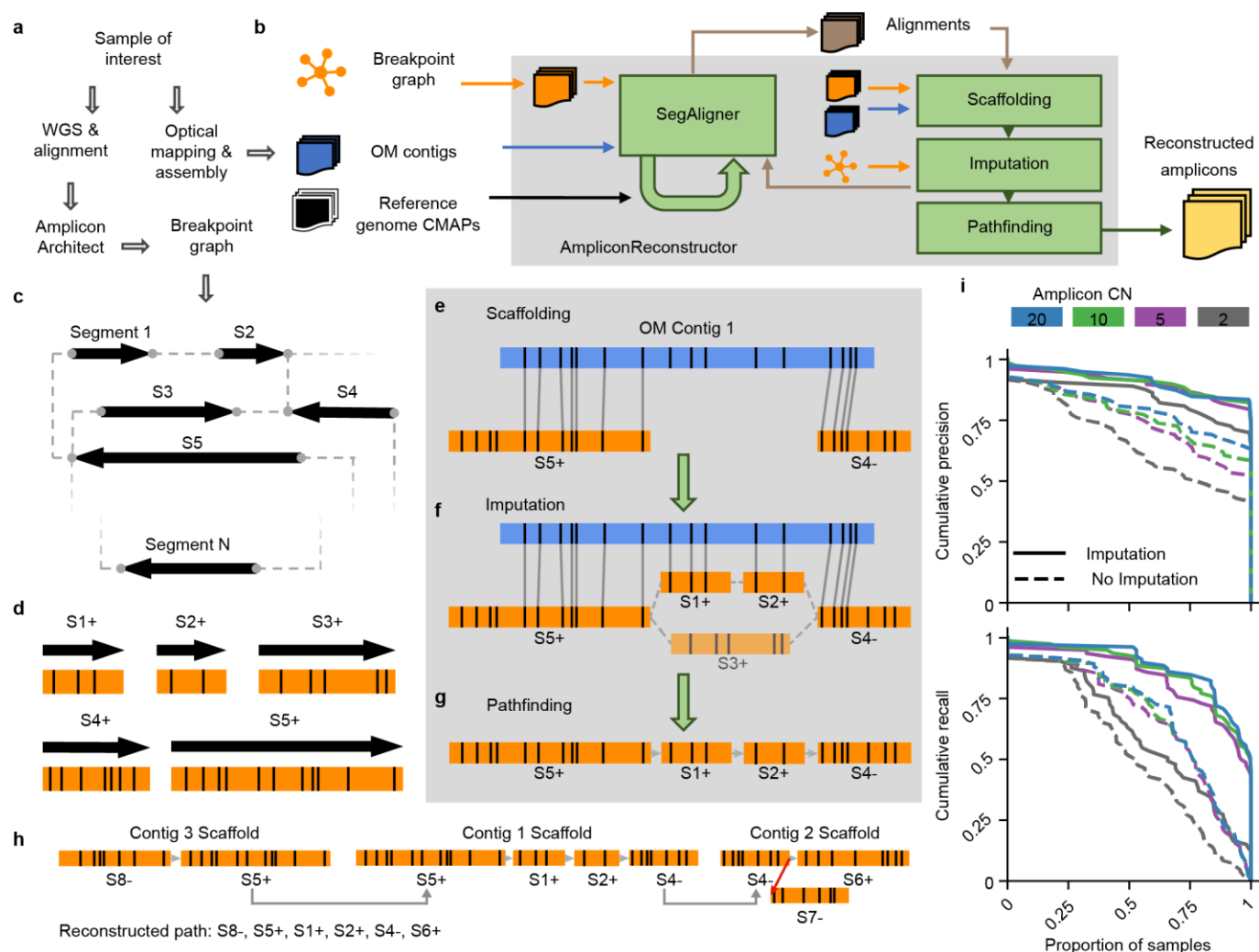1058       2740–2742 (2017).
1059

Figure 1: **[AmpliconReconstructor (AR) overview.] a**, Workflow to produce the necessary inputs for AR. AR accepts OM data in the consensus map (CMAP) format. **b**, High-level overview of the AR method, where the inputs and outputs are shown outside the grey box representing the AR wrapper. The green loop-back arrow on the SegAligner module represents the identification of reference segments not encoded in the breakpoint graph. **c**, A breakpoint graph with *N* segments. **d**, *In silico* digestion of breakpoint graph segments (orientation given by +/-) from **c** to produce graph OM segments. **e**, Alignment of graph OM segments to OM contigs produces a scaffold of segment-contig alignments. **f**, AR uses the structure of the breakpoint graph to identify paths between scaffold alignment endpoints which are also paths in the breakpoint graph. AR generates composite optical maps from combined path segments to score each candidate path against the gap in the scaffold. **g**, AR identifies a candidate path with maximum score out of the possible imputed paths between two alignments. **h**, AR links individual scaffolds sharing overlap between graph segments. The resulting graph has two types of edges, allowed (grey) and forbidden (red). **i**, Cumulative precision and recall curves based on simulated OM data for AR using SegAligner, calculated with the Length (bp) LCS metric. Line color indicates the copy number (CN) of the simulated amplicon.
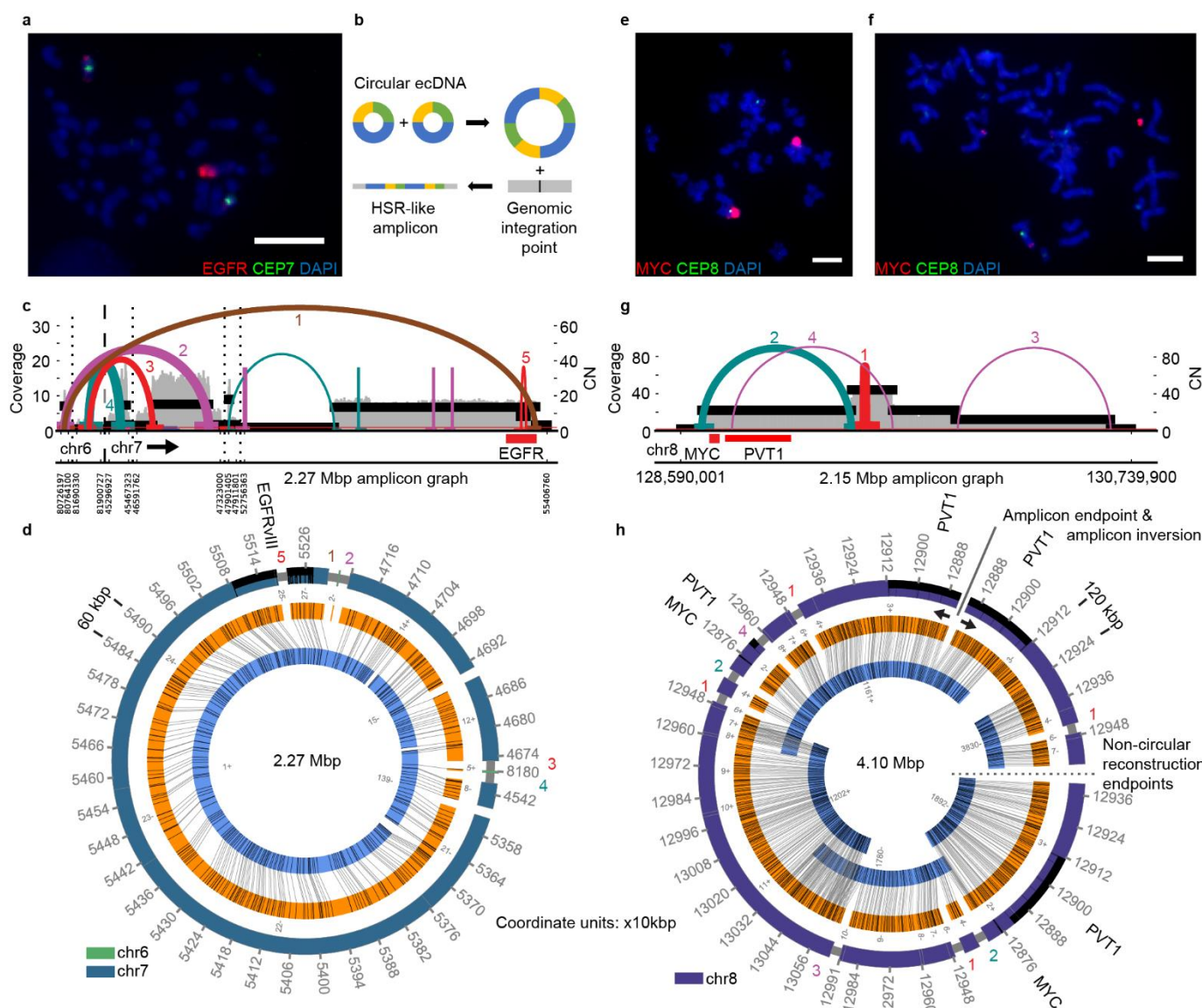
Figure 2: **[Reconstruction of extrachromosomal DNA (ecDNA)] a**, FISH with DAPI (4′,6-diamidino-2-phenylindole)-stained metaphase chromosomes in HK301 showing an HSR-like amplicon containing EGFR. Scale bar indicates 10 µm. **b**, Theoretical model for the integration of circular extrachromosomal DNA into HSR-like amplicons, preserving the structure of breakpoint graph. **c**, AA-generated breakpoint graph for HK301 containing EGFR and also segments from chr6. The coloring of the graph edges represents the orientation of the junction between the two segments. Edge thickness indicates AA-estimated breakpoint copy number. Vertical dashed lines separate segments from different chromosomes while dotted lines indicate distinct genomic regions from the same chromosome. Numbering of breakpoint edges corresponds with AR reconstruction breakpoint numbering. **d**, Cyclic AR reconstruction of HK301 amplicon containing EGFRvIII. Breakpoint graph edges supported by the AA graph are numbered in a manner corresponding to the numbering in panel **c**. **e**, FISH with DAPI-stained metaphase chromosomes in NCI-H460 shows HSR-like MYC amplicon. Scale bar indicates 7.3 µm. **f**, FISH with DAPI-stained metaphase chromosomes in NCI-H460 showing extrachromosomal MYC amplicon. Scale bar indicates 7.3 µm. **g**, AA-generated breakpoint graph for

NCI-H460 containing MYC and PVT1. **h**, AR reconstruction of the NCI-H460 amplicon. Indicated in this figure is an amplicon inversion point (top right) where the reconstruction explaining the full amplicon ends, and then the structure begins to repeat in the opposite direction (solid line & opposing black arrows). Also indicated is an endpoint for the non-circular reconstruction (center right) where the AR reconstruction and full amplicon structure both stop (dotted line).
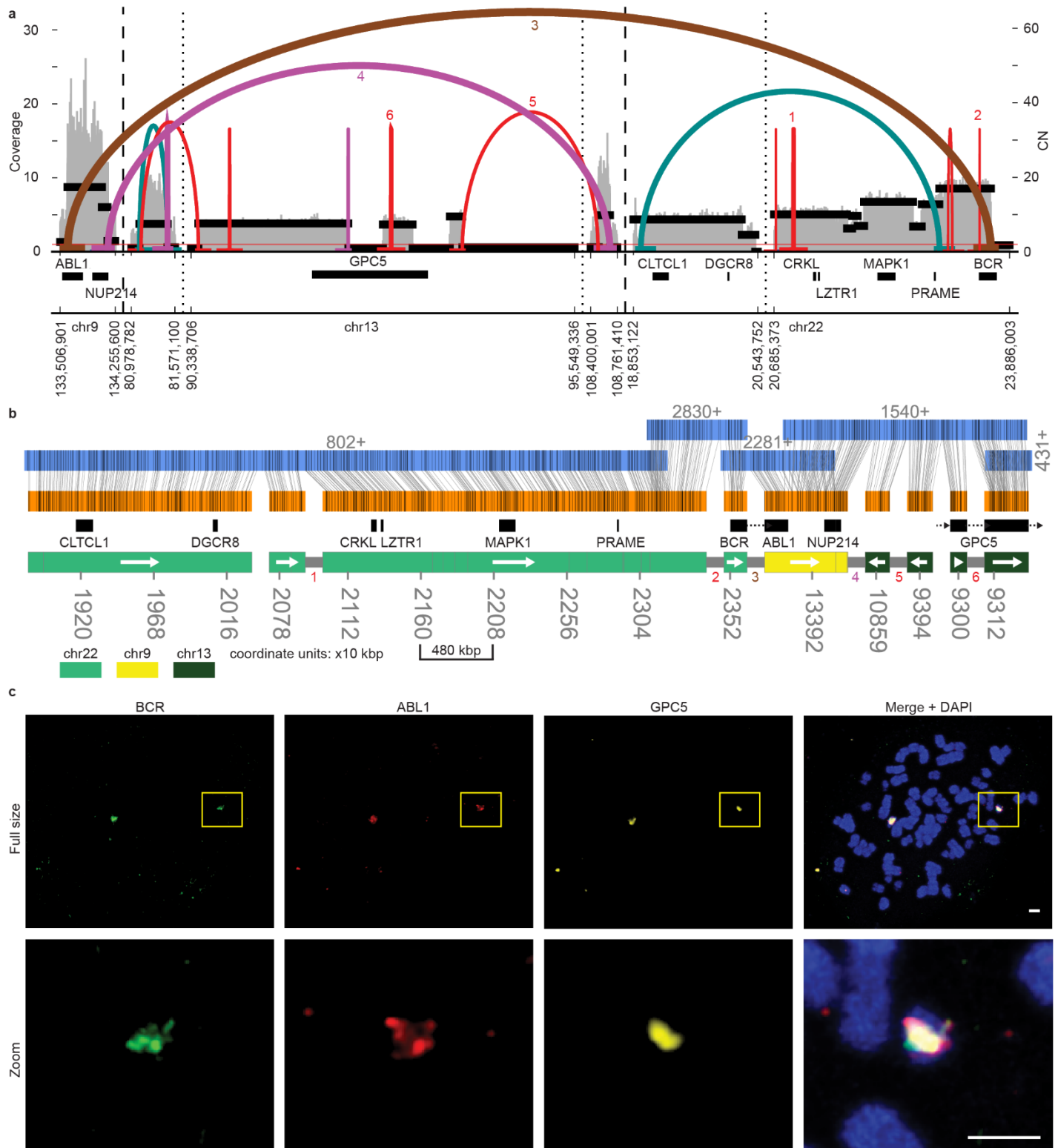
Figure 3: [**Reconstruction of a complex Philadelphia chromosome] a**, AA-generated breakpoint graph for K562. Estimated copy number (CN), coverage, discordant reads forming breakpoint graph edges, and a subset of the genes in these regions are shown. **b**, AR reconstruction of an 8.5 Mbp focal amplification which was supported by both Irys and Saphyr reconstructions. The tracks from top to bottom are: OM contigs (with contig ID and direction indicated above), graph segments (alignments shown with vertical grey lines), gene subset and color-coded reference genome bar with genomic coordinates (scaled as 10 kbp units). Grey half-height bars between individual segments on the reference genome bar indicate support from edges in the AA breakpoint graph. White arrows inside the chromosome color bar indicate direction of genomic segment(s). Colored numbers correspond to numbered breakpoint graph edges in panel **a**. **c**, Multi-FISH using probes against BCR, ABL1 and GPC5 with DAPI-stained metaphase chromosomes. Scale bars indicate 2 μm in both "Full size" and "Zoom" rows.
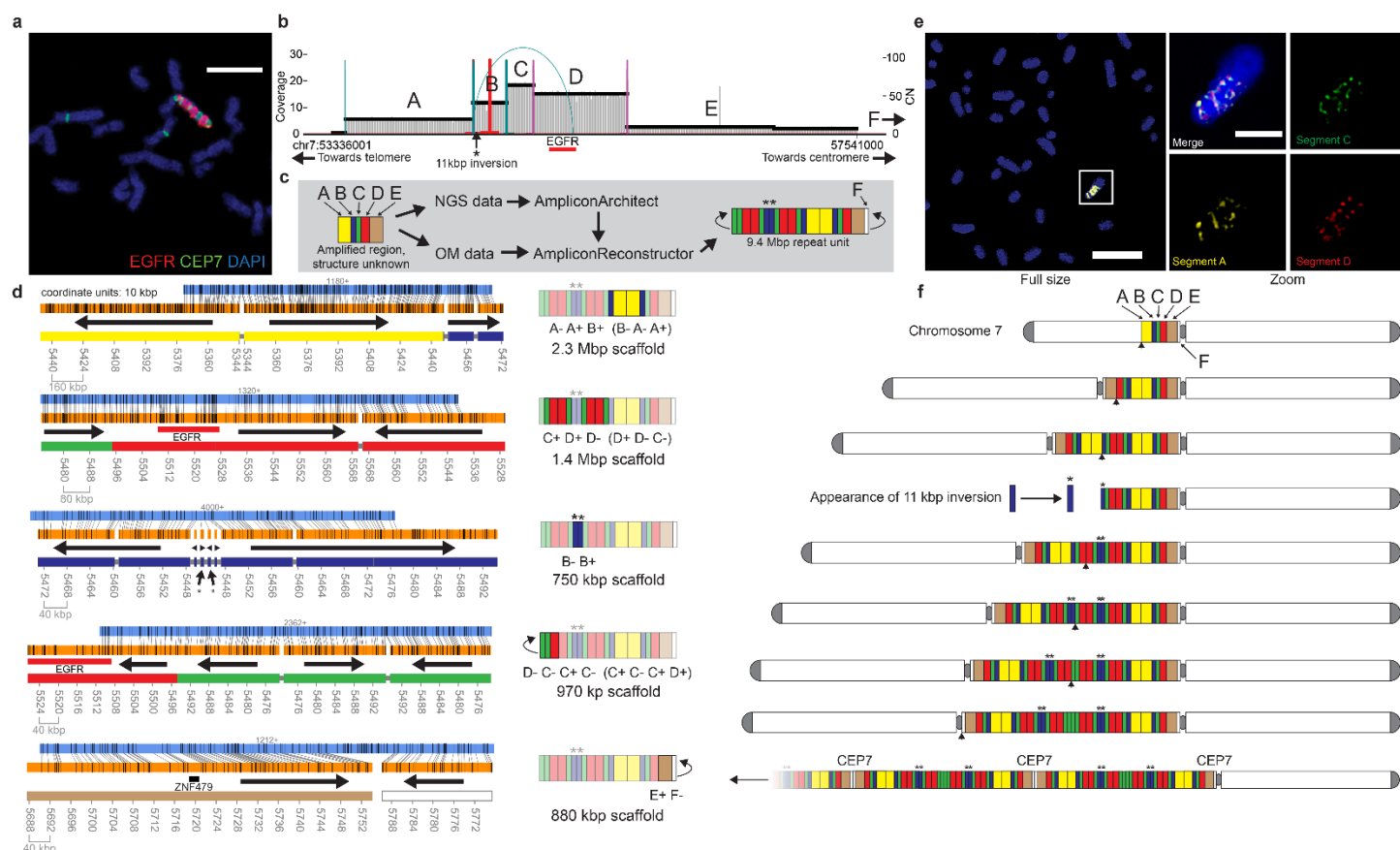
Figure 4: **[Reconstruction of Breakage-Fusion-Bridge.] a**, FISH confocal microscopy of DAPI-stained metaphase chromosomes in HCC827 showing multiple distinct bands of EGFR and CEP7 (chr7 centromeric repeat probe). Scale bar indicates 6µm. **b,** AA-generated breakpoint graph for amplified EGFR region in HCC827. Asterisk ('*') symbol indicates presence of 11 kbp inversion at 5' end of segment B. **c**, Workflow for analysis of amplified EGFR region in HCC827 to reveal BFB repeat unit structure. Amplified intervals detected by AA are labeled A-E and are colored yellow, blue, green, red and brown, respectively. "F" indicates a region identified by AR but not AA. **d**, Visualization of the AR-generated scaffolds (left column) and cartoon illustration of reconstructed region(s) of the BFB (right column), including segment sequence. Black arrows in the scaffold column indicate segment directionality. **e**, Multi-FISH for BFB segments using super-resolution confocal microscopy on DAPI-stained metaphase chromosomes in HCC827. FISH probes used for segments "A", "C", and "D" were RP11-64M3, RP11-117I14, and EGFR, respectively. Scale for full size image indicates 11 µm. Scale bar for zoomed images indicates 3 µm. Brightness was decreased using ImageJ between full size and zoomed images. **f**, Theoretical model of formation for HCC827 EGFR BFB. Each row indicates a prefix inversion and duplication characteristic of BFB, alongside other SVs. Black arrowheads beneath the intermediate step in each row indicates the breakpoint of the BFB chromosome. The bottom row shows multiple duplications of the BFB unit along with a pericentromeric region.