

Robust identification of extrachromosomal DNA and genetic variants using multiple genetic abnormality sequencing (MGA-Seq)

Da Lin^{2*#}, Yanyan Zou^{1,5#}, Jinyue Wang^{1,4,6}, Qin Xiao^{1,4,6}, Fei Lin⁷, Ningyuan Zhang⁷, Zhaowei Teng⁸, Shiyi Li^{9,10}, Yongchang Wei^{10,11}, Fuling Zhou¹², Rong Yin¹², Siheng Zhang^{1,3}, Chengchao Wu^{1,3}, Jing Zhang¹³, Sheng Hu¹³, Shuang Dong¹³, Xiaoyu Li¹³, Shengwei Ye¹⁴, Haixiang Sun^{7*}, Gang Cao^{1,3,4,5*}

Affiliations

¹State Key Laboratory of Agricultural Microbiology, Huazhong Agricultural University, Wuhan, China.

²Precision Research Center for Refractory Diseases, Institute for Clinical Research, Shanghai General Hospital, Shanghai Jiao Tong University School of Medicine, Shanghai, China.

³College of Veterinary Medicine, Huazhong Agricultural University, Wuhan, China.

⁴College of Bio-Medicine and Health, Huazhong Agricultural University, Wuhan, China.

⁵College of Informatics, Huazhong Agricultural University, Wuhan, China.

⁶College of Life Science and Technology, Huazhong Agricultural University, Wuhan, China.

⁷Reproductive Medical Center, Nanjing Drum Tower Hospital, The Affiliated Hospital of Nanjing University Medical School, Nanjing, China.

⁸The First People's Hospital of Yunnan Province, Affiliated Hospital of Kunming University of Science and Technology, Kunming, China.

⁹Baylor College of Medicine, Houston, Texas, USA.

¹⁰Department of Radiation & Medical Oncology, Zhongnan Hospital of Wuhan University, Wuhan, China.

¹¹Hubei Key Laboratory of Tumor Biological Behaviors, Zhongnan Hospital of Wuhan University, Wuhan, China.

¹²Department of Hematology, Zhongnan Hospital of Wuhan University, Wuhan, China.

¹³Department of Medical Oncology, Hubei Cancer Hospital, Tongji Medical College, Huazhong University of Science and Technology, Wuhan, China.

¹⁴Department of Gastrointestinal Surgery, Hubei Cancer Hospital, Tongji Medical College, Huazhong University of Science and Technology, Wuhan, China.

[#]These authors contributed equally: Da Lin, Yanyan Zou.

*Correspondence:

Da Lin, Ph.D, dalin870427@gmail.com

Haixiang Sun, Ph.D, Professor, stevensunz@163.com

Gang Cao, Ph.D, Professor, gcao@mail.hzau.edu.cn

38

39 SUMMARY

40 Genomic abnormalities, including structural variation (SV), copy number variation
 41 (CNV), single-nucleotide polymorphism (SNP), homogenously staining regions (HSR)
 42 and extrachromosomal DNA (ecDNA), are strongly associated with cancer, rare
 43 diseases and infertility. A robust technology to simultaneously detect these genomic
 44 abnormalities is highly desired for clinical diagnosis and basic research. In this study,
 45 we developed a simple and cost-effective method – multiple genetic abnormality
 46 sequencing (MGA-Seq) – to simultaneously detect SNPs, CNVs, SVs, ecDNA and
 47 HSRs in a single tube. This method has been successfully applied in both cancer cell
 48 lines and clinical tumour samples and revealed that focal amplification in tumour
 49 tissue is substantially heterogeneous. Notably, we delineated the architecture of focal
 50 amplification and the ecDNA network by MGA-Seq, which facilitated the exploration
 51 of the regulation of gene expression in ecDNA. This method could be extensively
 52 applied for diagnosis and may greatly facilitate the investigation of the genomic
 53 mechanism for genetic diseases.

54

55

56

57

58

59

60

INTRODUCTION

Genomic abnormalities, including structural variation (SV), copy number variation (CNV), focal amplification (FA) (1), and single-nucleotide polymorphisms (SNPs), are strongly associated with the development and progression of cancer (2, 3), rare diseases (RDs) (4) and infertility (5, 6). Accumulating data have demonstrated that numerous cancer cells contain extrachromosomal DNA (ecDNA), a form of FA (7). The copy number of oncogenes can be highly elevated by ecDNA-based amplification. Moreover, the chromatin architecture of ecDNA is usually highly accessible (8), which dramatically increases the expression level of oncogenes. ecDNAs can be spatially close to each other and cluster together to form ecDNA hubs (8-10), which perform enhancer-like functions and increase the expression of proto-oncogenes through intermolecular interactions (8, 9, 11). Intriguingly, in response to antitumor drug treatment, ecDNA can reintegrate back into the chromosome in another form of FA, homogenously staining regions (HSRs), via a myriad of mechanisms (12). Increasing evidence suggests that ecDNA is associated with cancer progression and can be used as a diagnostic marker (13, 14). However, there is no method thus far to simultaneously detect diverse types of genomic abnormalities, which greatly hampers the precise diagnosis and understanding of the molecular mechanism of cancer and genetic disease.

Second-generation sequencing-based whole exome sequencing (WES) and whole genome sequencing (WGS) can efficiently detect single-nucleotide variants (SNVs) and small indels (< 50 bp). However, due to the limitations of short read length, it is extremely challenging to identify larger inversions, translocations, insertions (>1 Mb) (15), and ecDNA. To improve the detection capability of complex genomic structural variation, several new technologies have been developed (15, 16). These technologies can be generally divided into two categories: one is based on single molecule long fragment sequencing or detection, such as Pacific Biosciences (PacBio) SMRT sequencing (17, 18), Oxford Nanopore Technologies (ONT) sequencing (16, 19, 20),

and Bionano (21); the other is based on long DNA sequence reconstruction using short read sequencing, such as strand-seq (22, 23), 10x Genomics linked-reads (24-26), and Hi-C (27, 28). Due to the high cost, tedious experimental steps, and large amount of initial sample, these technologies are mostly applied in scientific research, such as genome assembly (29-31), full-length transcriptome sequencing (32), and gene transcription regulation (33), but not for clinic diagnosis.

Based on WGS datasets, researchers developed the FA prediction software AmpliconArchitect (34) and delineated the focal amplifications and general structure of ecDNA in different types of tumours (35). Due to the natural disadvantage of the short read length of next-generation sequencing datasets, the accuracy of AmpliconArchitect prediction results is limited, and there is no spatial structural information of ecDNA hubs. Recently, a multiomics strategy based on second-generation sequencing, third-generation sequencing, and Hi-C has been developed to decode the spatial architecture of ecDNA hubs in detail (9, 36). This integrated analysis strategy can effectively decode the circular structure and spatial mobility of ecDNA. However, this strategy requires expensive multiple sequencing library construction and sequencing from the same sample, which limits its clinical application for precise diagnosis. Thus, a simple method for the simultaneous detection of different types of genomic abnormalities is crucial and highly desired for precise diagnosis and understanding the molecular mechanism of cancer and genetic disease.

In this study, we sought to develop an efficient and cost-effective method, multiple genetic abnormality sequencing (MGA-Seq), to simultaneously detect SNPs, CNVs, SVs, and the spatial architecture of FA and distinguish ecDNA from HSR. Using MGA-Seq, we successfully identified SNPs, CNVs, specific chromosomal translocation types and breakpoints with single-base resolution in cancer cells and blood samples from infertile patients. As MGA-Seq can locate the approximate location of genomic structural variation, it can facilitate breakpoint searching. We

demonstrated that MGA-Seq can indeed distinguish HSR and ecDNA and construct spatial structure and interaction networks of focal amplification regions that could be extensively applied for precise diagnosis and the investigation of the molecular mechanism of cancer and genetic disease.

Results

Overview of MGA-Seq

To maintain the spatial architecture of the genome, the nuclei are first fixed by formaldehyde in multiple genetic abnormality sequencing (MGA-Seq). The genome is digested *in situ* by restriction endonuclease followed by 5' DNA overhang fill-in by DNA polymerase I. Next, the spatially adjacent chromatin fragments are proximity ligated using T4 DNA ligase and then fragmented into a high-throughput sequencing library (**Fig. 1A**). This library contains two kinds of sequencing reads. The reads without proximity ligation junctions were used to detect SNPs, CNVs, small inserts and deletions (< 50 bp), focal amplification (FA), and genomic breakpoints (**Fig. 1A** and **fig. S1**). As the reads with proximity ligation junctions contain spatially adjacent chromosome fragment contact information of the genome, they can be used to decode chromosome structure. Thus, the integrated analysis of all the sequencing reads can identify large chromosome structural variation, such as balanced and unbalanced translocations, extrachromosomal DNA (ecDNA), and intrachromosomal homogeneously staining regions (HSRs). Notably, all MGA-Seq steps are carried out in the same tube and do not require buffer replacement, which takes only 9 hours and costs just 56 dollars (**Fig. 1B**).

Identification of SNPs and indels by MGA-Seq

To evaluate the SNP and indel detection capability, we performed MGA-Seq on the colorectal cancer cell line SW480 as described in **Fig. 1A**. After sequencing, we obtained 194,167,430 read pairs, of which 2,982,113 (1.5%) read pairs contained

149 “AAGCTAGCTT” ligation junction sequences. To avoid false-positives caused by
150 ligation junctions, we filtered out this part of the reads for SNP and indel detection
151 (see Methods) and analysed the remaining reads by the Genome Analysis Toolkit
152 (GATK). To evaluate the SNP and indel variation calling efficacy, we used the SW480
153 cell line to generate standard WGS datasets and downloaded the SW480 *in situ* Hi-C
154 datasets (37) for comparison with the same parameters (see Methods). As shown in
155 **Fig. 2A**, MGA-Seq identified 2,722,682 variants, including 2,446,823 SNPs, 130,087
156 insertions, and 145,772 deletions. A total of 82.8% of these variants were consistent
157 with WGS (**Fig. 2B**). Hi-C found only 1,166,315 variants, which is much lower than
158 that identified by MGA-Seq and WGS (**Fig. 2A**). Furthermore, the sequencing
159 coverage and depth of MGA-Seq were also much higher than those of Hi-C (**fig. S2**,
160 **A-C**).

161

162 **Detection of chromosome copy number variation by MGA-Seq**

163 To test the CNV detection capability of MGA-Seq, we plotted the log₂ ratio of
164 average read depths in 50 Kb bins across the genome, as shown in **Fig. 2C**. Our data
165 showed that the genome coverage and uniformity of MGA-Seq are highly consistent
166 with the gold standard WGS datasets and much higher than those of the Hi-C datasets.
167 After zooming in on chromosome 3, we observed that Hi-C roughly divided
168 chromosome 3 into two CNV intervals, whereas MGA-Seq accurately identified all
169 the small copy number variation across the whole chromosome (**Fig. 2D**). Next, we
170 systematically analysed the size and number of CNVs identified by these three
171 methods (**Fig. 2, E-G**) and found that it was extremely difficult to detect CNVs less
172 than 10 Mb by Hi-C (**Fig. 2, E and F**). In this scenario, the CNV detection capability
173 of MGA-Seq is much better than that of Hi-C, especially for micro-CNVs (<1 Mb),
174 which is highly consistent with WGS (**Fig. 2, E and G**).

175

176 **Identification of chromosomal translocations and breakpoints by MGA-Seq with** 177 **single base-pair resolution**

178 By using SW480 MGA-Seq sequencing datasets, we obtained the genome-wide

179 chromosome contact matrix. As shown in **Fig. 3, A and B**, we identified 8
180 translocations and 1 inversion. Although MGA-Seq only used 190 million raw reads,
181 the structural variants detected by MGA-Seq were completely consistent with *in situ*
182 Hi-C with 300 million raw reads (**Fig. 3A**). To further identify the chromosomal
183 translocation types and breakpoints of these translocations, we combined chromosome
184 contact matrix, CNV, and split read information from MGA-Seq datasets and
185 performed integrated analysis. Taking T(2;12)(q35;q12) as an example, from the CNV
186 data, we observed that the copy number of chromosome 12 was increased, whereas
187 the copy number of chromosome 2 was decreased downstream of the chromosome
188 breakpoint (**Fig. 3C**), suggesting that unbalanced translocation occurred between
189 chromosomes 2 and 12.

190 Based on the split reads information in MGA-Seq, we further identified that the
191 translocation breakpoint is located at chr2: 220,857,416 and chr12: 43,120,970 (**Fig.**
192 **3C**). In contrast, due to the low genome coverage and depth of Hi-C, it is not feasible
193 to precisely determine the type and breakpoint of translocation (**Fig. 3, B and D**). In
194 this scenario, MGA-Seq identified that all 8 chromosomal translocations in the
195 SW480 cell line were unbalanced translocations. Notably, we were able to pinpoint
196 the breakpoints of 6 out of 8 translocations sites at single-base resolution (75.0 %).
197 We also used WGS data with the same sequencing depth as MGA-Seq to identify the
198 translocations. As there is no chromosome interaction information in this dataset,
199 none of the chromosomal translocations were found (**Fig. 3B**). Moreover, we verified
200 the T(2;12)(q35;q12) translocation by two-colour DNA fluorescence *in situ*
201 hybridization (FISH). As shown in **Fig. 3E**, chromosomes 2 and 12 were indeed fused
202 together in SW480 cells, supporting the integrity of MGA-Seq.

203 Furthermore, to test the chromosomal translocation detection capability of MGA-Seq
204 in clinical samples, we collected peripheral blood from two infertile patients with
205 known translocation sites and constructed an MGA-Seq library. By combining the
206 chromosome interaction matrix and CNV data, we detected a T(10;22)(p12;q13)
207 translocation in sample 1 (**fig. S3A**) and a T(9;11)(q21;p14) translocation in sample 2
208 (**fig. S3B**), which are consistent with the known translocation sites identified by

209 karyotyping. In addition, based on the split reads, we pinpointed the precise location
210 of the breakpoints with single base-pair resolution (**fig. S3, A and B**). Next, we
211 analysed the CNV information of these two samples based on the MGA-Seq data to
212 determine the translocation type. Our data showed that there are no chromosome copy
213 number changes around the translocation breakpoint, meaning that both infertile
214 patients carry balanced translocations. Together, these data demonstrated that MGA-
215 Seq can detect specific chromosomal translocation types and the corresponding
216 breakpoint with high efficacy and low cost.

217

218 **Detection of ecDNA and HSR by MGA-Seq**

219 There are two types of focal amplifications, extrachromosomal DNA (ecDNA) and
220 intrachromosomal HSRs. Due to the high mobility and dramatic amplification amount
221 of ecDNA, we speculated that ecDNA can randomly interact with each chromosome
222 with a significantly higher interaction frequency than the normal interchromosome
223 interaction, while HSRs only interact strongly within the specific chromosomes (**Fig.**
224 **4A**). To prove this hypothesis, we selected the ecDNA-positive cell line COLO320-
225 DM (7) and the HSR-positive cell line SW480 (38) for MGA-Seq analysis. First,
226 MYC amplifications in the form of ecDNA in COLO320-DM cells and in the form of
227 HSR in SW480 cells were confirmed by DNA FISH (**Fig. 4, B and C**). In comparison
228 to HSR-positive SW480 cells, ecDNA-positive COLO320-DM cells showed MYC
229 amplification throughout the nucleus (**Fig. 4, D and E**). Furthermore, CNV analysis
230 based on the MGA-Seq dataset accurately located the MYC amplification regions in
231 these two cell lines (**Fig. 4, F-I**).

232 Next, we constructed the chromatin interaction matrix using MGA-Seq data. Since the
233 amplified ecDNAs were randomly distributed in the nucleus (**Fig. 4D**), the ecDNA
234 fragments were unbiasedly ligated to all the chromatin fragments upon proximity
235 ligation and thus presented a strip-like structure in the whole chromatin contact matrix
236 (**Fig. S4A**). In contrast, as HSR is amplified on specific chromosomal regions (**Fig. 4,**
237 **C and E**), it only shows strong interchromosomal interactions on certain
238 chromosomes (**fig. S4B**), which is consistent with our hypothesis (**Fig. 4A**). In

239 addition, we observed the same interchromosomal interaction pattern in ecDNA-
240 positive cell lines TR14 and SUN16 (7, 9) (**fig. S4, C and D**). From the
241 interchromosomal interaction matrix of SW480, we found that the MYC focal
242 amplification region has a strong interaction with 19q13.3, indicating that MYC is
243 likely to be amplified on chr19 (**fig. S4, E and F**). This finding is consistent with a
244 previous report (38).
245 Since the judgement dependent on the naked eye is subjective and differs among
246 individuals, we performed genome-wide interaction fluctuation analysis (GWIFA) on
247 the focal amplification regions (**Fig. 4, J-O**) for more objective identification of HSR
248 and ecDNA (see Methods). First, we divided the genome into fixed-size bins and
249 calculated the interaction intensity between the amplified region and each bin (**Fig. 4,**
250 **J and K**). The cumulative interaction intensity curve was then plotted as shown in **Fig.**
251 **4, L and M**. Next, second-order backwards difference (SOBD) analysis was applied
252 to evaluate the fluctuation of the cumulative interaction intensity curve (**Fig. 4, N and**
253 **O**). As HSR is amplified on the specific chromosome, the value of SOBD fluctuates
254 dramatically at specific genomic locations (**Fig. 4O**). However, ecDNA has strong
255 interactions with distinct strengths across the whole genome. Thus, the value of
256 SOBD fluctuates greatly throughout the whole genome (**Fig. 4N, and fig. S4, G and**
257 **H**).

258

259 **Delineation of the architecture of focal amplification in K562 cells**

260 Our MGA-Seq analysis of K562 cells identified an abnormal increase in chromosome
261 copy number on specific regions on chromosomes 9, 13, and 22 (**fig. S5A**). After
262 zooming in on the abnormally amplified regions, we identified six precisely amplified
263 subregions, one on chromosome 9, four on chromosome 13, and one on chromosome
264 22, which were named “A” to “F”, respectively (**Fig. 5A**). Based on genome-wide
265 interaction fluctuation analysis (GWIFA), we found that these regions were amplified
266 in K562 cells in the form of HSR rather than ecDNA (**fig. S5B**). Notably, we observed
267 strong interactions between these amplified regions, suggesting that these regions are
268 spatially close together, which likely originate from the same HSR (**Fig. 5B**). Taking

the “B”, “C”, and “D” amplified regions of chromosome 13 as examples, these three regions are in high contact with each other and form a high-density topologically associating domain (TAD)-like structure (39) (**Fig. 5C**). Such abnormal genome amplification and TAD-like structures were absent in healthy human peripheral blood cells (**Fig. 5D**).

Since the MGA-Seq dataset contains whole-genome sequencing information, we extracted the split reads located at the boundaries of these six amplified regions (table S1) and assembled the structure of HSR. In the K562 cell line, *ABL1* in the “A” amplification region, *GPC5* in the “B” amplification region, *GPC6* in the “D” amplification region, and *DGCR8* and *BCR* in the “F” amplification region were spliced to form a repeating HSR (**Fig. 5E and fig. S6**). To validate the HSR structure predicted by MGA-Seq, we compared our predicted results with published K562 third-generation sequencing data (36). Our analysis showed that *ABL1* in chromosome 9, *GPC5* and *GPC6* in chromosome 13, and *DGCR8* and *BCR* in chromosome 22 indeed come from the same scaffold, which is highly consistent with our results. Finally, we applied DNA FISH to verify the spatial location of the *ABL1* amplification region on Chr9 and the *BCR* amplification region on chromosome 22. As shown in **Fig. 5F**, *ABL1* and *BCR* indeed come from the same HSR.

287

288 Identification of focal amplification in tumour tissue

Next, we applied MGA-Seq to tumour samples and detected 40 focal amplification regions in one renal cancer tissue (**fig. S7 and table S2**). The length distribution of these regions varies from 4.2 Kb to 2.53 Mb (**table S2**). These amplified regions contain a large number of immune genes, oncogenes, and enhancers, such as *CHD1L*, *BCL6*, *JAK2*, *PD-L1*, and *CDK4* (**Fig. 6A and table S2**). In addition, the RNA transcription level of these genes within the amplified region was significantly higher than that of the normal kidney tissue control (**Fig. 6A and fig. S8A**). Through MGA-Seq chromatin contact matrix and GWIFA (**fig. S8, B and C**), we identified that these FA regions are amplified in the form of HSR. Of note, these amplified regions are not independent but contact each other at the spatial level (**Fig. 6B**).

299 FA in tumour tissue is highly heterogeneous compared to single-cell-derived cell lines.
 300 Taking the PD-L1 amplification region on chromosome 9 in this tumour as an
 301 example, this region can be spliced with multiple FA regions, as indicated by the split
 302 reads in the MGA-Seq dataset, suggesting that multiple types of HSR coexist in this
 303 heterogeneous tumour tissue (**Fig. 6C**). To verify this result, we performed single-
 304 molecule nanopore sequencing on the same tumour sample, which revealed highly
 305 consistent inter- and intrachromosomal structural variation as with the MGA-Seq
 306 dataset (**fig. S8D**). The inter- and intrachromosomal interaction analysis of chr1, chr9,
 307 and chr12 amplification regions based on the MGA-Seq chromatin contact matrix
 308 (**Fig. 6D**) identified highly complicated and heterogeneous spatial architectures of
 309 these FA regions. For example, chr1 and chr12 show an “uneven amplification”
 310 pattern, meaning that in a certain chromosome interval, only some regions were
 311 amplified, such as the regions containing proto-oncogenes, immune genes, and some
 312 regulatory elements (**Fig. 6, A and D**). These genes and regulatory elements are
 313 spliced together and eventually form a variety of HSRs (**Fig. 6E**). Based on the
 314 chromatin contact information and split reads, we constructed an interaction network
 315 of these amplified oncogenes, in which the TNFSF18 region was connected with 11
 316 amplified regions as supported by the split reads (**Fig. 6F**).

317 Discussion

318 The occurrence of tumours, infertility, and rare diseases are closely related to focal
 319 amplification (34) and structural variation (4, 5). These genetic diseases affect
 320 hundreds of millions of people around the world and have become a major human
 321 health concern. An effective multiple genetic abnormality detection method is highly
 322 desired for clinical diagnosis. In this study, we developed multiple genetic
 323 abnormality sequencing (MGA-Seq) to simultaneously detect SNPs, CNVs, SVs, and
 324 the spatial architecture of FAs and distinguish ecDNA from HSR. Taking advantage of
 325 the versatility of reaction buffers, all the MGA-Seq library construction steps are
 326 carried out in a single tube, which can minimize sample loss due to buffer exchanges
 327 and simplify the operation. Notably, MGA-Seq takes only 9 hours and costs just 56
 328 dollars to complete all the sequencing library construction steps. It demonstrated
 329 robust detection capability for both small (SNPs and INDELS) and large (CNV, SV,
 330 HSR, and ecDNA) genomic abnormalities and has great potential for clinical and
 331 scientific research applications.

332
 333 ecDNA is prevalent in at least 30 different cancer types, is closely associated with
 334 cancer progression (11, 40), and might be used as a potential prognostic marker.
 335 However, there is still a lack of an unbiased and efficient detection method in clinical
 336 practice. While AmpliconArchitect can be used for ecDNA prediction (34), the
 337 identification of ecDNA based on the WGS dataset generally has a high false-positive
 338 rate. For instance, in the cell line K562 in this study, due to the head-to-tail tandem
 339 duplication HSR structure (**Fig. 5E**), a large number of split reads also presented a
 340 circle junction-like structure. Circle-Seq can effectively analyse the structure of
 341 circular DNAs (41, 42). However, the DNA extraction process of this method can
 342 easily destroy the circular structure of ecDNA. Moreover, Circle-Seq is based on
 343 rolling-circle DNA amplification; it preferentially amplifies smaller circular DNAs,
 344 resulting in biased amplification results. Here, we demonstrated that MGA-Seq can
 345 unbiasedly detect the presence of ecDNA in both cell lines and clinical samples.
 346 Importantly, we proposed an ecDNA detection algorithm, GWIFA, and successfully

differentiated ecDNA and HSR. Of note, MGA-Seq can reveal trans interactions between ecDNA and the genome and decode the ecDNA network, which could facilitate the exploration of the potential regulation of the expression inside of the ecDNA.

Chromosomal translocations can be divided into unbalanced and balanced translocations. Unbalanced translocation usually occurs with an altered chromosomal copy number at the breakpoint (gain or loss of genetic material), resulting in abnormal gene expression. A large number of unbalanced translocations have been found in cancer cells (43, 44), especially in blood tumour genomes (45, 46). Balanced translocations do not have any genetic material changes. These translocation carriers usually have normal phenotypes and intelligence but can produce various unbalanced rearranged gametes during germ cell meiosis, resulting in infertility, abortion, stillbirth, and multiple malformations (6, 47, 48). Thus far, it is still challenging to precisely identify the specific translocation types by a simple and cost-effective method. As MGA-Seq contains CNV and chromatin contact information, it can guide translocation breakpoint searching and facilitate to identifying translocation types and breakpoints. Here, we revealed the translocation types and breakpoints of infertile couples by MGA-Seq. With this important information, high-quality blastocysts can be quickly screened by PCR before blastocyst transfer during *in vitro* fertilization, which greatly reduces the cost and time of traditional whole genome sequencing for each blastocyst.

Together, we developed a simple, cost-effective and robust MGA-Seq to simultaneously detect SNPs, CNVs, SVs, and the spatial architecture of FA and distinguish ecDNA from HSR in a single tube experiment. We successfully identified small SNPs/INDELs and large genomic structural variations in clinical and cell line samples, decoded the focal amplification spatial structure in SW480, COLO320-DM, and K562 cell lines, and constructed interaction networks of the amplified proto-oncogenes in a clinical kidney cancer tissue sample. Our data revealed that focal

377 amplification is highly diverse in tumour tissue compared to single-cell-derived
378 cancer cell lines. In the future, it would be important to develop single-cell MGA-Seq
379 for diverse ecDNA detection in single cells or highly heterogeneous cancer cells. With
380 its multifunctional and cost-effective advantages, we expect MGA-Seq to be
381 extensively applied for the diagnosis of cancer, infertility, and rare diseases and may
382 greatly facilitate the investigation of the genomic mechanism for genetic diseases.
383
384

Figure legend

Figure 1. Experimental procedure, time, and cost of multiple genetic abnormalities sequencing (MGA-Seq). (A) Flowchart of MGA-Seq. Nuclei were cross-linked with 0.5% formaldehyde and then digested with HindIII. 5' DNA overhangs of digested chromatin fragments were filled in by DNA polymerase and then proximity ligated by T4 DNA ligase. The proximity ligation products were fragmented and then subjected to high-throughput sequencing library construction. After sequencing, all the reads were used to generate chromatin contact matrix for genome structural variation calling. In the sequencing library, the reads without ligation junction "AAGCTAGCTT" were used for the detection of CNV, SNP, small indels (< 50bp), region of focal amplification, and genome breakpoints. By combining all information, the types and breakpoints of structural variation can be decoded. Notably, MGA-Seq can distinguish ecDNA and HSR, predict the structure of simple focal amplification regions, and construct the interaction network of focal amplified genes. (B) The main steps, time, and cost of MGA-Seq.

Figure 2. Detection of SNPs, indels, and CNVs by MGA-Seq. (A) Comparison of the numbers of SNPs and indels (< 50bp, include insertions and deletions) identified by WGS, MGA-Seq, and Hi-C. (B) Overlap of the SNPs and indels between MGA-Seq, WGS, and Hi-C. (C) Comparison of log₂ copy ratios calculated using reads coverage between Hi-C, MGA-Seq, and WGS. (D) Comparison of the CNVs on chromosome 3 identified by Hi-C, MGA-Seq, and WGS. (E) Statistics of the number and size distribution of CNVs identified by Hi-C, MGA-Seq, and WGS. (F) Consistency of the CNV segments (categorized by size) detected by Hi-C and WGS. Overall, Hi-C cannot detect CNV with length less than 20 Mb. (G) Consistency of the CNV segments detected by MGA-Seq and WGS. The number and size distribution of CNV segments detected by MGA-Seq and WGS are highly consistent, especially for micro-CNVs (< 1Mb).

Figure 3. Identification of translocation types and breakpoints in SW480 at single base-pair resolution by MGA-Seq. (A) Identification of translocation in the SW480 cell line by genomic contact matrix constructed with MGA-Seq datasets. The detected structural variations are indicated by arrows. (B) Translocation types and breakpoint information defined by MGA-Seq. (C) Application of integrated chromatin contact matrix, CNVs, and split reads analysis to identify translocation types and breakpoints between chr 2 and chr 12 at single base-pair resolution using MGA-Seq datasets. (D) Identification of translocation types and breakpoints between chr 2 and chr 12 using Hi-C datasets. (E) Validation of the T(2;12)(q35;q12) translocation in SW480 cells by DNA FISH. FISH probes for 12q12 and 2q35 were directly labeled with Alexa Fluor 555 (red) and Alexa Fluor 488 (green), respectively. K562 cells without the T(2;12)(q35;q12) translocation were used as a control.

Figure 4. Identification of ecDNA and HSR by MGA-Seq. (A) Putative diagram of inter-chromosomal interaction pattern differences between ecDNA and HSR positive cell line. (B-E) Validation of *MYC* amplification in COLO320-DM and SW480 cell lines by DNA FISH. The red signal represents *MYC* and the green signal represents the centromere of chr 8. (F and G) Copy number variation analysis of chr 8 in COLO320-DM and SW480 cell lines. Gains and losses of copy number are shown in red and blue, respectively. (H and I) Location of the *MYC* amplification region in COLO320-DM and SW480 cell lines. (J and K) Interaction intensity between the focal amplified region and whole genome. (L and M) Cumulative interaction intensity curve of COLO320-DM and SW480 cell lines. The x-axis represents the genome position, 100 kb bin size. The Y axis represents the accumulation of interaction intensity. (N and O) Plotted the second order backward difference (SOBD) value across the genome of COLO320-DM and SW480 cell lines in 100-kb bin size.

Figure 5. Deciphering the spatial structure of homogeneously staining region (HSR) in K562 cell line. (A) Location of the amplification region on chr 9, 13, and 22. (B) Circos plots of the chromatin interactions mediated by amplification regions

444 across all 23 chromosomes in K562 cell lines. The interactions between chromosomes
445 9, 13, and 22 are marked with red lines. **(C and D)** Comparison of chromatin contact
446 matrix of amplification region (Chr13:90423781-92475244, Chr13:92943122-
447 93351872, and Chr13:93848028-94027981) between K562 cell line and healthy
448 human peripheral blood cells (control). **(E)** Assembling the amplified regions from
449 "A" to "F" with split reads. The breakpoint of the amplification regions is marked in
450 the figure. **(F)** Metaphase analysis and DNA FISH to validate the location of the
451 *ABL1* amplification region and the *BCR* amplification region in K562 cell line. FISH
452 probes for the *ABL1* amplification region and the *BCR* amplification region were
453 directly labeled with Alexa Fluor 555 (red) and Alexa Fluor 488 (green), respectively.

454

455 **Figure 6. Heterogeneity of focal amplification in renal cancer tissue.** **(A)**
456 Sequencing reads coverage and RNA expression level in typical focal amplification
457 regions of a renal cancer tissue sample. **(B)** Circos plots of the chromatin interactions
458 mediated by focal amplification regions across all 23 chromosomes in renal cancer
459 tissue. **(C)** Circos plots of the split reads mediated by focal amplification regions
460 across all 23 chromosomes. The split reads aligned to the PD-L1 amplified region are
461 marked with red lines. **(D)** Chromatin contact matrix between the amplified regions of
462 chr1, chr9, and chr12, and sequencing reads coverage within these amplified regions.
463 **(E)** Diverse structures of HSR in the renal cancer tissue sample. **(F)** Interaction
464 network of amplified oncogenes in the renal cancer tissue sample. Different amplified
465 oncogenes are assembled by split reads. The thickness of the line indicates the
466 chromatin contact strength.

467

468 **Figure S1. Flow-chart of MGA-Seq data analysis.** After sequencing, all sequencing
469 reads were used to generate chromatin contact matrix by juicer pipeline. The reads
470 without proximity ligation junction were used to detect small indels and SNPs (<
471 50bp), CNVs, split reads, and genomic amplification regions. With the integrated
472 analysis of chromatin contact matrix, these datasets can be used to decode the type
473 and breakpoints of translocations, distinguish ecDNA from HSR, predict the focal

474 amplification structure, and construct FA region interaction network.

475

476 **Figure S2. Comparison of the sequencing depth and coverage of MGA-Seq, WGS,**

477 **and Hi-C. (A)** Scatter plot of sequencing depth and coverage for each chromosome.

478 Blue points represent MGA-Seq, yellow points represent MGA-Seq, green points

479 represent Hi-C. X-axis represents coverage, and Y-axis represents sequencing depth.

480 **(B)** Histogram of coverage for each chromosome. Blue represents MGA-Seq, yellow

481 represents MGA-Seq, and green represents Hi-C. **(C)** Histogram of sequencing depth

482 for each chromosome. Blue represents MGA-Seq, yellow represents MGA-Seq, and

483 green represents Hi-C.

484

485 **Figure S3. Identification of translocation types and breakpoints by MGA-Seq. (A)**

486 Identification of balance translocation T(10;22)(p12;q13) and genome breakpoint in

487 patient 1. **(B)** Identification of balance translocation T(9;11)(q21;p14) and genome

488 breakpoint in patient 2.

489

490 **Figure S4. Chromatin contact matrix and genome-wide interaction fluctuation**

491 **analysis (GWIFA) of ecDNA-positive cell lines. (A-D)** Genome-wide chromatin

492 contact matrix of COLO320-DM, SW480, TR14, and SUN16 cell lines. The

493 amplified regions are marked with arrows. **(E)** The chromatin interaction matrix of

494 SW480 cell line between chr 8 and chr 19. The *MYC* amplified region is marked with

495 a dashed line in the figure. **(F)** *MYC* is amplified in the form of HSR on chr 19. **(G**

496 **and H)** The second order backward difference (SOBD) value across the genome of

497 TR14 and SUN16 cell lines in 100-kb bin size.

498

499 **Figure S5. Copy number variation (CNV) analysis of K562 cell line. (A)** CNV

500 analysis of chromosomes 9, 13, and 22 in K562 cell line. Gains and losses of copy

501 number are shown in red and blue, respectively. Representative genes located in

502 amplification region are marked with arrows. **(B)** The second order backward

503 difference (SOBD) value across the genome of K562 cell line in 100-kb bin size.

504

505 **Figure S6. Sequence and breakpoints of split reads used to assemble HSR in**
506 **K562 cells.**

507

508 **Figure S7. CNV analysis of renal cancer tissue.** CNV analysis of chromosomes
509 with abnormal amplification in renal cancer tissue. Gains and losses of copy number
510 are shown in red and blue, respectively. Representative genes located in amplification
511 region are marked with arrows.

512

513 **Figure S8. Verification of inter and intra chromosomal interaction between focal**
514 **amplification regions in renal cancer tissue by nanopore. (A)** Volcano plots of
515 differential expression genes between renal cancer tissue and normal kidney tissue
516 control. **(B)** Genome-wide chromatin contact matrix of renal cancer tissue. Potential
517 HSR regions are marked with arrows in the figure. The inter-chromosomal contacts
518 between the focal amplification regions and Chr1 and Chr12 are zoomed in. **(C)** The
519 second order backward difference (SOBD) value across the genome of the renal
520 cancer tissue in 100-kb bin size. **(D)** Validation of split reads and chromatin
521 interactions across focal amplification regions with Nanopore long reads.

522

523 **Figure S9. Distribution of fluctuation score (FS) in different cell lines.** Blue bars
524 indicate HSR-positive cell lines and yellow bars indicate ecDNA-positive cell lines.
525 The Y axis represents the value of FS.

526

527 **Table S1. The split reads located at the boundaries of focally amplified regions in**
528 **K562 cells.**

529

530 **Table S2. The regions of focal amplification in the renal cancer tissue.**

531

532 **METHODS**

533 **MGA-Seq library construction**

534 **1. Preparation of cell suspension**

535 For tumor tissue, 0.5 cm³ tissue blocks were used and minced through a 40 µm
536 strainer to obtain single cell suspension. For blood samples, we directly took 1 ml of
537 anticoagulated whole blood, and centrifuged at 1500 g/min for 10 min to collect blood
538 cells.

539 **2. Nuclei preparation**

540 Cells were cross-linked with 0.5 % formaldehyde (Sigma) for 10 mins. The cross-
541 linking reaction was terminated by glycine at a final concentration of 200 mM and
542 lysed in lysis buffer (PBS contain 0.2% SDS) at room temperature for 5 min. After
543 incubation, the nuclei were pelleted by centrifugation at 2,000 g/min for 5 min. The
544 nuclei were transferred to 1.5 ml tubes and washed twice with PBS.

545 **3. *In situ* digestion**

546 For *in situ* restriction enzyme digestion, 140 µl of ddH₂O, 20 µl of 10% Triton X-100,
547 20 µl of 10× NEBuffer 2.1, and 20 µl of HindIII (NEB, 20 units/µl) were added to the
548 nuclei pellet and digested for 1.5 h at 37 °C in thermomixer (Eppendorf) with rotation
549 at 1000 r.p.m.

550 **4. End filling-in**

551 Add 5 µl of dNTP mix (10 mM each) and 5 µl of DNA polymerase I Klenow
552 fragment (NEB, M0210) to the reaction system, place the sample in thermomixer with
553 rotation at 37 °C at 1000 r.p.m for 30 mins.

554 **5. *In situ* proximity ligation**

555 Add 27.5 µl of H₂O, 3 µl of ATP (adenosine-triphosphate, 10mM), and 10 µl of T4
556 DNA ligase (Thermo, EL0011) to the reaction system, and placed the tube on the
557 rotating mixers for 2 h at room temperature with rotation at 20 r.p.m.

558 **6. Reversal of cross-linking and DNA purification**

559 Add 20 µl of proteinase K (20 µg/ml) to the proximity ligation system, and then
560 incubate at 60 °C for 2 hours. After digestion, the DNA was directly extracted using
561 PCR Purification Kits (Zymo, D4013).

562 **7. Sequencing library construction**

563 DNA sequencing libraries were prepared using the VAHTS Universal Plus DNA
564 Library Prep Kit (NDM627) according to the manufacturer's protocol.

565

566 **Metaphase analysis and DNA fluorescence in situ hybridization (FISH) assay**

567 SW480 and COLO320-DM cell lines were treated with colchicine at final
568 concentration 8 µg/ml for 24 hours. After cultivation, cells were collected by
569 centrifugation at 1000 g/min for 10 minutes. Next, 10 ml of hypotonic KCl solution
570 (0.075 M) were added to the cell pellet to resuspend the cells. After 30 min incubation
571 at 37 °C, 2 ml of fixative (3:1 methanol:glacial acetic acid) were added to the cell
572 suspension. The cell pellet was re-collected by centrifugation at 1000 g/min for 10
573 min and then resuspend in 5 ml of fixative (3:1 methanol:glacial acetic acid). After 5
574 min incubation, the cell pellet was re-collected by centrifugation at 1000 g/min for 10
575 minutes and resuspend in 1 ml of fixative. After fixation, 10 µl of the suspension were
576 dropped on the glass slide and incubated in the prewarmed 2x SSC at 60 °C for 30 min.
577 The cells were dehydrated sequentially in 70%, 85%, 100% ethanol solution. After
578 ethanol dehydration, the cells were heated on a hot plate at 82 °C for 10 min in 80 %
579 formamide (Sigma) and 2×SSC for DNA denaturation. Next, cells were incubated for
580 12 hours in hybridization solution with 2 µM DNA probes (MYC and CEP8, Spatial
581 FISH Co. Ltd.) in the presence of 50 % formamide, 8% dextran sulfate sodium salt
582 (Sigma), and 2× SSC. After hybridization, the cells were washing for three times with
583 20 % formamide and 3 times with 2×SSC. Finally, the slides were stained with DAPI
584 (Life Technologies) and observed under super-resolution microscope (Nikon, N-SIM).

585

586 **RNA-Seq library preparation**

587 RNA was extracted using the RNAiso Plus (Takara, 9109) according to the
588 manufacturer's protocol. Sequencing libraries were prepared using the VAHTS

589 Stranded mRNA-Seq Library Prep Kit (Vazyme, NR602-02) according to the
590 manufacturer's protocol.

591

592

593 **Identification of SNP, indel, split reads, and CNV using MGA-Seq datasets**

594 **1. Pre-analysis**

595 FastQC (version: 0.11.5) (49) was used to assess the quality of raw reads. FASTP
596 (version: 0.23.2) (50) was used to filter out the low-quality bases and adapter
597 sequences. The clean read pairs which contained proximity ligation junction
598 sequences "AAGCTAGCTT" were filtered out by Linux command line utility "grep".
599 The remaining reads were used for SNPs, indels, split reads, and CNV calling.

600 **2. SNP and indel calling**

601 The remaining reads were aligned to the reference genome (hg19) and generated
602 BAM file using BWA-MEM (version 0.7.17) (51). The BAM file was sorted by
603 SAMtools (version 1.15.1) (52) and deduplicated by Sambamba (53) (version 0.6.6).
604 Next, we used BaseRecalibrator (GATK, version 4.2.2) (54) to calibrate the base
605 quality scores, and HaplotypeCaller (GATK, version 4.2.2) to detect SNPs and indels.

606 **3. Split reads calling**

607 The deduplicated BAM file generated in SNP and indel calling step were used to
608 identify split reads. The split alignment reads were extracted by SAMtools (version
609 1.15.1) with command line "samtools view test_deduplicated.bam | grep SA >
610 test_split_reads.txt".

611 **4. CNV calling**

612 BIC-seq2(55) was used to derive CNV segments from reads coverage data. For more
613 details, refer to the software manual "http://www.compbio.med.harvard.edu/BIC-
614 seq/". For the segmentation step, parameters were designed as binsize=50,000 bp
615 and $\lambda = 2$ to determine the final CNV breakpoints.

616

617 **Construction of genome-wide chromatin interaction matrix using MGA-Seq** 618 **datasets**

619 FastQC (version: 0.11.5) was used to assess the quality of raw reads. FASTP (version:
620 0.23.2) was used to filter out the low-quality bases and adapter sequences. All the
621 remaining read pairs were used to generate the chromatin contacts matrix file (.hic)
622 using Juicer software(56). For more details, refer to the software manual
623 “<https://github.com/aidenlab/juicer>”.

624

625 **Identification of translocations types and breakpoints using MGA-Seq datasets**

626 The chromatin contacts matrix file (.hic) was imported into Juicebox (version: 1.9.8,
627 <https://github.com/aidenlab/Juicebox>) software for visualization. The translocations
628 and large structural variations were identified according to the inter-/intra-
629 chromosome interaction patterns (15, 57). The types and breakpoints of translocations
630 were identified according to the split reads and CNV information. For unbalanced
631 translocations, the chromosomal copy number at the breakpoint were usually altered,
632 while balanced translocations do not have any chromosomal copy number changes.

633

634 **Identification of SNP and indel using *in situ* Hi-C datasets**

635 The *in situ* Hi-C datasets of SW480 cell line (37) were downloaded from Gene
636 Expression Omnibus (GEO Accession: GSM3930294 and GSM3930295). The Hi-C
637 ligation junction sequence “GATCGATC” and bases behind ligation junction were
638 removed by FASTP (version: 0.23.2). An example command line is as follows:

639 1) fastp -i insitu_sw480_1.fq -o trim_sw480_1.fq -w 15 --adapter_sequence

640 GATCGATC

641 2) fastp -i insitu_sw480_2.fq -o trim_sw480_2.fq -w 15 --adapter_sequence

642 GATCGATC

643 Trimmed reads1 and reads2 were merged together by command line “cat
644 trim_sw480_1.fq trim_sw480_2.fq > sw480_1_2.fq”. The merged reads file was used
645 to identify SNPs and indels using the same parameters as MGA-Seq.

646

647 **Identification of CNV and translocation by *in situ* Hi-C datasets**

648 All the raw Hi-C read pairs were used to detect CNVs. FastQC (version: 0.11.5) was
649 used to assess the quality of raw reads, and FASTP (version: 0.23.2) was used to filter
650 out the low-quality bases and adapter sequences. The CNV calling was carried out by
651 BIC-seq2 (55). The observed values were the residuals from GAM Poisson regression,
652 and the expected values were set to zero. Translocation detection was performed by
653 HINT-TL as implemented in HINT (58), a computational method for detecting CNVs
654 and translocations based on Hi-C data.

655

656 **Identification of SNP, indel, and translocations using WGS datasets**

657 FastQC (version: 0.11.5) was used to evaluate the quality of raw reads. FASTP
658 (version: 0.23.2) was used to filter out the low-quality bases and adapter sequences.
659 The trimmed reads pairs were used to identify SNPs and indels. The parameters are
660 exactly the same as MGA-Seq.

661 Structural variation identification were carried out using Delly2 (59) (version: 0.8.6)
662 and Gridss (60) (version: 2.12.2) with default parameters. One WGS data from
663 healthy person was served as control. Translocations that passed the internal quality
664 control were merged with SURVIVOR (version: 1.0.7, parameters: 1000 1 1 1 0 30)
665 (61). Only translocations supported by at least one definite split alignment read were
666 retained.

667

668 **Genome-wide interaction fluctuation analysis (GWIFA)**

669 According to the inter-chromosomal interaction feature of ecDNA and HSR, we
670 designed a genome-wide interaction fluctuation analysis (GWIFA) to further
671 characterize the inter-chromosomal interaction fluctuation of the focal amplification
672 regions and defined a fluctuation score (FS) to distinguished ecDNA from HSR.

673 Firstly, we divided genome into fixed-sized bins (100 kb), and calculated the
674 cumulative interaction intensity (CII) between the focal amplified regions and the
675 whole genome (**Fig. 4, J-M**).

$$CII_x = \sum_{i=1}^x C_i$$

676 In the formula, x represents the genome position measured by the number of bin, and
677 C_i represents the number of contact counts inside the i th bin. We recommend a linear
678 fit on CII, which can eliminate the abnormal fluctuations caused by uneven
679 sequencing.

680 Next, second order backward difference (SOBD) was introduced to further
681 characterize the fluctuation of interactions across the genome (**Fig. 4, N and O**).
682 Denoting OBD of CII as OBDc.

$$OBDc_x = \frac{(\sum_{i=x-h+1}^x C_i - \sum_{j=x-2h+1}^{x-h} C_j)}{h^2}$$

683 In the formula, h is a customizable space, default is 3.

684 We then defined a fluctuation score (FS) to distinguished ecDNA from HSR.

$$FS = \frac{\sum_{i=1}^n S_i}{\sum_{j=1}^n S_j}$$

685 In the formula, S is descending sorted distribution of |OBDc| (absolute value of
686 OBDc), n is the quantity of OBDc, T is a customizable parameter ($T < 1$). After
687 multiple rounds of testing, our suggested T is 0.6 (**fig. S9**).
688 ecDNA and HSR can be distinguished as follow:

689

$$Discrimination = \begin{cases} ecDNA, & FS < T \\ HSR, & FS > T \end{cases}$$

690

691 The complete analysis pipeline is available in:

692 <https://github.com/yanyanzou0721/GWIFA>.

693

694 **Long-read sequencing (Nanopore) data analysis**

695 The nanopore sequencing reads with quality score more than 7 were mapped to the
696 reference genome hg19 using minimap2 (version: 2.17 , -ax map-ont) (62). Structural
697 variants were called using NanoSV(63) (version: 1.2.4) with default parameters. Only
698 SV supported by at least one definite split alignment read were retained for
699 subsequent statistics.

700

701 **RNA-seq data analysis**

702 FastQC (version: 0.11.5) was used to assess the quality of the raw reads. FASTP
703 (version: 0.23.2) was used to filter out the low-quality bases and adapter sequences.
704 The clean reads were aligned to the hg19 using BWA-MEM (version: 0.7.17) with
705 default parameters and sorted by Samtools (version: 1.15.1). Gene expression levels
706 were assessed using featureCounts (version: 2.0.0) (64). Differential gene expression
707 analysis was performed using DEseq2 (version: 1.20.0) (65). The complete analysis
708 pipeline is available in [https://github.com/GangCaoLab/NGS-pipelines/tree/master/RNA-](https://github.com/GangCaoLab/NGS-pipelines/tree/master/RNA-Seq)

709 [Seq.](#)

710

711 **Data availability**

712 Data have been deposited in the Gene Expression Omnibus (GEO). To review GEO
713 accession GSE205293, go to
714 <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE205293>, and enter token
715 “epgtysmkzvufmp” into the box.

716

717 **DECLARATION OF INTERESTS**

718 The authors declare no competing interests.

719

720

- 721 1. J. A. Engelman *et al.*, MET amplification leads to gefitinib resistance in lung cancer by
722 activating ERBB3 signaling. *science* **316**, 1039-1043 (2007).
- 723 2. J. R. Dixon *et al.*, Integrative detection and analysis of structural variation in cancer
724 genomes. *Nature genetics* **50**, 1388-1398 (2018).
- 725 3. L. Yang *et al.*, Diverse mechanisms of somatic structural variations in human cancer
726 genomes. *Cell* **153**, 919-929 (2013).
- 727 4. J. M. Holt *et al.*, Identification of pathogenic structural variants in rare disease patients
728 through genome sequencing. *BioRxiv*, 627661 (2019).
- 729 5. M. Zorrilla, A. N. Yatsenko, The genetics of infertility: current status of the field.
730 *Current genetic medicine reports* **1**, 247-260 (2013).
- 731 6. G. L. Harton, H. G. Tempest, Chromosomal disorders and male infertility. *Asian*
732 *journal of andrology* **14**, 32 (2012).
- 733 7. S. Wu *et al.*, Circular ecDNA promotes accessible chromatin and high oncogene
734 expression. *Nature* **575**, 699-703 (2019).
- 735 8. Y. Zhu *et al.*, Oncogenic extrachromosomal DNA functions as mobile enhancers to
736 globally amplify chromosomal transcription. *Cancer cell* **39**, 694-707. e697 (2021).
- 737 9. K. L. Hung *et al.*, ecDNA hubs drive cooperative intermolecular oncogene expression.
738 *Nature* **600**, 731-736 (2021).
- 739 10. E. Yi *et al.*, Live-cell imaging shows uneven segregation of extrachromosomal DNA
740 elements and transcriptionally active extrachromosomal DNA hubs in cancer. *Cancer*
741 *discovery*, (2021).
- 742 11. E. van Leen, L. Brückner, A. G. Henssen, The genomic and spatial mobility of

743 extrachromosomal DNA and its implications for cancer therapy. *Nature Genetics*, 1-8
744 (2022).

745 12. K. Song *et al.*, Plasticity of extrachromosomal and intrachromosomal BRAF
746 amplifications in overcoming targeted therapy dosage challenges. *Cancer discovery*,
747 (2021).

748 13. T. Wang, H. Zhang, Y. Zhou, J. Shi, Extrachromosomal circular DNA: a new potential
749 role in cancer progression. *Journal of translational medicine* **19**, 1-16 (2021).

750 14. R. G. Verhaak, V. Bafna, P. S. Mischel, Extrachromosomal oncogene amplification in
751 tumour pathogenesis and evolution. *Nature Reviews Cancer* **19**, 283-288 (2019).

752 15. S. S. Ho, A. E. Urban, R. E. Mills, Structural variation in the sequencing era. *Nature*
753 *Reviews Genetics* **21**, 171-189 (2020).

754 16. G. A. Logsdon, M. R. Vollger, E. E. Eichler, Long-read human genome sequencing
755 and its applications. *Nature Reviews Genetics* **21**, 597-614 (2020).

756 17. M. O. Carneiro *et al.*, Pacific biosciences sequencing technology for genotyping and
757 variation discovery in human data. *BMC genomics* **13**, 1-7 (2012).

758 18. M. A. Quail *et al.*, A tale of three next generation sequencing platforms: comparison
759 of Ion Torrent, Pacific Biosciences and Illumina MiSeq sequencers. *BMC genomics*
760 **13**, 1-13 (2012).

761 19. T. Laver *et al.*, Assessing the performance of the oxford nanopore technologies
762 minion. *Biomolecular detection and quantification* **3**, 1-8 (2015).

763 20. J. L. Weirather *et al.*, Comprehensive comparison of Pacific Biosciences and Oxford
764 Nanopore Technologies and their applications to transcriptome analysis.

765 *F1000Research* **6**, (2017).

766 21. E. T. Lam *et al.*, Genome mapping on nanochannel arrays for structural variation
767 analysis and sequence assembly. *Nature biotechnology* **30**, 771-776 (2012).

768 22. E. Falconer, P. M. Lansdorp, in *Seminars in cell & developmental biology*. (Elsevier,
769 2013), vol. 24, pp. 643-652.

770 23. A. D. Sanders, E. Falconer, M. Hills, D. C. Spierings, P. M. Lansdorp, Single-cell
771 template strand sequencing by Strand-seq enables the characterization of individual
772 homologs. *Nature Protocols* **12**, 1151-1176 (2017).

773 24. P. Marks *et al.*, Resolving the full spectrum of human genome variation using Linked-
774 Reads. *Genome research* **29**, 635-645 (2019).

775 25. G. X. Zheng *et al.*, Haplotyping germline and cancer genomes with high-throughput
776 linked-read sequencing. *Nature biotechnology* **34**, 303-311 (2016).

777 26. F. Zhang *et al.*, Haplotype phasing of whole human genomes using bead-based
778 barcode partitioning in a single tube. *Nature biotechnology* **35**, 852-857 (2017).

779 27. D. Lin *et al.*, Digestion-ligation-only Hi-C is an efficient and cost-effective method for
780 chromosome conformation capture. *Nature genetics* **50**, 754-763 (2018).

781 28. L. Harewood *et al.*, Hi-C as a tool for precise detection and characterisation of
782 chromosomal rearrangements and copy number variation in human tumours.
783 *Genome biology* **18**, 1-11 (2017).

784 29. H. Lu, F. Giordano, Z. Ning, Oxford Nanopore MinION sequencing and genome
785 assembly. *Genomics, proteomics & bioinformatics* **14**, 265-279 (2016).

786 30. S. Yeo, L. Coombe, R. L. Warren, J. Chu, I. Birol, ARCS: scaffolding genome drafts

787 with linked reads. *Bioinformatics* **34**, 725-731 (2018).

788 31. S. B. Kingan *et al.*, A high-quality de novo genome assembly from a single mosquito
789 using PacBio sequencing. *Genes* **10**, 62 (2019).

790 32. A. Rhoads, K. F. Au, PacBio sequencing and its applications. *Genomics, proteomics*
791 *& bioinformatics* **13**, 278-289 (2015).

792 33. S. Heinz *et al.*, Transcription elongation can affect genome 3D structure. *Cell* **174**,
793 1522-1536. e1522 (2018).

794 34. V. Deshpande *et al.*, Exploring the landscape of focal amplifications in cancer using
795 AmpliconArchitect. *Nature communications* **10**, 1-14 (2019).

796 35. H. Kim *et al.*, Extrachromosomal DNA is associated with oncogene amplification and
797 poor outcome across multiple cancers. *Nature genetics* **52**, 891-897 (2020).

798 36. J. Luebeck *et al.*, AmpliconReconstructor integrates NGS and optical mapping to
799 resolve the complex structures of focal amplifications. *Nature communications* **11**, 1-
800 14 (2020).

801 37. S. E. Johnstone *et al.*, Large-scale topological changes restrain malignant
802 progression in colorectal cancer. *Cell* **182**, 1474-1489. e1423 (2020).

803 38. T. Knutsen *et al.*, Definitive molecular cytogenetic characterization of 15 colorectal
804 cancer cell lines. *Genes, Chromosomes and Cancer* **49**, 204-223 (2010).

805 39. J. R. Dixon *et al.*, Topological domains in mammalian genomes identified by analysis
806 of chromatin interactions. *Nature* **485**, 376-380 (2012).

807 40. O. S. Chapman *et al.*, The landscape of extrachromosomal circular DNA in
808 medulloblastoma. *bioRxiv*, (2021).

809 41. R. P. Koche *et al.*, Extrachromosomal circular DNA drives oncogenic genome
810 remodeling in neuroblastoma. *Nature genetics* **52**, 29-34 (2020).

811 42. H. D. Møller *et al.*, Genome-wide purification of extrachromosomal circular DNA from
812 eukaryotic cells. *JoVE (Journal of Visualized Experiments)*, e54239 (2016).

813 43. S. Heim, F. Mitelman, *Cancer cytogenetics: chromosomal and molecular genetic*
814 *aberrations of tumor cells*. (John Wiley & Sons, 2015).

815 44. Q. An *et al.*, Variable breakpoints target PAX5 in patients with dicentric chromosomes:
816 a model for the basis of unbalanced translocations in cancer. *Proceedings of the*
817 *National Academy of Sciences* **105**, 17050-17054 (2008).

818 45. J. Pedersen-Bjergaard, J. D. Rowley, The balanced and the unbalanced chromosome
819 aberrations of acute myeloid leukemia may develop in different ways and may
820 contribute differently to malignant transformation. (1994).

821 46. F. Dicker, S. Schnittger, T. Haferlach, W. Kern, C. Schoch, Immunostimulatory
822 oligonucleotide-induced metaphase cytogenetics detect chromosomal aberrations in
823 80% of CLL patients: a study of 132 CLL cases with correlation to FISH, IgVH status,
824 and CD38 expression. *Blood* **108**, 3152-3160 (2006).

825 47. S. J. Morin, J. Eccles, A. Iturriaga, R. S. Zimmernan, Translocations, inversions and
826 other chromosome rearrangements. *Fertility and sterility* **107**, 19-26 (2017).

827 48. R. M. Gardner, G. R. Sutherland, L. G. Shaffer, *Chromosome abnormalities and*
828 *genetic counseling*. (OUP USA, 2011).

829 49. S. Andrews. (Babraham Bioinformatics, Babraham Institute, Cambridge, United
830 Kingdom, 2010).

831 50. S. Chen, Y. Zhou, Y. Chen, J. Gu, fastp: an ultra-fast all-in-one FASTQ preprocessor.
832 *Bioinformatics* **34**, i884-i890 (2018).

833 51. H. Li, R. Durbin, Fast and accurate short read alignment with Burrows–Wheeler
834 transform. *bioinformatics* **25**, 1754-1760 (2009).

835 52. H. Li *et al.*, The sequence alignment/map format and SAMtools. *Bioinformatics* **25**,
836 2078-2079 (2009).

837 53. A. Tarasov, A. J. Vilella, E. Cuppen, I. J. Nijman, P. Prins, Sambamba: fast
838 processing of NGS alignment formats. *Bioinformatics* **31**, 2032-2034 (2015).

839 54. A. McKenna *et al.*, The Genome Analysis Toolkit: a MapReduce framework for
840 analyzing next-generation DNA sequencing data. *Genome research* **20**, 1297-1303
841 (2010).

842 55. R. Xi, S. Lee, Y. Xia, T.-M. Kim, P. J. Park, Copy number analysis of whole-genome
843 data using BIC-seq2 and its application to detection of cancer susceptibility variants.
844 *Nucleic acids research* **44**, 6274-6286 (2016).

845 56. N. C. Durand *et al.*, Juicer provides a one-click system for analyzing loop-resolution
846 Hi-C experiments. *Cell systems* **3**, 95-98 (2016).

847 57. K. Kim, M. Kim, Y. Kim, D. Lee, I. Jung, in *Seminars in Cell & Developmental Biology*.
848 (Elsevier, 2021).

849 58. S. Wang *et al.*, HiNT: a computational method for detecting copy number variations
850 and translocations from Hi-C data. *Genome biology* **21**, 1-15 (2020).

851 59. T. Rausch *et al.*, DELLY: structural variant discovery by integrated paired-end and
852 split-read analysis. *Bioinformatics* **28**, i333-i339 (2012).

853 60. D. L. Cameron *et al.*, GRIDSS: sensitive and specific genomic rearrangement
854 detection using positional de Bruijn graph assembly. *Genome Res* **27**, 2050-2060
855 (2017).

856 61. D. C. Jeffares *et al.*, Transient structural variations have strong effects on quantitative
857 traits and reproductive isolation in fission yeast. *Nat Commun* **8**, 14061 (2017).

858 62. A. Zirkel *et al.*, HMGB2 loss upon senescence entry disrupts genomic organization
859 and induces CTCF clustering across cell types. *Mol. Cell* **70**, 730–744.e736 (2018).

860 63. P. Euskirchen *et al.*, Same-day genomic and epigenomic diagnosis of brain tumors
861 using real-time nanopore sequencing. *Acta Neuropathol* **134**, 691-703 (2017).

862 64. Y. Liao, G. K. Smyth, W. Shi, featureCounts: an efficient general purpose program for
863 assigning sequence reads to genomic features. *Bioinformatics* **30**, 923-930 (2014).

864 65. M. I. Love, W. Huber, S. Anders, Moderated estimation of fold change and dispersion
865 for RNA-seq data with DESeq2. *Genome Biol* **15**, 550 (2014).

866

(A) Flowchart of MGA-Seq. Nuclei were cross-linked with 0.5% formaldehyde and then digested with HindIII. 5' DNA overhangs of digested chromatin fragments were filled in by DNA polymerase and then proximity ligated by T4 DNA ligase. The proximity ligation products were fragmented and then subjected to high-throughput sequencing library construction. After sequencing, all the reads were used to generate chromatin contact matrix for genome structural variation calling. In the sequencing library, the reads without ligation junction "AAGCTAGCTT" were used for the detection of CNV, SNP, small indels (< 50bp), region of focal amplification, and genome breakpoints. By combining all information, the types and breakpoints of structural variation can be decoded. Notably, MGA-Seq can distinguish ecDNA and HSR, predict the structure of simple focal amplification regions, and construct the interaction network of focal amplified genes. **(B)** The main steps, time, and cost of MGA-Seq.

Figure 2

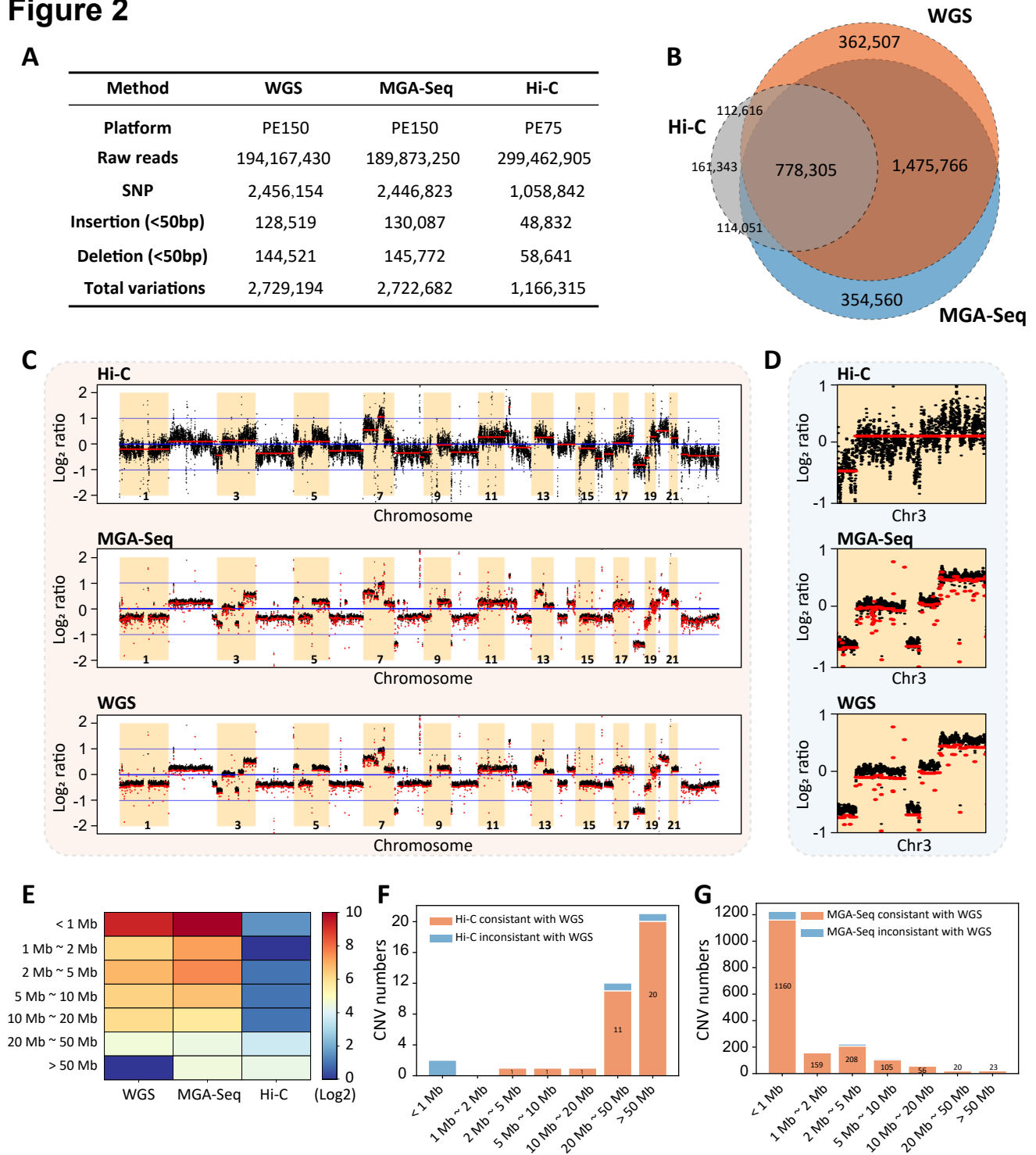
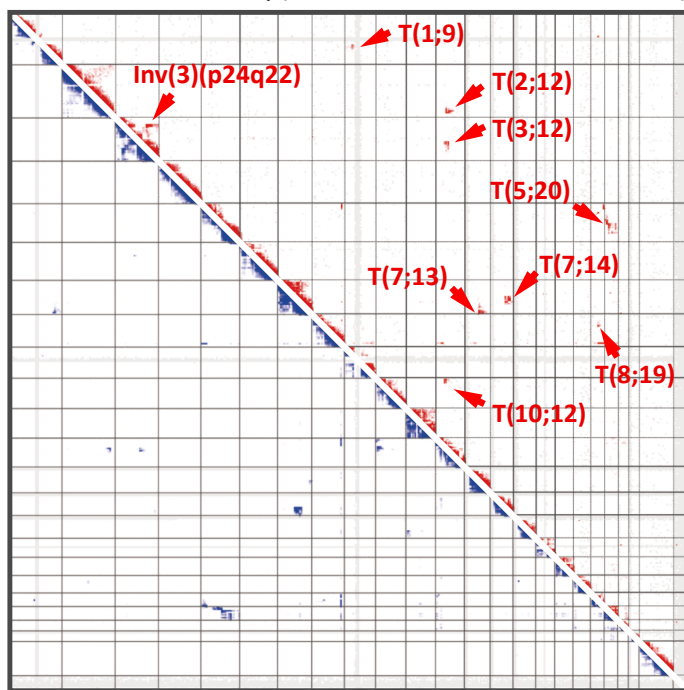
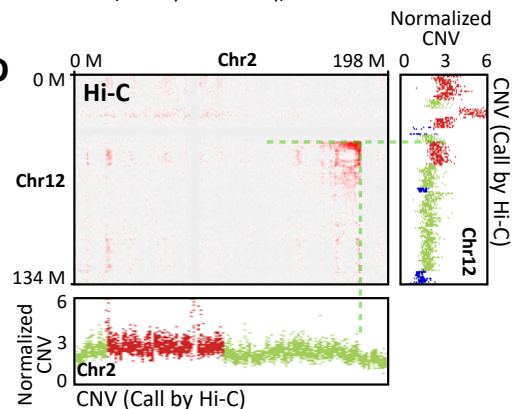
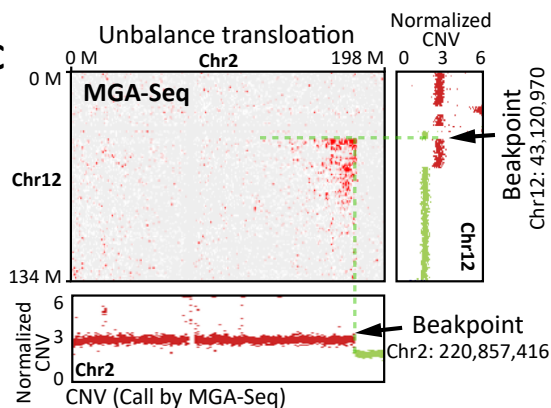


Figure 2. Detection of SNPs, indels, and CNVs by MGA-Seq. (A) Comparison of the numbers of SNPs and indels (< 50bp, include insertions and deletions) identified by WGS, MGA-Seq, and Hi-C. **(B)** Overlap of the SNPs and indels between MGA-Seq, WGS, and Hi-C. **(C)** Comparison of \log_2 copy ratios calculated using reads coverage between Hi-C, MGA-Seq, and WGS. **(D)** Comparison of the CNVs on chromosome 3 identified by Hi-C, MGA-Seq, and WGS. **(E)** Statistics of the number and size distribution of CNVs identified by Hi-C, MGA-Seq, and WGS. **(F)** Consistency of the CNV segments (categorized by size) detected by Hi-C and WGS. Overall, Hi-C cannot detect CNV with length less than 20 Mb. **(G)** Consistency of the CNV segments detected by MGA-Seq and WGS. The number and size distribution of CNV segments detected by MGA-Seq and WGS are highly consistent, especially for micro-CNVs (< 1Mb).

A



C



B

SV types	WGS	Hi-C			MGA-Seq		
	Identification by Delly2 & Gridss	Identification by heatmap	Translocation type identify by CNV	Translocation Breakpoint	Identification by heatmap	Translocation type identify by CNV	Translocation Breakpoint
Inv(3)(p24q22)	✗	✓	N/A	Not detected	✓	N/A	Not detected
T(1;9)(q22;p13)	✗	✓	✗	Not detected	✓	Unbalance transloation	Not detected
T(2;12)(q35;q12)	✗	✓	✗	Not detected	✓	Unbalance translocation	Chr2: 220,857,417 Chr12: 43,120,970
T(3;12)(q13;q14)	✗	✓	✗	Not detected	✓	Unbalance translocation	Chr3:109,985,783 Chr12:61,209,043
T(7;13)(q36;q21)	✗	✓	✗	Not detected	✓	Unbalance translocation	Chr7: 153,635,126 Chr13: 60,073,622
T(7;14)(q11;q23)	✗	✓	✗	Not detected	✓	Unbalance translocation	Chr7: 76,973,295 Chr14: 67,665,593
T(10;12)(p15;q12)	✗	✓	✗	Not detected	✓	Unbalance translocation	Not detected
T(5;20)(q15;p12)	✗	✓	✗	Not detected	✓	Unbalance translocation	Chr5: 93,821,190 Chr20: 17,313,991
T(5;19)(p14;q13)	✗	✓	✗	Not detected	✓	Unbalance translocation	Chr5: not detect Chr19: 59,060,484

E

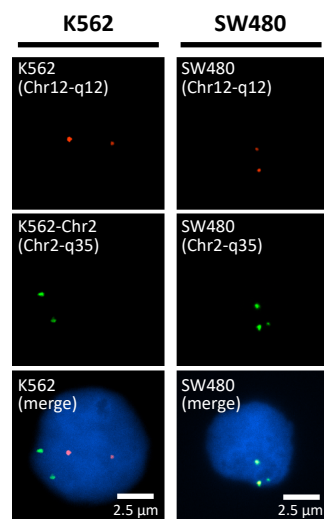


Figure 3. Identification of translocation types and breakpoints in SW480 at single base-pair resolution by MGA-Seq. (A) Identification of translocation in the SW480 cell line by genomic contact matrix constructed with MGA-Seq datasets. The detected structural variations are indicated by arrows. **(B)** Translocation types and breakpoint information defined by MGA-Seq. **(C)** Application of integrated chromatin contact matrix, CNVs, and split reads analysis to identify translocation types and breakpoints between chr 2 and chr 12 at single base-pair resolution using MGA-Seq datasets. **(D)** Identification of translocation types and breakpoints between chr 2 and chr 12 using Hi-C datasets. **(E)** Validation of the T(2;12)(q35;q12) translocation in SW480 cells by DNA FISH. FISH probes for 12q12 and 2q35 were directly labeled with Alexa Fluor 555 (red) and Alexa Fluor 488 (green), respectively. K562 cells without the T(2;12)(q35;q12) translocation were used as a control.

Figure 4

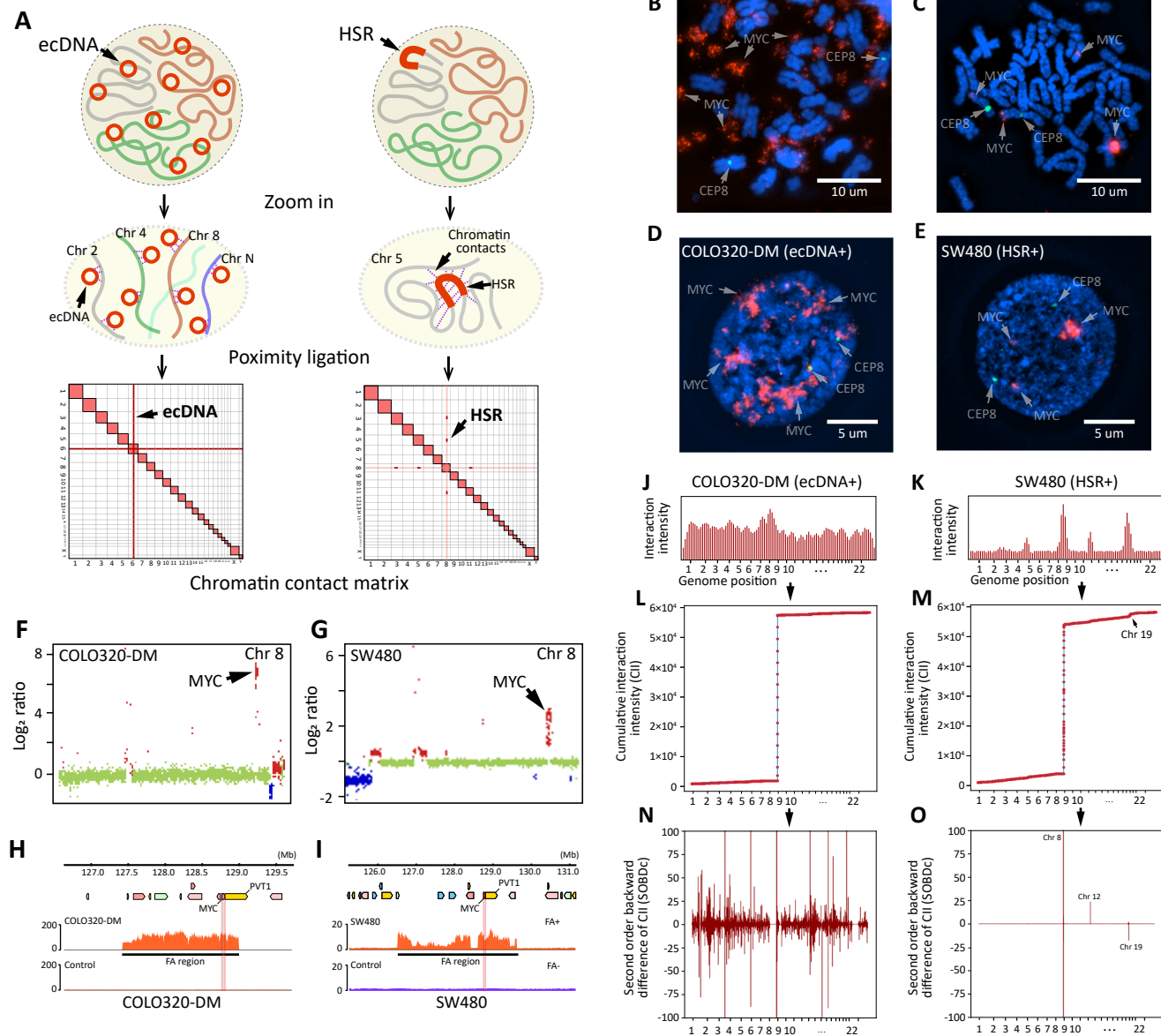
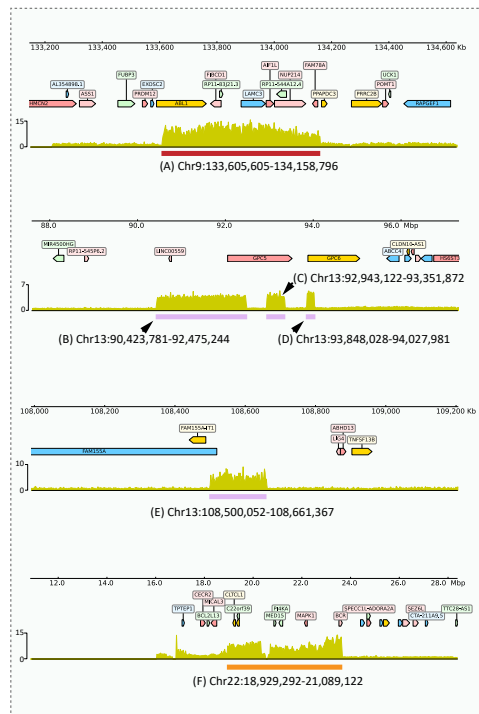


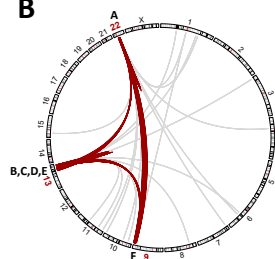
Figure 4. Identification of ecDNA and HSR by MGA-Seq. (A) Putative diagram of inter-chromosomal interaction pattern differences between ecDNA and HSR positive cell line. (B-E) Validation of MYC amplification in COLO320-DM and SW480 cell lines by DNA FISH. The red signal represents MYC and the green signal represents the centromere of chr 8. (F and G) Copy number variation analysis of chr 8 in COLO320-DM and SW480 cell lines. Gains and losses of copy number are shown in red and blue, respectively. (H and I) Location of the MYC amplification region in COLO320-DM and SW480 cell lines. (J and K) Interaction intensity between the focal amplified region and whole genome. (L and M) Cumulative interaction intensity curve of COLO320-DM and SW480 cell lines. The x-axis represents the genome position, 100 kb bin size. The Y axis represents the accumulation of interaction intensity. (N and O) Plotted the second order backward difference (SOBD) value across the genome of COLO320-DM and SW480 cell lines in 100-kb bin size.

Figure 5

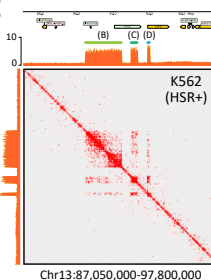
A



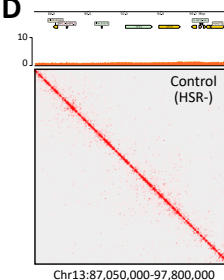
B



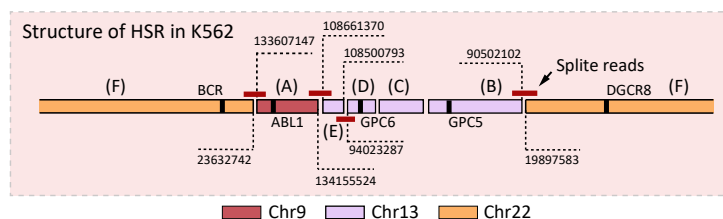
C



D



E



F

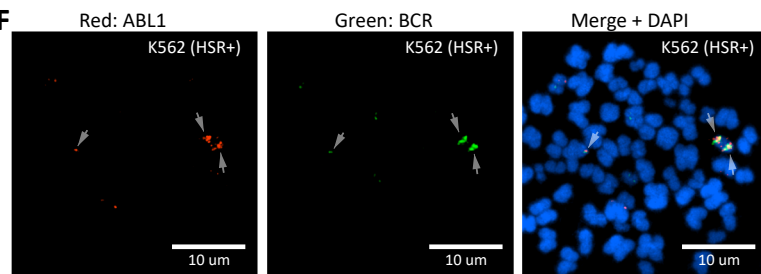


Figure 5. Deciphering the spatial structure of homogenously staining region (HSR) in K562 cell line. (A) Location of the amplification region on chr 9, 13, and 22. **(B)** Circos plots of the chromatin interactions mediated by amplification regions across all 23 chromosomes in K562 cell lines. The interactions between chromosomes 9, 13, and 22 are marked with red lines. **(C and D)** Comparison of chromatin contact matrix of amplification region (Chr13:90423781-92475244, Chr13:92943122-93351872, and Chr13:93848028-94027981) between K562 cell line and healthy human peripheral blood cells (control). **(E)** Assembling the amplified regions from "A" to "F" with split reads. The breakpoint of the amplification regions is marked in the figure. **(F)** Metaphase analysis and DNA FISH to validate the location of the ABL1 amplification region and the BCR amplification region in K562 cell line. FISH probes for the ABL1 amplification region and the BCR amplification region were directly labeled with Alexa Fluor 555 (red) and Alexa Fluor 488 (green), respectively.

Figure 6

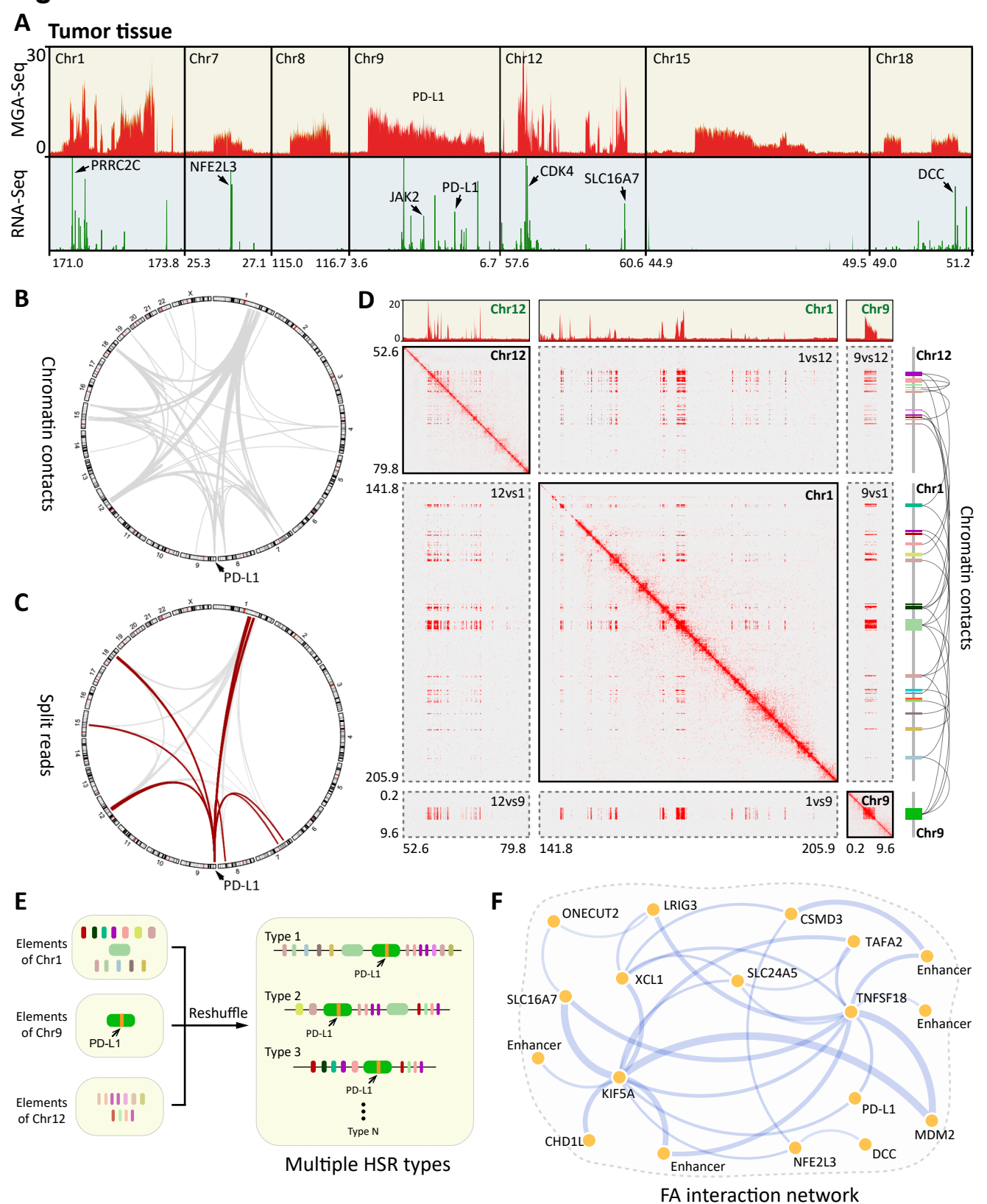


Figure 6. Heterogeneity of focal amplification in renal cancer tissue. **(A)** Sequencing reads coverage and RNA expression level in typical focal amplification regions of a renal cancer tissue sample. **(B)** Circos plots of the chromatin interactions mediated by focal amplification regions across all 23 chromosomes in renal cancer tissue. **(C)** Circos plots of the split reads mediated by focal amplification regions across all 23 chromosomes. The split reads aligned to the PD-L1 amplified region are marked with red lines. **(D)** Chromatin contact matrix between the amplified regions of chr1, chr9, and chr12, and sequencing reads coverage within these amplified regions. **(E)** Diverse structures of HSR in the renal cancer tissue sample. **(F)** Interaction network of amplified oncogenes in the renal cancer tissue sample. Different amplified oncogenes are assembled by split reads. The thickness of the line indicates the chromatin contact strength.