

# Genomics, Proteomics & Bioinformatics

## Long-Range NGS Linked-Reads and Applications of Hi-C, 10x, Haplotagging and TELL-Seq Platforms --Manuscript Draft--

<b>Manuscript Number:</b>	GPB-D-22-00474
<b>Article Type:</b>	Review Article
<b>Keywords:</b>	Long-range NGS reads, Hi-C, 10x, haplotagging, TELL-Seq, genome assembly and quality assessment
<b>Corresponding Author:</b>	Zemin Ning, PhD The Wellcome Sanger Institute ambridge, UNITED KINGDOM
<b>First Author:</b>	Zemin Ning, PhD
<b>Order of Authors:</b>	Zemin Ning, PhD
<b>Abstract:</b>	<p>Long-range reads grant insight into additional genetic information, from the original DNA samples, far beyond what can be accessed by short reads, or even modern long-read technology. Several new sequencing technologies have become available for long-range "linked reads" with high-throughput and high-resolution genome analysis. These long-range technologies are rapidly advancing the fields of genome assembly, genome scaffolding, and allowing more comprehensive variant identification. In this article, we focus on four major long-range sequencing technologies: Hi-C, 10x Genomics Linked-Reads, Haplotagging and TELL-Seq. We detail the mechanisms and data products of these four platforms, introduce their most important applications, evaluate the quality of sequencing data from different platforms, and discuss the currently available bioinformatics tools. We hope this work will benefit the selection of appropriate long-range technology for specific biological studies.</p>
<b>Suggested Reviewers:</b>	<p>Yingguang Frank Chan, PhD Professor, Max Planck Institute for Biology Tübingen frank.chan@tue.mpg.de Dr Chan invited Haplotagging and is a leading expert in long-range data analysis.</p> <p>Joana Isable Meier, PhD Group Leader, University of Cambridge and Sanger jm2276@cam.ac.uk Haplotagging expert</p> <p>Chenxi Zhou, PhD Postdoc, University of Cambridge cz370@cam.ac.uk Dr Zhou developed YaHS, a popular genome scaffolding tool which has been widely used.</p> <p>Arang Rhie, PhD Staff scientist, NIH/NHGRI arang.rhir@nih.gov Dr Rhie developed a number of widely used informatics tools, including Merqury for assembly and data QC.</p>
<b>Opposed Reviewers:</b>	



Dear Editor,

We are pleased to submit our manuscript entitled “Long-Range NGS Linked-Reads and Applications of Hi-C, 10x, Haplotagging and TELL-Seq Platforms” by Jiang *et al.*, for publication as a review/research article in Genomics, Proteomics & Bioinformatics.

With rapid development in sequencing technologies, long-range data types and their applications have gained substantial momentum in the genomic community. Long-range reads grant insight into additional genetic information, from the original DNA samples, far beyond what can be accessed by short reads, or even modern long-read technology. The roles of these datasets in genome scaffolding, consensus base polishing, phasing as well as structural variation detection are simply unreplacable, without which many high-profile projects such as VGP, Darwin Tree of Life and Earth BioGenome won't be able to achieve the targeted objectives. Currently, there is not a single published paper which discusses the common features, data evaluations and applications of long-range data, when looking at literature. We hope this paper fills this important gap and is of great interest in both genomics and bioinformatics.

This is an invited review paper and we haven't submitted this work to other journals. All the authors have seen and approved the submitted version of the manuscript. Please feel free to contact me, should you have any questions.

Dr Zemin Ning

A handwritten signature in black ink, appearing to be "Zemin Ning".

Senior Scientific Manager  
High Performance Algorithms  
The Wellcome Sanger Institute  
Wellcome Genome Campus  
Hinxton, Cambridge CB10 1SA  
UK  
Tel: (44) 1223 494705

E-mail: [zn1@sanger.ac.uk](mailto:zn1@sanger.ac.uk)



Dear Editor,

We would like to suggest these outstanding researchers with expertise in individual areas as referees for this manuscript:

1. Haplotagging  
Professor Yingguang Frank Chan  
Max Planck Institute in Tuebingen  
[frank.chan@tue.mpg.de](mailto:frank.chan@tue.mpg.de)
2. Haplotagging:  
Dr Joana Isabel Meier\*  
University of Cambridge  
[jm2276@cam.ac.uk](mailto:jm2276@cam.ac.uk)
3. Hi-C and genome scaffolding  
Dr Chenxi Zhou  
University of Cambridge  
[cz370@cam.ac.uk](mailto:cz370@cam.ac.uk)
4. Hi-C and 10x  
Dr Arang Rhie  
Genome Informatics Section at NIH/NHGRI,  
[arang.rhir@nih.gov](mailto:arang.rhir@nih.gov)

\*Dr Meier also has a group at Sanger now with Tree of Life Project. But work wise, I have no overlaps with her. Haplotagging is new technology and there are not so many experts in the community right now. Given this is a review paper, I think it is legitimate to invite referees who are from the same institute.

This is an invited review paper and we haven't submitted this work to other journals. All the authors have seen and approved the submitted version of the manuscript. Please feel free to contact me, should you have any questions.

Dr Zemin Ning

A handwritten signature in black ink, appearing to be "Zemin Ning", written over a light blue rectangular background.

Senior Scientific Manager  
High Performance Algorithms

The Wellcome Sanger Institute  
Wellcome Genome Campus  
Hinxton, Cambridge CB10 1SA  
UK

Tel: (44) 1223 494705

E-mail: [zn1@sanger.ac.uk](mailto:zn1@sanger.ac.uk)

# 1 Long-Range NGS Linked-Reads and Applications of Hi-C, 2 10x, Haplotagging and TELL-Seq Platforms

3

4 Libo Jiang<sup>1</sup>, Michael A. Quail<sup>2</sup>, Jack Fraser-Govi<sup>1,2</sup>, Haipeng Wang<sup>1</sup>, Xuequn Shi<sup>3</sup>, Karen Oliver<sup>2</sup>,  
5 Esther Mellado Gomez<sup>2</sup>, Fengtang Yang<sup>1</sup> and Zemin Ning<sup>2+</sup>

6

7 <sup>1</sup>School of Biological Sciences, Shandong University of Technology, Zibo, 255049, Shandong, China.

8 <sup>2</sup>The Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SA, UK

9 <sup>3</sup>College of Food Science and Technology, Hainan University, Hainan, 570228, China

10

11 <sup>+</sup>Corresponding author. E-mail: zn1@sanger.ac.uk (Ning Z).

12

## 13 **Abstract**

14 Long-range reads grant insight into additional genetic information, from the original  
15 DNA samples, far beyond what can be accessed by short reads, or even modern long-  
16 read technology. Several new sequencing technologies have become available for  
17 long-range “linked reads” with high-throughput and high-resolution genome analysis.  
18 These long-range technologies are rapidly advancing the fields of genome assembly,  
19 genome scaffolding, and allowing more comprehensive variant identification. In this  
20 article, we focus on four major long-range sequencing technologies: Hi-C, 10x  
21 Genomics Linked-Reads, Haplotagging and TELL-Seq. We detail the mechanisms  
22 and data products of these four platforms, introduce their most important applications,  
23 evaluate the quality of sequencing data from different platforms, and discuss the  
24 currently available bioinformatics tools. We hope this work will benefit the selection  
25 of appropriate long-range technology for specific biological studies.

26

## 27 **Keywords**

28 Long-range NGS reads, Hi-C, 10x, haplotagging, TELL-Seq, genome assembly and  
29 quality assessment

## 30 Introduction

31 Next-Generation Sequencing (NGS) technologies have revolutionized the field of  
32 genomics and genetics, providing low-cost and high-throughput data at an  
33 unprecedented scale. However, most NGS technologies make an underlying  
34 assumption that all relevant genetic information can be reconstructed from the smaller  
35 fragments that make up both short (100-250bp) and long (>10,000bp) reads. Such  
36 reads are ‘short range’ or ‘local’, because they contain only information about the  
37 genetic sequences of the reads, in contrast to ‘long-range’, ‘non-local’ or ‘linked’  
38 reads, which retain additional contextual information regarding the origin of the read  
39 within the complex, 3-dimensional physical structure of the DNA within and between  
40 chromosomes.

41 We emphasise that, despite the similar terminology, long-range reads are  
42 conceptually distinct from long reads. Although the size of long reads provides a large  
43 quantity of information, it is inherently local, relating only to the sequence without  
44 containing additional information about the origin of the fragment. In contrast, long-  
45 range reads provide additional non-local information, and can take the forms of both  
46 short and long reads, though in practice most long-range technologies currently use  
47 short reads. To avoid confusion, we suggest deviating from the literature standard and  
48 instead defining ‘long-range’ reads as either non-local reads, or linked reads.

49 Without the additional context of non-local information, for example, it remains  
50 challenging to reliably identify structural variation (SV) with short reads. Although  
51 short reads can identify SVs to base-pair resolution, utilizing only short-range  
52 information suffers from a higher false discovery rate than long reads [1]. It is also  
53 difficult to phase many millions of short reads to a haplotype-resolved genome,  
54 particularly for highly repetitive sequences, complex heterozygosity, and large  
55 polyploid genomes [2, 3].

56 Local long reads can sidestep many of the issues associated with local short reads  
57 although they contain only local information, because the large size of the reads  
58 makes it much easier to uniquely localize them within the genome [4, 5] [6,7].  
59 Currently, two major long-read technologies: Pacific Biosciences (PacBio) single-  
60 molecule real-time (SMRT) sequencing and Oxford Nanopore Technologies (ONT)  
61 nanopore sequencing are used for long-read genome analysis [4]. However, these

62 long-read methods have two drawbacks: (i) higher costs and lower throughput and (ii)  
63 higher DNA input requirements compared to short-read sequencing.

64 Whilst long-range information can be used in isolation for *de novo* assembly,  
65 long-range reads have already proven vastly more powerful since the large,  
66 contiguous reads make referenceless assembly much easier. However, the throughput  
67 and cost issues associated with long reads mean that using them as the sole means of  
68 long-range information in the *de novo* assembly of tens of thousands of genomes will  
69 likely be prohibitive. This is especially true if chromosome-level assemblies are  
70 desired, since long reads are still much smaller than chromosomes and hence do not  
71 carry chromosome-scale context. Therefore, a cheaper method for inferring long-  
72 range information is needed.

73 Several mechanisms for storing long-range information within short reads and  
74 hence the necessary context to reconstruct a single long molecule of DNA have been  
75 developed, including the “Pair-Linked Reads” (PLR) chromosome conformation  
76 capture-based Hi-C and “Chain-Linked Reads” (CLR) technologies [8]. 10x  
77 Genomics provides perhaps the best known chained read strategy, which can generate  
78 long-range information from standard approaches based on short reads [9]. In recent  
79 years, 10x Genomics Chain-Linked Read technology has been widely used, but a  
80 variety of other barcode-based methods such as TruSeq, BGI’s Long Fragment Reads,  
81 TELL-Seq, LoopSeq and haplotagging have been developed with ultralow DNA  
82 input, high per-base resolution, and low costs [10, 11, 12, 13].

83 Here, we focus on the four major long-range reads sequencing technologies, i.e.,  
84 Hi-C, 10x Genomics Linked-Reads, haplotagging and TELL-Seq. Firstly, we detail  
85 the protocols and mechanisms of the four platforms’ function. Secondly, we propose  
86 some criteria to evaluate the quality of sequencing data on different platforms, and  
87 apply these criteria to discuss the characteristics of datasets either downloaded from  
88 public resources or generated by us. Thirdly, we review the practical applications of  
89 these technologies in efforts such as genome scaffolding, *de novo* assemblies and  
90 variation screens. Finally, we provide a list of software tools which are commonly  
91 used for genome analysis with long-range reads and discuss their strengths and  
92 weaknesses.

## 93 Platforms

94 In this section, we briefly detail the four platforms of interest: Hi-C, 10x, Haplotgging  
95 and TELL-Seq, focusing on the protocols used to generate the long-range data, and  
96 how such long-range information is manifested in the data products. Across these four  
97 platforms, long-range, non-local information is stored in ‘linked reads’ in one of two  
98 ways: either in Pair-Linked Reads, in which two reads are coupled together to indicate  
99 a relationship between them, or Chain-Linked Reads, in which reads are tagged or  
100 labelled in some ways to indicate their origins. This qualitatively changes the non-  
101 local information provided by the platform, and hence informs which platform is  
102 suitable for a given application.

### 103 Hi-C

104 Hi-C is a Pair-Linked Read technology, and the culmination of several generations of  
105 Chromosome Conformation Capture technologies, which uses PLR to probe  
106 chromosome conformation – the spatial organization of chromatin within a cell – at a  
107 genome-wide scale [14], granting access to 3D proximity information within the  
108 nucleus. Since chromatin is a complex 3-dimensional structure, this information  
109 allows researchers to detect long-range interactions between segments within a  
110 chromosome or between different chromosomes. Since homologous chromosomes  
111 each tend to occupy distinct territories in nuclei [15], this feature enables the use of  
112 Hi-C data to improve *de novo* assembly, and phase heterozygous genome variants  
113 onto haplotypes,

114 Hi-C technology follows this protocol (see graphical summary in Fig.1A):

115

- 116 ● The nuclear chromatin is crosslinked using formaldehyde. By design, these  
117 crosslinks occur preferentially between strands that are close together in 3D space.
- 118 ● Crosslinked chromatin is solubilized and fragmented with a restriction enzyme
- 119 ● The crosslinked segment ends are repaired by filling in with biotin-labeled  
120 nucleotides.
- 121 ● DNA ligase is used to cyclize the blunt-end components, the proteins that bind the  
122 DNA fragments are degraded, and then the circular crosslinked fragments are



123 randomly broken again using sonication or other methods.  
124 ● DNA is purified and sheared. The biotin-labeled DNAs are captured with  
125 streptavidin-conjugated beads and amplified before sequencing.

126

127 The final result of this protocol is a large number of “deliberate chimeric”, paired  
128 short reads, with each end of the read originating from one of the crosslinked strands,  
129 which are potentially very far away from each other in the linear genome – and even  
130 on different chromosomes entirely. The generated paired reads are then mapped to a  
131 contig assembly of the genome and used to create a high-resolution interaction map  
132 within and between chromosomes: regions where larger numbers reads are found to  
133 have been crosslinked are then inferred to be regions of close contact between the  
134 DNA strands.

### 135 **10x Genomics linked-reads**

136 10x Genomics Linked-Reads (henceforth simply ‘10x’) are a product formerly  
137 provided by 10x Genomics. In 10x sequencing, long-range information is retained by  
138 combining 3’ barcoding with standard short-read sequencing [16], producing short  
139 Chain-Linked Reads with a ‘memory’ of the larger scale locality where they were  
140 derived from, and hence making it easier to assemble the resulting dataset. The  
141 resulting reads can improve the quality of genome assembly by expanding the range  
142 of linking information along the chromosome to define haplotypes. The 10x protocol  
143 (demonstrated in Fig.1B) is as follows:

- 144 ● First, high molecular weight (HMW) DNA is prepared and sheared into long DNA  
145 fragments (ideally > 100kb)
- 146 ● At two microfluidic junctions, tHMW-DNA is then combined with an oil-surfactant  
147 solution, enzymes and gel beads loaded with random primers and barcode  
148 sequences to produce “Gel Beads in Emulsion” (GEMs). Each GEM captures  
149 around 10 HMW-DNA fragments.
- 150 ● The GEMs are then isolated in partitions and the beads dissolved, releasing the  
151 barcodes and primers uniquely to the HMW-DNA fragments in that partition.
- 152 ● Each partition is then sheared, extended with both the barcodes and the primers,  
153 and then amplified and sequenced – in this case, by Illumina paired-end sequencing.

154 The end result of this process is a number of short reads preappended by a unique  
155 barcode identifying the GEM bead they originated from: all reads sharing a common  
156 barcode are called ‘Linked-Reads’ (which we distinguish as being distinct from the  
157 general term linked reads: under our terminology these are Chain-Linked Reads, a  
158 specific form of linked reads). The key statistic is that, since each GEM captures so  
159 few HMW-DNA fragments, the odds that a second fragment which shares the same  
160 barcode also originates from nearby in the genome is very small (see Supplementary  
161 Figure 2), and so the barcode acts to (nearly) uniquely group sets of reads together as  
162 being spatially co-located. This, for example, makes it much easier to phase short  
163 reads as the entire barcoded molecule must be simultaneously phased.

164 Although 10x sequencing can reconstruct multi-megabase phase blocks by  
165 assembling short reads with barcode information, it still has some drawbacks, such as  
166 relatively high costs in library preparation, and that the 10x platform performs  
167 counterintuitively when faced with smaller genomes, showing a marked performance  
168 degradation. This is because the partitions get saturated by the smaller genome size,  
169 and the statistics begin to favour ‘barcode collisions’ much more frequently. The 10x  
170 platform is optimized for the human genome size, and modifications such as smaller  
171 DNA samples would need to be made for non-human cases.

172 Most significantly, however, this product was been withdrawn and discontinued  
173 in 2020 [17]. However, we include this platform in our analysis for continuity with  
174 previous benchmarking and comparison efforts, and since future 10x Genomics  
175 products may be comparable to this previous iteration.

176

## 177 **Haplotagging**

178 Several other technologies have been developed to provide an alternative form of  
179 CLR in the absence of 10x. Haplotagging, as a simple and relatively low-cost Chain-  
180 Linked Read sequencing technique was developed by Meier et al. [13]. This technique  
181 allows high-throughput sequencing without losing haplotype information while  
182 maintaining the power, accuracy, and scalability of standard Illumina sequencing.

183 Haplotagging, like TELL-Seq mentioned below, is a transposon bead-based  
184 technology that employs transposomes containing bead-specific barcoded adaptors.  
185 These technologies utilise the tendency of segments of HMW-DNA to wrap around

186 microbeads, providing many points of contact between the bead and the DNA. The  
187 full protocol is as follows (Fig.1C):

- 188 ● As in the 10x protocol, HMW-gDNA (ideally >100kb) is prepared.
- 189 ● The HMW-gDNA is mixed with the barcoded beads. Each bead carries a standard  
190 Illumina Nextera Tn5 transposon adaptor, augmented with one of 85 million  
191 barcodes, and each bead captures only a single DNA fragment
- 192 ● Transposition transfers the barcoded adaptors into the long DNA fragments, before  
193 PCR amplification to generate a sequencing-ready library.
- 194 ● Finally, the libraries were sequenced using an Illumina platform.

195

196 The result is that the initial HMW-DNA fragments are broken into smaller units, each  
197 containing a unique barcode, that can be sequenced on short read sequencers.

198 Subsequently, all the reads originating from the same HMW-DNA fragment can  
199 grouped by their barcode, and hence correctly mapped to the same fragment.

200 The key difference between haplotagging technology and 10x is that DNA molecules  
201 tend to interact only with a single bead, instead of the approximately 10 (for humans)  
202 fragments-per-bead that 10x relied on. In addition, each bead is tagged with four  
203 barcode fragments that are distributed in the standard i5/7 index positions of the  
204 Illumina Nextera adaptor design. Thus, library preparation and barcoding are  
205 performed simultaneously within the same tube, making the process cheap and easy to  
206 produce using standard molecular biology equipment. The data output is very similar  
207 to that of the 10x platform: a series of short reads preappended by a barcode,  
208 indicating which reads originated from a similar vicinity. However, as mentioned  
209 above, the process is much cheaper (the original work claimed a 99% cost reduction);  
210 and since the fragment/bead interaction is close to 1:1, instead of approximately 10:1,  
211 each fragment is genuinely uniquely barcoded, resulting in fewer barcode collisions,  
212 as demonstrated in Supplementary Figure 2. In addition, the 4-fragment nature of the  
213 barcode is designed to allow for error-correction in the barcode reads, allowing for  
214 more robust identification of the barcode. However, the fragments are prone to  
215 display PCR duplication errors [18] and the product is not yet at the stage of  
216 commercial deployment.

## 217 **TELL-Seq**

218 TELL-Seq<sup>1</sup> is another CLR sequencing technology which functions very similarly to  
219 the Haplotagging platform but is currently commercially available through Sage  
220 Science. The TELL-seq technology workflow is as follows (Fig.1D):

- 221 ● Genomic DNA (0.1-5ng), the barcoded TELL beads (3-10 million) and  
222 transposomes are mixed in a PCR tube.
- 223 ● The transposomes and DNA segments interact to form a strand transfer  
224 complexes (STCs), which is connected with the barcode sequence on the TELL  
225 bead surface.
- 226 ● The transposase is removed, the DNA fragment is cut into two parts in the STC,  
227 and the beads removed, leaving a DNA fragment, connected to a transposon, which  
228 is in turn connected to a barcode.
- 229 ● The barcoded DNA molecules are amplified with P5 and P7 adaptors before  
230 illumina sequencing.

231 The library preparation for TELL-Seq differs from haplotagging in mostly minor  
232 ways, with the primary distinction being in the form that the barcode takes: TELL-Seq  
233 uses a simpler 18-base barcode, rather than the 4x6 method of Haplotagging. This  
234 allows for a larger number of unique barcodes – and hence reduced likelihood of a  
235 collision - but lacks the error-correction feature.

## 236 **Data features and quality assessment**

237 Before we discuss data applications, we first introduce metrics on quality assessment  
238 and then use the metrics to evaluate datasets sequenced for this study. Our focus will  
239 be on Hi-C, 10x and Haplotagging, which are currently or previously available in the  
240 market.

---

<sup>1</sup> We note that the acronym TELL-Seq (Transposase Enzyme Linked Long-Read Sequencing) falls afoul of the terminology confusion referenced earlier. Under the terminology we have enforced, the linked reads produced are *long range*, but they are not *long reads*.

## 241 **Data Metrics**

242 In order to provide a robust analysis of the relative performance of the platforms, we  
243 must first derive numerical metrics by which to judge them. Since the Pair-Linked  
244 Platforms platforms differ significantly in the mode of operation from the Chain-  
245 Linked Reads, the metrics used will be slightly different, but our design aims to  
246 enable as valid a comparison as possible.

247

### 248 *Metric 1: Association*

249 Association is the ability for long range information to be communicated by the  
250 platform, or equivalently, the amount of non-local information contained within a  
251 read. Datasets with a higher association contain more and longer-range information  
252 than those with a lower association. In the context of trying to use long-range  
253 information as an assembly tool, a stronger association is preferable.

254 For the PLR platforms, association is measured by the distribution of Link-  
255 Separation Distance, the distance on the linear genome between the two ends of  
256 paired reads which have been linked together. If the first end of the pair aligns to a  
257 location  $i$ , and the second end to  $j$ , then the genomic distance is  $|i - j|$ . If large values  
258 of  $|i - j|$  are found to occur more often, then the dataset has a stronger association.  
259 Whilst we should therefore favour platforms which have a higher proportion of reads  
260 with large  $|i - j|$ , we note also that there is an expected pattern at higher distances: if  
261 the linkage probability is inversely proportional to some power ( $b$ ) of the physical  
262 distance between the reads, and at large linear distances genomic distance and  
263 physical distance are approximately the same, then we expect the frequency to fall as

$$264 \quad p(\text{link } i, j) \propto \frac{1}{|i - j|^b} \quad \rightarrow \quad f(|i - j|) \approx A|i - j|^{-b}$$

265 Where  $A$  is an arbitrary scaling parameter. On a log-log scale, this manifests as a  
266 linear relationship between the separation distance, and the observed frequency.

267 Deviations from this pattern indicate problems with the library preparation and can  
268 result in the failure of any statistical inference based on the dataset. We should  
269 therefore prefer datasets which i) exhibit a power-law relationship in frequency at  
270 high separation distances and ii) Have a smaller exponent, resulting in a longer tail,  
271 and hence more long-range information.

272 In the case of the Chain-Linked Read platforms, the long-range information is  
273 conveyed by labelling reads as originating from a larger molecule via a tag shared by  
274 all fragments of that molecule. The association should therefore be measured by the  
275 size of the molecules from which the labelled reads are drawn.

276 It is clear that having a larger molecule is generally better: each barcode  
277 delineates a larger spatial region, so the information is longer-range. There is,  
278 however, an upper limit at which point increasing the molecule size gives decreasing  
279 returns: for example, if the molecules were chromosome scale, then the barcoding  
280 would simply inform us which chromosome the read is from: useful, but not  
281 beneficial for assembling the reads within a given chromosome. Of critical concern,  
282 however, is that increasing the molecule size increases the chances of barcode  
283 collisions, behaviour demonstrated in Supplementary Figure 2. Generally, the size at  
284 which collision rates become untenable is significantly below the genome size, and  
285 hence should be treated as the limiting factor on the molecule size. We should  
286 therefore favour platforms which generate larger molecule lengths, but which still  
287 have a small collision rate.

288

## 289 *Metric 2: Accessibility*

290 Accessibility is the fraction of the data which is unique, unambiguous and useable.  
291 Datasets which have a low accessibility may still contain useful scientific data, but  
292 much more data would be required to achieve the same level of significance. We  
293 should therefore prefer platforms which produce highly accessible data. For example,

294 both CLR and PLR suffer from potential PCR duplication – the overamplification of  
295 some portions of the genome through the library preparation process. A high PCR  
296 duplication rate is indicative of a poor accessibility, and vice versa. Complex factors  
297 underlying the library preparation can also lead to reads which cannot be mapped to  
298 the reference genome (and the rate of unmapped reads is notably higher in Long  
299 Range platforms than normal Illumina short reads), or which contain no linking  
300 information (‘singletons’). Such ‘unmapped’ reads contain no useful information, and  
301 so they too should be excluded from further analysis.

302 In addition, PLR explicitly allows inter-chromosomal interactions to be mapped.  
303 Whilst this is useful in general in 3D genomics, for the purposes of the applications  
304 discussed in section 0 this represents unusable data, as assembly should occur on a  
305 per-chromosome basis. In order to have the maximum amount of usable information,  
306 we should therefore prefer the platforms which have a smaller number of linkages  
307 between chromosomes: a smaller translocation rate.

308 Assuming that other sources are negligible (or, equal between platforms), the  
309 total accessibility of the dataset can therefore be computed from the PCR duplication  
310 rate  $D$ , the translocation rate  $T$  and the unmapped rate,  $U$ :

311 
$$A = 1 - D - T - U$$

312 A higher value of  $A$  indicates a dataset which contains more useful information.

### 313 *Metric 3: Evenness*

314 Evenness is the measure of statistical validity in the coverage of the genome. A high  
315 coverage is evidently preferred, as it means that more of the genome was sampled and  
316 there is a smaller chance of missing portions of the genome, however, it is also  
317 important to ensure that the coverage was not biased onto some portions of the  
318 genome over others: there should be an equal likelihood of a read being generated  
319 anywhere on the genome. Datasets which deviate from this pattern are uneven, and

320 likely to be biased in complex and unpredictable ways. We should instead seek out  
321 datasets with a higher level of evenness.

322 Under the standard statistical assumptions, if the genome is sampled at a uniform  
323 rate everywhere, the coverage should follow a Poisson distribution,  $\mathcal{P}(k|\lambda)$ . However,  
324 it is easy to show that the coverage of any platform exhibits a significantly greater  
325 dispersion than a Poisson distribution with the correct mean. This is generally  
326 interpreted as being indicative that there is not just one rate,  $\lambda$ , at which the genome is  
327 sampled, instead there are multiple values, over which the distribution is marginalised  
328 [19].

329 In Supplementary 0, we use this information to generate the following  
330 unevenness metric:

$$331 \quad \mathcal{U} = \frac{\text{Var}(\text{coverage}) - \langle \text{coverage} \rangle}{\text{Var}(\text{coverage})}$$

332 Where  $\text{Var}(\text{coverage})$  and  $\langle \text{coverage} \rangle$  are the standard statistical variance and mean of  
333 the non-zero coverage<sup>2</sup> distributions respectively. This value is zero when the  
334 coverage distribution is a perfect Poisson distribution, and is arbitrarily large for  
335 distributions which have many values of  $\lambda$  contributing to them. We should therefore  
336 favour platforms which generate smaller values of  $\mathcal{U}$ .

337

#### 338 *Metric 4: Capability*

339 Capability is the measure of usefulness of the dataset, the ability for the dataset to  
340 improve the outcome of a genetic inquiry than would otherwise be achieved without  
341 long-range information. A more capable platform produces data which allows the  
342 assembly to be vastly improved, and should therefore be preferred.

343 We measure the capability by comparing the N50 and N90 metrics of a  
344 scaffolding with and without the assistance of long range information. The N50

---

<sup>2</sup> We focus on the non-zero coverage distribution since the designs of the Hi-C and 10x protocols mean higher zero-coverage is to be expected, but the non-zero coverage should be unaffected.



345 metric is the standard measure of ‘completeness’, it is the length of the shortest  
346 continuous sequence such that all longer sequences make up more than 50% of the  
347 genome. N90 is defined similarly, but encapsulating 90%. Larger values for  $N_x$  are  
348 preferred, as this indicates that more of the genome has been grouped into larger,  
349 contiguous fragments.

350

## 351 **Pair-Linked Reads**

352 Of the three companies that have commercialized Hi-C; Cantata (formerly Dovetail),  
353 Arima, and Phase, the most widely applied technologies are OmniC (Cantata) and  
354 Arima. In this study, we only carry out analysis on Arima Hi-C reads and  
355 comparisons are performed between the two generations of Arima technology (V1  
356 and V2), to reveal characteristics and library improvement by the platform. In total,  
357 we obtained three human datasets, two from V2 (NA24385-AJ and NA12878-CEU)  
358 and one from V1 (NA12878-CEU; see Table 1).

359 Hi-C maps are shown in Figure 2 for three human samples by mapping the reads  
360 to the human reference assembly GRCh38. In these plots, regions of high density  
361 indicate real-space collocation of the genome – though there are some notable  
362 deviations from this; in particular, highly repetitive regions can cause spurious over-  
363 and-under densities, characterised by a cross-shape running through, i.e. the  
364 centromere of each chromosome. To explore the quality of these datasets in more  
365 detail, we present Figure 3 which consists of three separate plots showing link-  
366 separation distance (association), translocation rates (accessibility) and base coverage  
367 (evenness) respectively. Tabulated information regarding the metrics is also presented  
368 for accessibility (Table 1) and evenness (Table 3).

369 Fig. 3A shows how the long-range information is distributed in the Hi-C  
370 datasets: as expected we see a peak of very strongly associated regions in the region  
371 of 100-500bp (probably due to topologically-associated domains, TADs), and a long  
372 power law tail for the three human datasets. For demonstration purposes, we also  
373 include an additional dataset – derived from Oak – which demonstrates a strong  
374 deviation from the power law structure. In assessing the association demonstrated  
375 here, we would say that the oak should be penalised due to this deviation whilst the

376 human datasets are comparatively much nicer. Aside from this, the plots demonstrate  
377 that V2 datasets have more information stored in longer length reads than the V1, and  
378 hence have a stronger association.

379 The Usability metric is shown graphically in Fig. 3b and in more detail in Table  
380 1. We find that the V1 dataset shows a consistently poorer mapping rate, PCR  
381 duplication rate and translocation rate over the V2 datasets, resulting in a usability of  
382 0.328, compared to 0.53 and 0.60 for the V2 data, though we do note from Fig. 3B  
383 that the difference between human datasets was, on some chromosomes, more  
384 pronounced than the difference between platforms.

385 Fig. 3C and Table 3 show the evenness statistics. Fig. 3B shows the raw  
386 coverage data for the Hi-C data, along with a standard Illumina sequencing of the  
387 NA12878-CEU sample for comparison: given that the Illumina data has been  
388 sequenced more directly, with fewer intervening biochemical alterations, we should  
389 expect it to be the “purest” sample. Visually, we can see that this is the case: the  
390 Illumina is tightly peaked and resembles a Poisson distribution. The V2 datasets –  
391 though sampling to slightly different depths – show a similar “fattened Poisson”  
392 distribution, and the V1 data seems to be the least pure sampling, showing a strong  
393 overdensity at low base coverage. These observations are carried through by the  
394 statistical metric developed in 0: the Illumina data was given a score of 2.7, whilst the  
395 V2 platforms both scored approximately 5, and the V1 platform scored 10, indicating  
396 a strongly uneven coverage. This would indicate that whilst there is some statistical  
397 bias in the V2 data, it is significantly less than that of the V1.

398 From the information presented here, we would conclude that the V2 platform  
399 produces data which robustly outperforms V1, with the two V2 datasets very close  
400 together in quality: V2 NA24385-AJ has a slightly higher mean base coverage, but V2  
401 NA12878-CEU scores slightly better on the global accessibility and evenness metric.  
402 In the next section, we will demonstrate how Hi-C data can be used to aid Genome  
403 Scaffolding, and hence assess the usability of these datasets.

#### 404 **Chain-Linked Reads**

405 Due to their similarity in mechanism and data output, we discuss the 10x and  
406 Haplotagging qualities together, presenting an analysis on five datasets, two from  
407 from 10x (human and hummingbird) and 3 from Haplotagging (human, rat and oak).

408 The human 10x dataset was downloaded from the 10x Genomics website and  
409 hummingbird dataset is part of the VGP project (see Data availability). The  
410 Haplotagging datasets of human, rat and oak were sequenced by the Sanger Institute  
411 as part of the Darwin Tree of Life project. We note that, since the data arises from  
412 wildly different species, we must take care with our inferences that the differences  
413 arise from the choice of platform rather than the choice of species: any strong  
414 comparisons should be based primarily on the human samples.

415 In Figure 4 we see the distribution of molecule lengths for the CLR platforms,  
416 which we use as a measure of the association, and in **Error! Reference source not**  
417 **found.** computes the associated barcode collision frequencies for these molecule  
418 distributions. We find that all of the platforms produce molecules which have  
419 collision rates below 0.1%, and have mean molecule lengths in the region of 50kb.  
420 We note that the 10x platforms have a more prominent tail at the high-length end of  
421 the distribution, most evident through the N50 values: the 10x N50 values exceed the  
422 mean length by 40kb, whilst the haplotagging N50 exceed the mean by only 20kb,  
423 indicating the 10x has a stronger tail of high-association data included, even if the  
424 bulk of the data has similar associations.

425 Table 2 shows the PCR duplication rates and the unmapped rates (and hence the  
426 accessibility) as well as N50 reads per barcode. We can see clearly that 10x has lower  
427 PCR duplications than haplotagging, although the inverse is true for the unmapped  
428 rate: this could be largely due to the differing tools used to analyze the datasets (EMA  
429 for haplotagging versus LongRanger for 10x), that the 10x data is several years old,  
430 whilst the haplotagging is state-of-the-art, and we note that the total accessibility is  
431 broadly the same. Probably of more importance is that 10x data exhibits a higher  
432 number of N50 reads per barcode than haplotagging.

433 In Figure 5, we demonstrate the coverage profiles of the CLR datasets, in their  
434 raw form in 5a, and, to remove the effects of differing sequencing depths, a  
435 normalised form in 5b, where the coverage is given as a fraction of the maximal  
436 value. In these figures we can see a number of clear features: firstly, it is clear that  
437 both the oak and the rat show extremely strong deviations from an even profile: we  
438 hypothesise that this may be due to the effects of highly repetitive regions (that rats  
439 are known to possess [20]), which causes some regions of the genome to be  
440 erroneously ‘covered’ thousands of times, whilst other regions are deprived of  
441 coverage. The human profiles are similar in shape to the Illumina curve. However .

442 we note that the high coverage end (as with the rat, likely a spurious tail due to over-  
443 coverage of repetitive regions) is suppressed relative to the Illumina sample, an  
444 indication of the long-range information allowing correct alignment of some repetitive  
445 regions. Nevertheless, we do see a stronger bias towards the low coverage end in both  
446 the 10x and haplotagging than in Illumina. In contrast, the hummingbird displays a  
447 remarkably Poisson-like shape, in large part due to the almost total absence of a  
448 repetitive high-coverage tail – likely a feature of a small, non-repetitive genome [21].

449 These visual conclusions are supported by the unevenness metric in Table 3,  
450 where we see both the rat and the oak scoring very poorly (11.5 and 43.6  
451 respectively), the humans scoring between 2 and 9 (illumina: 2.7, haplotagging: 5.1  
452 and 10x: 8.9), and the hummingbird with the lowest score of 1.4. We note that the  
453 haplotagging's improved score over the 10x platform is likely a feature of the more  
454 powerful suppression of the over-coverage of repetitive regions, rather than of an  
455 improvement at the low-coverage end: this might indicate that the haplotagging is  
456 more successful in allowing alignment of repetitive regions than 10x.

457 From the metrics presented here, we conclude that, for base polishing, 10x data is  
458 superior to that of Haplotagging due to its slightly higher association strength.  
459 However, Haplotagging has a larger number of unique barcodes, resulting in a much  
460 lower collision rate, and this means that it is more efficient to handle large number of  
461 samples when low coverage data is targeted. In addition, we recall that the statistical  
462 properties of the haplotagging platform indicate that haplotagging allows better  
463 alignment of highly repetitive regions than 10x.

#### 464 *TELL-Seq*

465 Like haplotagging, TELL-Seq is a promising successor to 10x technology, and so we  
466 wish we could have a similar analysis of the platform against the metrics we have  
467 formulated here. However, the authors were unable to produce a TELL-Seq library of  
468 sufficient quality to provide a viable comparison. We must therefore rely on the  
469 literature (i.e. [12]) when discussing the properties of TELL-Seq.

## 470 Applications

471 In this section we briefly outline some of the main applications for long-range data,  
472 and discuss how this has been applied in the literature.

### 473 **Genome scaffolding**

474 Genome scaffolding is the process by which a number of continuous sequences  
475 ('contigs') generated from overlapping reads are linked together into a single structure  
476 (a scaffold) of known sequences, separated by gaps of unknown sequences but where  
477 the length of the gap is relatively well constrained. This forms a critical step in  
478 genome assembly [22], but conventional means are both laborious and  
479 computationally intensive, though recent advances in long-range sequencing  
480 technologies have improved the continuity of genome scaffolds [23], for example, the  
481 assembly quality thresholds proposed by Vertebrate Genome Project (VGP) are that  
482 contig N50 > 10Mb and the scaffold N50 is the chromosome length [6], indicating  
483 that chromosome-scale scaffolding is now routinely possible.

### 484 *Pair-Linked Reads*

485 The Hi-C protocol provides a fast and lower-cost way of constructing scaffolding  
486 from the contigs, given that the spatial information within Hi-C Pair-Linked Reads  
487 can identify whether contigs come from the same chromosome and infer the correct  
488 orders of the contigs within each chromosome based on the relative proximity  
489 between bases in each contig [14]. This technology is widely used to assemble the  
490 contigs of eukaryotic genomes into chromosome-scale scaffolds [6, 22], and has  
491 recently been applied to assemble the giant and complex genome of Chinese Pine into  
492 a chromosome-level assembly [24]. To further improve the quality of genomic  
493 assembly, some studies evaluated the different sample preparation kits/protocols and  
494 computational programs and identified the optimal conditions for Hi-C scaffolding  
495 [25].

496 To demonstrate how Hi-C data can improve the quality of scaffolding, we  
497 applied the above methods to the human genome. 54X HiFi reads from HG002 were  
498 downloaded from GIBA and an assembly was obtained using Hifiasm [26] with

499 contig N50 at 45.1 Mb – this represents the baseline shown in Figure 6A, which is a  
500 standard Hi-C proximity map. However, without scaffolding, the proximity map  
501 shows a high degree of fragmentation. After removing haplotype duplications, the  
502 contigs were further assembled to chromosome-level scaffolds using ~30X Hi-C reads  
503 and the YaHS scaffolding tool [27]. The resulting maps for Arima V1 (Figure 6B) and  
504 Arima V2 (Fig. 6C) are shown: after scaffolding, chromosome blocks are clearly seen  
505 and the fragmentation visibly reduced. Table 4 details how the lengths of the  
506 assembled fragments vary after applying the Hi-C data: we find that although V1  
507 produces slightly higher N50 scaffolds and a larger maximum length scaffold, the V2  
508 platform has a higher mean length, indicating a significant reduction in the number of  
509 poorly-scaffolded contigs, consistent with our earlier analysis of these platforms.

510 Detailed instructions on genome assemblies are provided in Supplementary.  
511 Assembly pipelines/instructions/recommendations can also be found in VGP, see  
512 Rhie, et al., [6].

### 513 *Chain-Linked Reads*

514 The Chain-Linked strategy provides an effective and less expensive alternative  
515 technology than NGS mate pair data for genome scaffolding [22], since the nature of  
516 these platforms groups reads by their proximity in the genome (absent any barcode  
517 collisions). The information retained from long stretch sequences can be used to link  
518 the different contigs to chromosome-level scaffolds, a strategy which is already  
519 widely used for genome scaffolding in a variety of complex and polyploid species.  
520 For example, in an early study, 10x reads were applied to assist in scaffolding of  
521 genome sequence of *Triticum urartu*, the progenitor of the A subgenome of tetraploid  
522 and hexaploid wheat [28], and Lee et al. [29] reported that 10x linked reads were  
523 successfully used to correct and scaffold the assembly for an allopolyploid rapeseed.

524 Currently, 10x-based approaches can no longer be used due to the withdrawal of  
525 the product, but alternative linked-read technologies have developed based on similar  
526 methods, such as Haplotagging and TELL-Seq [12, 13]. As we have seen, the  
527 association strength between these platforms is similar, and hence these platforms all  
528 offer the ability to produce high quality scaffolding, an assumption which was  
529 validated when Chen et al. [12] completed a comprehensive assessment of TELL-Seq  
530 using the sequenced data from sample NA12878 and NA24385. The TELL-Seq data

531 of NA12878 produced a *de novo* assembly with a scaffold N50 of 31.5 Mb, the  
532 longest contig 109.2Mb and the longest alignment of 23.6 Mb.

533

## 534 **De novo genome assembly**

535 *De novo* genome assembly is the fundamental process in reconstructing a genome  
536 from sequencing reads without a reference sequence [30]. A whole-genome assembly  
537 with high level of completeness, continuity and accuracy is the key, which can  
538 significantly enhance the reliability of the downstream analyses. In general, the  
539 primary step for *de novo* assembly from the collection reads consists of three phases,  
540 contig assembly, scaffolding and gap filling. As we saw in 0, both Pair- and Chain-  
541 Linked Reads can connect and order contigs into a ‘scaffold’ in the second phase,  
542 however, Chain-Linked Reads also offer additional support for the other procedures  
543 involved in *de novo* assembly, due to their inherently more structured nature.

544 10x technology has been used extensively for the *de novo* assembly of the  
545 eukaryotic and prokaryotic genomes [3]. For instance, the first complete genome  
546 sequence for the mound-building mouse, *Mus spicilegus*, was generated with 10x  
547 reads and resulted in the *de novo* assembly of a 2.50 Gbp genome with a scaffold N50  
548 of 2.27 Mbp [31]. Using 10x data and the Supernova assembler, Ozerov et al. [32]  
549 assembled a ~0.8Gb draft genome of the *Silurus glanis*, an important species for  
550 freshwater ecosystem balance. It has also been demonstrated that 10x data can be used  
551 to assemble high accuracy contigs and scaffolds, even for large, highly similar  
552 repetitive sequences, polyploid plant genomes [3, 33].

553 As comparatively newer technologies, the non-10x CLR platforms have seen less  
554 ubiquitous use in *de novo* assembly, though Chen et al. [12] did introduce the TELL-  
555 Seq platform by immediately providing *de novo* assembly of bacteria (*Escherichia*  
556 *coli*, *Campylobacter jejuni*, *Rhodobacter sphaeroides*) and humans (NA12878). The  
557 human assembly showed “longer aligned contig length and at least 28% and 71%  
558 fewer misassemblies than other linked-read or nanopore methods, respectively” [12].  
559 In addition, some attempts have been made to use haplotagging for *de novo* assembly,  
560 though the success has been more limited [34].

561 Although *de novo* genome assembly can be performed by using CLR  
562 technologies alone, most current studies adopt a hybrid strategy of multiple

563 technologies to complete genome assembly. For example, Batra et al. [35] performed  
564 a *de novo* genome assembly of the olive baboon using a hybrid sequencing approach  
565 of 10x sequencing, Oxford Nanopore sequencing, Illumina paired-end sequencing and  
566 Hi-C, which have complementary advantages. Lind et al. [36] generated a high-  
567 resolution *de novo* chromosome-scale genome assembly for the Komodo dragon  
568 *Varanus komodoensis* using data from different platforms, including 10x Genomics  
569 linked-reads, Oxford Nanopore long reads, PacBio long reads and Bionano optical  
570 mapping.

## 571 **Variation detections**

572 One of the most fundamental goals in genetics is to link genomic variations and the  
573 evolution of traits between populations or species. DNA polymorphisms are  
574 widespread genomic variations among individuals and include single-nucleotide  
575 variants (SNVs), small insertions and deletions (Indels; <50 bp), and structural  
576 variations (SVs). Many methods have been proposed to test DNA changes across the  
577 genome from different sequencing technologies, but there are still considerable  
578 limitations on what can be achieved in SV detection due to technical difficulties of the  
579 standard short-read platform. The long-range information provided by CLR and PLR  
580 platforms can improve detection for haplotype-specific deletion and large SV [37,  
581 38].

### 582 *Pair-Linked Reads*

583 Since Hi-C technology detects regions of high interaction probability in a genome,  
584 this intrinsically makes it particularly useful for detecting SVs. One of the main  
585 advantages of Hi-C is that it can accurately detect SVs with low-depth sequencing  
586 data. This feature provides a higher chance of identifying SVs at repetitive regions in  
587 complex genomes.

588 As a result, Hi-C has been demonstrated to be a promising technology to  
589 precisely detect SVs, including chromosomal rearrangements and copy number  
590 variation in plant and human genomes [39, 40]. In recent years, several research  
591 projects have shown the ability of Hi-C to support identifying three-dimensional  
592 genome organization alterations as a result of SVs in the human cancer genome [41,  
593 42, 43]. Hi-C has also been applied to screen the complex genomic rearrangements



594 associated with the development of disease in humans. For example, Melo et al. [44]  
595 used Hi-C to investigate the genetic variation that causes developmental disorders,  
596 and Hi-C was used to detect multi-megabase polymorphic inversions in wheat and  
597 barley [45, 46].

### 598 *Chain-Linked Reads*

599 Recent work has used SNVs detected by 10x sequencing technology to draw the  
600 landscape of meiotic recombination in plant population at the genome-scale resolution  
601 [47, 48], and Rommel Fuentes et al. [48] pinpointed meiotic crossovers of  
602 interspecific hybrid F1 tomato pollen at the SNV resolution level by using 10x data.  
603 This technique also has been a powerful tool for detecting genomic variants  
604 associated with human diseases. A number of novel and important SVs associated  
605 with metastatic castration-resistant prostate cancer were identified by 10x whole-  
606 genome sequencing [49], and CLR sequencing validated the inverted rearrangement  
607 in the triple-negative breast cancer sample TN-19 [50]. A 2020 study confirmed that  
608 10x sequencing provides a cost-efficiency way of mining genomic variants at  
609 moderate depth and population scale [51], and it was also reported that 10x  
610 technology could be used to screen nucleotide resolution of the structural variants  
611 linked with potential risk loci in small and rare disease cohorts [52].

612 Haplotagging is particularly suitable for constructing the original haplotype, and  
613 as a result has been successfully applied to construct the genome haplotypes in the  
614 two butterfly species, and detect the genetic markers controlling the distinct wing  
615 color patterns [13], indicating that haplotagging might be a promising method to  
616 identify the superior haplotype alleles in the diverse plant or animal populations for  
617 model and non-model species. Bhat et al. [18] thought that this technique would  
618 provide important support for haplotype-based breeding for crop improvement.

619 The utility of the TELL-seq protocol, for detecting genome variations has not  
620 been nearly so widely used in the plants or animals, though the study by Chen et al.  
621 [12] demonstrated that linked-read data generated by TELL-seq could be used to  
622 screen genetic variation using an analysis pipeline developed for the 10x technology..  
623 Although this means TELL-seq also could be used to detect SVs, the initial study  
624 found that it missed some deletions in the NA12878 sample. The authors thought that  
625 two factors (the short library insert length and different barcoding chemistry) might be

626 responsible, and they encourage the research community to further develop and  
627 optimize analytical tools to improve the ability to detect SV using linked-read data  
628 [12]. More extensive validation studies are therefore needed to prove whether TELL-  
629 seq can accurately detect genome-wide variation as an alternate method for the 10x  
630 platform.

## 631 **Other Applications**

### 632 *Phasing*

633 Phasing – the assignment of alleles to either the maternal or paternal haplotype – is  
634 another potential application for long-range reads, since even long reads can struggle  
635 to accurately identify heterozygosity and correctly assign differences to haplotypes.  
636 Along with *de novo* assembly, Chen *et al.* demonstrated how TELL-seq can be used  
637 as a powerful tool for phasing the genome :TELL-Seq phasing results on NA12878  
638 and NA24385 samples showed that the highest heterozygous rate is 99.9% and  
639 99.8%, the phasing block N50 is 16.1Mb and 13.4Mb, the longest phasing block is  
640 67.5Mb and 59.2Mb, and adjusted N50 (1.24 Mb), and the lowest switch error rate is  
641 0.04% and 0.08, respectively [12].

642 Most recently, a study compared the performance and accuracy of genome  
643 phasing between Hi-C and 10x Genomics Linked Read in Hanwoo Cattle [53]. The  
644 results of this study showed that the phasing strategy with 10x linked-read technology  
645 and Long Ranger software displayed the best phasing performance. The best strategy  
646 had the highest phasing rate (89.6%), longest adjusted N50 (1.24 Mb), and lowest  
647 switch error rate (0.07%). Moreover, the phasing accuracy and yield of the best  
648 strategy stayed over 90% for distances up to 4 Mb and 550 Kb, respectively.

### 649 *Metagenomics*

650 Another application of Chain-Linked Read sequencing technology is assembling  
651 high-quality metagenome of microbial species, which is able to improve continuity  
652 and accuracy in *de novo* assembly using barcode information, as comprehensively  
653 evaluated by Zhang *et al* [54]. This study showed that 10x reads significantly  
654 improved the metagenome assemblies when compared with Illumina short-reads,

655 although both were outperformed by PacBio CCS long-reads. Due to the low cost and  
656 the high base quality, sequencing the metagenomes using Chain-Linked Read  
657 technology remains persuasive. Recently, Roodgar et al., [55] explored the  
658 longitudinal trajectories of gut microbiome for a single individual using linked-read  
659 metagenomic sequencing in 10x Genomics Chromium platform.  
660

## 661 **Selection of data platforms**

662 With rapid development of long read technologies for longer read length and better  
663 base accuracy, high profile projects have been launched such as Vertebrate Genomes  
664 Project (VGP) which aims to sequence all the vertebrate species [6] and Darwin Tree  
665 of Life (DTOL) which plans to generate de novo assemblies for the 70,000 eukaryotic  
666 genomes in Britain and Ireland ( <https://www.darwintreeoflife.org> ). More  
667 ambitiously the Earth BioGenome Project was proposed to decode ~1.5 million  
668 eukaryotic species, including animals, plants and microbiomes [7]. If targeting  
669 chromosome-level assemblies, Hi-C data sequencing should be planned, either with  
670 Arima V2 or OminC. When sequencing Hi-C in large volumes of data with various  
671 species, the assessment metrics and methods presented in the study could be used for  
672 data QC. In the cases where there are choices of platforms, data assessment and  
673 comparisons are essential in order to ensure proper Hi-C libraries are prepared. For  
674 small research groups, contracting Hi-C sequencing is one of the options while most  
675 sequencing companies provide Hi-C QC reports.

676  
677 Chain-Linked reads, such as 10x, Haplotagging and TELL-seq can be used for  
678 consensus polishing to improve the quality of genome assembly and enhance the  
679 detection of genomic structure variants. We do note that, as of June 30, 2020, 10x  
680 Genomics discontinued the sale of Chromium Genome and Exome product lines – the  
681 most prominent CLR platform, on which a significant portion of the literature was  
682 focused. Various alternatives have been suggested: in this work we studied Sage  
683 Science’s TELL-Seq platform and Haplotagging. Whilst we found that haplotagging  
684 data was in some cases of a higher quality than 10x, haplotagging beads are not (yet)

685 commercially available, being obtainable from the Chan Lab at the Max Planck lab in  
686 Tuebingen only via academic collaboration. Commercial supply of these reagents  
687 could make haplotagging a powerful tool, as the beads are potentially inexpensive,  
688 which would allow haplotagging to be used widely in genetic population sequencing  
689 studies. Additionally, more work on analyzing and processing non-10x data would  
690 futher enable the community to make use of these potentially powerful platforms.

## 691 Software tools

692 To date, a large number of tools have been developed to analyze data generated from  
693 long-range sequencing technologies [54, 56, 57]. Here, we highlight recent  
694 developments in software tools used for genome scaffolding, de novo assembly and  
695 variation detection based on the long-range linking information.

### 696 Hi-C Analysis Tools

#### 697 *Genome Scaffolding with Hi-C*

698 Several scaffolding methods have been developed for assembling contigs to scaffolds  
699 based on Hi-C data, examples of which are shown in Table 5. There are many  
700 different approaches which can be taken in designing these tools, which we broadly  
701 split into three categories: *deterministic*, *probabilistic*, and *improver*.

702 Deterministic tools use algorithms which always return the single result which  
703 optimises some underlying metric. Some examples of deterministic algorithms  
704 include:

- 705 • Heirarchical clustering

706 Typically using an agglomerative approach, such as in the early tools  
707 LACHESIS [58] and dnaTri [59] (both no longer actively developed), and  
708 more recently by ALLHIC [60], a framework particularly designed for  
709 scaffolding autoploid or heterozygous diploid genomes.

- 710 • Best-Neighbour

711 Though deterministic, best-neighbour methods return only an approximation  
712 to the desired answer, at the benefit of vastly increased speed. 3D-DNA [61]  
713 used this approach after correcting the input contigs. The best-neighbour

714 approach then assembles the contigs into one megascaffold, before it is then  
715 cuts to a number of chromosomes on the basis of Hi-C contact matrix.

716 • Maximal-Matching

717 This approach is used in the SALSA1 [62] tool, which first corrects  
718 misassemblies derived from the input contig using a low Hi-C mapping rate  
719 as the signal for error and then orients and orders the corrected contigs to  
720 generate scaffolds using a maximal matching algorithm.

721 • Novel approaches

722 Some more novel solutions include SALSA2 [56], an overhaul of the  
723 SALSA1 program that can take advantage of all the interaction information  
724 from the Hi-C map to reduce assembly errors using a novel iterative  
725 scaffolding method, as well as the newly developed YaHS [27], which  
726 introduced a novel algorithm to establishing the contact matrix to obtain the  
727 more accurate inferences of contig joins.

728 Probabilistic approaches, in contrast, return results which are not exact, but are good  
729 approximations to a desired solution where direct computation would be prohibitive.

730 We identify two main classes of probabilistic algorithm.

731 • Markov Chain Monte Carlo

732 MCMC methods are a class of algorithms which attempt to efficiently  
733 approximate drawing values from an underlying (unknown) distribution  
734 function. This is used by GRAAL [63] which uses a MCMC algorithm to  
735 generate scaffolds from the Hi-C data. Recently, Baudry et al. [64] developed  
736 instaGRAAL, an upgrade of the GRAAL version, which can be used to  
737 assemble large genomes.

738 • Maximum Likelihood

739 Maximum likelihood methods use Bayesian formulations to derive a  
740 probability of observing given results, given a hypothesized original state. By  
741 optimizing this function, the original state can be inferred. This approach is  
742 used by HiRise, the tool developed by Dovetail Genomics for their Hi-C  
743 service [65].

744 Finally, we note a class of tools which we dub *improvers*, these tools do not perform  
745 the assembly themselves, but act to improve the quality of assemblies performed  
746 using other tools. Examples include HIC-Hiker [66], a probabilistic and dynamic

747 programming approach which can improve the quality of scaffolds produced by other  
748 Hi-C scaffolding software, and the recently developed EndHic [67], which can reduce  
749 the error rate of assembly using only the the Hi-C contacts from the end regions of the  
750 contigs.

751 Several studies have evaluated the performance of different scaffolders for  
752 scaffolding accuracy [60, 67, 68, 27]. For example, a recent study evaluated the  
753 performance of five Hi-C scaffolders including LACHESIS, HiRise, 3D-DNA,  
754 SALSA2, and ALLHiC; the results found that the HiRise and LACHESIS display the  
755 best performance on average under all tested scripts [68]. However, with all the  
756 available software, it remains challenging to correctly assemble large contigs into  
757 chromosomes, and manual checking and curation are often necessary. The selection  
758 of suitable tools therefore often remains an exercise in trial-and-error by the  
759 researcher.

#### 760 *Variation Detection with Hi-C*

761 There exist several computational tools which have been developed to identify SVs  
762 from chromatin interaction data. We divide these by the kinds of SV which they can  
763 identify.

764 Tools which can identify Copy Number Variations (CNVs) include HiCNV [41],  
765 OneD [69] and HiNT-CNV [70]. Generally speaking, these tools use Bayesian  
766 information criteria (and in the case of HiCNV and OneD, Hidden Markov Models,  
767 HMM) to identify the location of CNVs. Similar methods can be used to identify  
768 interchromosomal translocations – the tools HiCTrans and HiNT-TL are packaged  
769 alongside HiCNV and HiNT-CNV respectively [41, 70]. Although the above  
770 algorithms were used to screen the SVs within Hi-C data, most of these methods can  
771 only detect interchromosomal translocations and long-range intrachromosomal SVs at  
772 a low resolution.

773 Some more specific tools include the HiTea [71] software, developed specifically  
774 for identifying mobile transposable element insertions in Hi-C data, as well as  
775 NeoLoopFinder [72]; developed for predicting SV-induced chromatin loops, though  
776 also capable of detecting complex SVs with Hi-C data. Wang and colleagues [40]  
777 have also presented a computational framework, EagleC, which integrates deep-

778 learning and ensemble-learning strategies to detect a full range of SVs at high  
779 resolution.

780 Overall, there are still strong demands for analysis tools that can use Hi-C data  
781 for high-resolution SV detections.

## 782 **Chain-Linked Read Analysis Tools**

### 783 *Genome Scaffolding with CLR*

784 Generally speaking, most CLR tools should be equally effective, regardless of which  
785 CLR platform was used – however due to its prominence, many tools were designed  
786 specifically for 10x, and so their applicability to another linked-read platform,  
787 including TELL-seq and haplotagging, still need to be further verified. Unlike the Hi-  
788 C tools where a wide variety of differing algorithms were used for scaffolding, CLR  
789 algorithms broadly follow the same approach: first attempting to unambiguously  
790 identify the HMW-DNA fragments each read originated from, before using these  
791 fragments as the basis for a scaffolding. The tool fragScaff was first developed for  
792 scaffolding the data from contiguity preserving transposase sequencing, but was one  
793 of the first tools to receive explicit support for 10x reads [73]. fragScaff uses an  
794 explicit threshold metric to determine barcode uniqueness, before constructing and  
795 traversing a scaffold graph. ARCS and ARKS are two closely related tools developed  
796 by the same team [74, 75]: ARCS is a stand-alone genome scaffolding developed  
797 specifically for 10x linked reads, whilst ARKS uses a kmer mapping strategy to align  
798 linked reads and contigs to improve computational efficiency, and is an optional  
799 additional mode for ARCS. Hiltunen et al. [76] presented a software package  
800 ARBitR, which is explicitly designed to work on multiple platforms beyond 10x. The  
801 main distinctive feature of the ARBitR is that it consider the overlaps between the  
802 involved contigs when splicing, so as to improve the genome scaffolding accuracy.

803 Other CLR tools include SLR-superscaffolder [77], which uses an inverted top-  
804 down approach, and Architect, which uses co-barcoding and paired-end information  
805 to improve the contiguity of genome scaffolding [78].

### 806 *De novo Assembly with CLR*

807 Although there is much mature software that can be applied to de novo assembly of  
808 genomes with short-read sequence data [30, 79, 80], only a few comparatively fewer  
809 tools have been developed for generating a *de novo* genome from CLR data.

810 Supernova [16] was developed specifically for *de novo* assembly of genomes that  
811 were deeply sequenced using 10x linked-read sequencing platform, by 10x Genomics.  
812 Compared to other methods, Supernova can generate phased diploid assemblies over  
813 very long distances. Moreover, despite being a 10x product, Supernova can also be  
814 used for the data generated on other CLR platforms, such as TELL-seq [12].

815 Other assembly tools often use a de Bruijn-type approach, for example,  
816 cloudSPAdes [81] (an extensible module of the SPAdes assembler) uses CLR<sup>3</sup> data  
817 to expand the de Bruijn graph, and can also be applied to metagenomic or hybrid  
818 assembly. The Ariadne [82] module uses a novel algorithm, based on de Bruijn  
819 Graphs, to handle the barcode deconvolution problem. In their introduction of the  
820 TELL-Seq platform, Chen et al. [12] presented TuringAssembler, another de Bruijn  
821 graph-based assembler.

822 Whilst not strictly related to *de novo* assembly, we also note that Bishara et al.  
823 presented an assembler, Athena, that use the tag information from linked-read  
824 sequencing to improve metagenome assembly [83].

825

## 826 *Variation Detection with CLR*

827 A number of tools developed to detect genetic variations in NGS data can also be  
828 used on CLR data without significant modification, such as GATK [84], SNVer [85],  
829 VarScan [86] and VarDict [87]. However, since these tools do not exploit the long-  
830 range information, genome-scale SV detection remains limited, tools which are aware  
831 of the long-range information promise much greater detection power.

832 Long Ranger [38] is the official program developed by 10x Genomics, which can  
833 screen variants and SVs, and combines a number of existing tools, such as BWA and  
834 GATK, augmented with long-range specific algorithms. GROC-SVs [37] adopt a  
835 similar strategy to Long Ranger for identifying SVs, but it performs local assembly on  
836 barcoded reads to test high-resolution complex SVs. Recently, a new structural

---

<sup>3</sup> These platforms use the terminology “Synthetic Long Reads” (SLR), which we have attempted to move away from, preferring instead to refer to them as Chain-Linked Reads.



837 variant calling software was presented, called LEVIATHAN [88], which can detect  
838 SVs in highly fragmented and heterozygosity genomes using similar methods.

839 A “split molecule” approach has also proven successful, by identifying  
840 molecules which are Chain-Linked together, but aligned to disjoint parts of the  
841 genome. VALOR [89] has been developed to discover large genomic inversions from  
842 linked-read data by an algorithm based on this “split molecule” signature and read  
843 pair signature, and an improved version, VALOR2 [90], can identify not only  
844 inversions but also other complex SVs involved in segmental duplications,  
845 translocations and deletions. LinkedSV [91] also uses split-molecule methods to  
846 simultaneously integrate barcode overlapping and enriched fragment endpoints to  
847 identify large SVs.

848 NAIBR [92] identifies SVs by combining the split-molecule approach with a  
849 probabilistic model, and similarly, Xia et al. [93] developed the ZoomX tool using  
850 probabilistic models SVs signals would be represented in CLR, meaning ZoomX can  
851 detect novel genomic junctions, and hence identify large rearrangements (>200kb).

## 852 Conclusions

853 Here, we discussed the methodologies and applications of long-range, non-local  
854 sequencing technologies, focussing on the Pair-Linked Read technology of Hi-C  
855 through the Arima V1 and V2 platforms, and the Chain-Linked Read platforms of  
856 10x, Haplotagging and TELL-Seq. Assessing the published literature, we found that  
857 Hi-C has been widely used in genome scaffolding to assemble the genome on a  
858 chromosomal level, using a wide variety of different algorithmic approaches. Hi-C  
859 technology has also been used for assembly curation as well as evaluation and recent  
860 efforts have been seen on structural variation detections. The various Chain-Linked  
861 Read platforms have been demonstrated to enhance the value of short reads for  
862 genome assembly and, in contrast to the PLR platforms, widely used for improved  
863 structural variation detection.

864 We also introduced metrics with which to assess the quality of the sequencing  
865 data produced by these platforms, and briefly demonstrated that these metrics  
866 provided a robust insight into the ability of the platforms to provide useful genomic

867 information to researchers finding, for example, that the Arima V2 platform produces  
868 significantly higher quality data than the V1 platform.

869 From our analysis of the existing literature and from our quality metrics, we have  
870 found that long-range protocols, including Hi-C and Chain-Linked Read methods,  
871 have already been demonstrated to significantly improve the quality of genome  
872 assembly and enhance the detection of genomic structure variants, and as NGS  
873 technologies and the associated software pipelines continue to develop further, these  
874 technologies will continue to move from strength to strength.

875 We have emphasised throughout this work the distinction between true long-read  
876 platforms and the long-range technologies which employ genome partitioning and  
877 barcoding to cluster reads into groups providing with much needed long-range  
878 information with only a modest cost increase over standard short-read sequencing.  
879 Whilst the development of Long-Read technologies would initially seem to make the  
880 short-read based technologies discussed here less attractive to researchers, we have  
881 demonstrated robustly that non-local information can help supplement Long-Read  
882 endeavours, and avoid some of the drawbacks of these emerging technologies, such  
883 that a combined long-read/long-range approach remains a cost-effective strategy for  
884 complex genome and pan-genome assembly, population genetics, and high-resolution  
885 analysis of complex traits.

## 886 References

887

- [1] R. Sethi, J. Becker, J. d. Graaf, M. Löwer, M. Suchan, U. Sahin and D. Weber, "Integrative analysis of structural variations using short-reads and linked-reads yields highly specific and sensitive predictions," *PLoS computational biology*, vol. 16, no. 11, p. e1008397, 2020.
- [2] S. Goodwin, J. D. McPherson and W. R. McCombie, "Coming of age: ten years of next-generation sequencing technologies," *Nature Reviews Genetics*, vol. 17, no. 6, pp. 333--351, 2016.
- [3] A. Ott, J. C. Schnable, C.-T. Yeh, L. Wu, C. Liu, H.-C. Hu, C. L. Dalgard, S. Sarkar and P. S. Schnable, "Linked read technology for assembling large complex and polyploid genomes," *BMC genomics*, vol. 19, no. 1, pp. 1--15, 2018.
- [4] G. A. Logsdon, M. R. Vollger and E. E. Eichler, "Long-read human genome sequencing and its applications," *Nature Reviews Genetics*, vol. 21, no. 10, pp. 597--614, 2020.
- [5] S. L. Amarasinghe, S. Su, X. Dong, L. Zappia, M. E. Ritchie and Q. Gouil, "Opportunities and challenges in long-read sequencing data analysis," *Genome biology*, vol. 21, no. 1, pp. 1--16, 2020.
- [6] A. Rhie, S. A. McCarthy, O. Fedrigo, J. Damas, G. Formenti, S. Koren, M. Uliano-Silva, W. Chow, A. Fungtammasan, J. Kim and others, "Towards complete and error-free genome assemblies of all vertebrate species," *Nature*, vol. 592, no. 7856, pp. 737--746, 2021.
- [7] H. A. Lewin, G. E. Robinson, W. J. Kress, W. J. Baker, J. Coddington, K. A. Crandall, R. Durbin, S. V. Edwards, F. Forest, M. T. P. Gilbert and others, "Earth BioGenome Project: Sequencing life for the future of life," *Proceedings of the National Academy of Sciences*, vol. 115, no. 17, pp. 4325--4333, 2018.

- [8] S. Selvaraj, J. R. Dixon, V. Bansal and B. Ren, "Whole-genome haplotype reconstruction using proximity-ligation and shotgun sequencing," *Nature biotechnology*, vol. 31, no. 12, pp. 1111–1118, 2013.
- [9] F. J. Sedlazeck, H. Lee, C. A. Darby and M. C. Schatz, "Piercing the dark matter: bioinformatics of long-range sequencing and mapping," *Nature Reviews Genetics*, vol. 19, no. 6, pp. 329–346, 2018.
- [10] A. Bankevich and P. A. Pevzner, "TruSPAdes: barcode assembly of TruSeq synthetic long reads," *Nature methods*, vol. 13, no. 3, pp. 248–250, 2016.
- [11] I. Wu, H. S. Kim and T. Ben-Yehezkel, "A single-molecule long-read survey of human transcriptomes using LoopSeq synthetic long read sequencing," *bioRxiv*, p. 532135, 2019.
- [12] Z. Chen, L. Pham, T.-C. Wu, G. Mo, Y. Xia, P. L. Chang, D. Porter, T. Phan, H. Che, H. Tran and others, "Ultralow-input single-tube linked-read library method enables short-read second-generation sequencing systems to routinely generate highly accurate and economical long-range sequencing information," *Genome research*, vol. 30, no. 6, pp. 898–909, 2020.
- [13] J. I. Meier, P. A. Salazar, M. Ku\vcka, R. W. Davies, A. Dr\u00e9au, I. Ald\u00e1s, O. Box Power, N. J. Nadeau, J. R. Bridle, C. Rolian and others, "Haplotype tagging reveals parallel formation of hybrid races in two butterfly species," *Proceedings of the National Academy of Sciences*, vol. 118, no. 25, p. e2015005118, 2021.
- [14] E. Lieberman-Aiden, N. L. Van Berkum, L. Williams, M. Imakaev, T. Ragoczy, A. Telling, I. Amit, B. R. Lajoie, P. J. Sabo, M. O. Dorschner and others, "Comprehensive mapping of long-range interactions reveals folding principles of the human genome," *science*, vol. 326, no. 5950, pp. 289–293, 2009.
- [15] J. O. Korbelt and C. Lee, "Genome assembly and haplotyping with Hi-C," *Nature biotechnology*, vol. 31, no. 12, pp. 1099–1101, 2013.
- [16] N. I. Weisenfeld, V. Kumar, P. Shah, D. M. Church and D. B. Jaffe, "Direct determination of diploid genome sequences," *Genome research*, vol. 27, no. 5, pp. 757–767, 2017.
- [17] 1. Genomics, "Linked Reads, Retrieved from <https://www.10xgenomics.com/products/linked-reads> on 2020-10-17.," 2020.
- [18] J. A. Bhat, D. Yu, A. Bohra, S. A. Ganie and R. K. Varshney, "Features and applications of haplotypes in crop breeding," *Communications biology*, vol. 4, no. 1, pp. 1–12, 2021.
- [19] M. S. Lindner, M. Kollock, F. Zickmann and B. Y. Renard, "Analyzing genome coverage profiles with applications to quality control in metagenomics," *Bioinformatics*, vol. 29, no. 10, pp. 1260–1267, 2013.
- [20] S. K. Shore, L. T. Bacheler, J. Kimball de Riel, L. R. Barrows and M. Lynch, "Cloning and characterization of a rat-specific repetitive DNA sequence," *Gene*, vol. 45, no. 1, pp. 87–93, 1986.
- [21] T. R. Gregory, C. B. Andrews, J. A. McGuire and C. C. Witt, "The smallest avian genomes are found in hummingbirds," *Proceedings of the Royal Society B: Biological Sciences*, vol. 276, no. 1674, pp. 3753–3757, 2009.
- [22] J. Luo, Y. Wei, M. Lyu, Z. Wu, X. Liu, H. Luo and C. Yan, "A comprehensive review of scaffolding methods in genome assembly," *Briefings in Bioinformatics*, vol. 22, no. 5, p. bbab033, 2021.
- [23] E. S. Rice, R. E. Green and others, "New approaches for genome assembly and scaffolding," *Annu Rev Anim Biosci*, vol. 7, no. 1, pp. 17–40, 2019.
- [24] S. Niu, J. Li, W. Bo, W. Yang, A. Zuccolo, S. Giacomello, X. Chen, F. Han, J. Yang, Y. Song and others, "The Chinese pine genome and methylome unveil key features of conifer evolution," *Cell*, vol. 185, no. 1, pp. 204–217, 2022.
- [25] K. Yamaguchi, M. Kadota, O. Nishimura, Y. Ohishi, Y. Naito and S. Kuraku, "Technical considerations in Hi-C scaffolding and evaluation of chromosome-scale genome assemblies," *Molecular Ecology*, vol. 30, no. 23, pp. 5923–5934, 2021.
- [26] H. Cheng, G. T. Concepcion, X. Feng, H. Zhang and H. Li, "Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm," *Nature methods*, vol. 18, no. 2, pp. 170–175, 2021.
- [27] C. Zhou, S. A. McCarthy and R. Durbin, "YaHS: yet another Hi-C scaffolding tool," *bioRxiv*, 2022.
- [28] H.-Q. Ling, B. Ma, X. Shi, H. Liu, L. Dong, H. Sun, Y. Cao, Q. Gao, S. Zheng, Y. Li and others, "Genome sequence of the progenitor of wheat A subgenome *Triticum urartu*," *Nature*, vol. 557, no. 7705, pp. 424–428, 2018.
- [29] H. Lee, H. S. Chawla, C. Obermeier, F. Dreyer, A. Abbadi and R. Snowdon, "Chromosome-scale assembly of winter oilseed rape *Brassica napus*," *Frontiers in plant science*, vol. 11, p. 496, 2020.
- [30] J.-i. Sohn and J.-W. Nam, "The present and future of de novo whole-genome assembly," *Briefings in bioinformatics*, vol. 19, no. 1, pp. 23–40, 2018.
- [31] M. B. Couger, L. Ar\u00e9valo and P. Campbell, "A high quality genome for *Mus spicilegus*, a close relative of house mice with unique social and ecological adaptations," *G3: Genes, Genomes, Genetics*, vol. 8, no. 7, pp. 2145–2152, 2018.
- [32] M. Y. Ozerov, M. Flaj\vshans, K. Noreikiene, A. Vasem\u00e4gi and R. Gross, "Draft genome assembly of the freshwater apex predator wels catfish (*Silurus glanis*) using linked-read sequencing," *G3: Genes, Genomes, Genetics*, vol. 10, no. 11, pp. 3897–3906, 2020.
- [33] A. M. Hulse-Kemp, S. Maheshwari, K. Stoffel, T. A. Hill, D. Jaffe, S. R. Williams, N. Weisenfeld, S. Ramakrishnan, V. Kumar, P. Shah and others, "Reference quality assembly of the 3.5-Gb genome of *Capsicum annuum* from a single linked-read library," *Horticulture research*, vol. 5, 2018.
- [34] F. Chan, *Private Communications*, 2022.
- [35] S. S. Batra, M. Levy-Sakin, J. Robinson, J. Guillory, S. Durinck, T. P. Vilgalys, P.-Y. Kwok, L. A. Cox, S. Seshagiri, Y. S. Song and others, "Accurate assembly of the olive baboon (*Papio anubis*) genome using long-read and Hi-C data," *Gigascience*, vol. 9, no. 12, p. giaa134, 2020.
- [36] A. L. Lind, Y. Y. Lai, Y. Mostovoy, A. K. Holloway, A. Iannucci, A. C. Mak, M. Fondi, V. Orlandini, W. L. Eckalbar, M. Milan and others, "Genome of the Komodo dragon reveals adaptations in the cardiovascular and chemosensory systems of monitor lizards," *Nature ecology & evolution*, vol. 3, no. 8, pp. 1241–1252, 2019.

- [37] N. Spies, Z. Weng, A. Bishara, J. McDaniel, D. Catoe, J. M. Zook, M. Salit, R. B. West, S. Batzoglou and A. Sidow, "Genome-wide reconstruction of complex structural variants using read clouds," *Nature methods*, vol. 14, no. 9, pp. 915--920, 2017.
- [38] P. Marks, S. Garcia, A. M. Barrio, K. Belhocine, J. Bernate, R. Bharadwaj, K. Bjornson, C. Catalanotti, J. Delaney, A. Fehr and others, "Resolving the full spectrum of human genome variation using Linked-Reads," *Genome research*, vol. 29, no. 4, pp. 635--645, 2019.
- [39] L. Harewood, K. Kishore, M. D. Eldridge, S. Wingett, D. Pearson, S. Schoenfelder, V. P. Collins and P. Fraser, "Hi-C as a tool for precise detection and characterisation of chromosomal rearrangements and copy number variation in human tumours," *Genome biology*, vol. 18, no. 1, pp. 1--11, 2017.
- [40] X. Wang, Y. Luan and F. Yue, "EagleC: A deep-learning framework for detecting a full range of structural variations from bulk and single-cell contact maps," *Science advances*, vol. 8, no. 24, p. eabn9215, 2022.
- [41] A. Chakraborty and F. Ay, "Identification of copy number variations and translocations in cancer cells from Hi-C data," *Bioinformatics*, vol. 34, no. 2, pp. 338--345, 2018.
- [42] J. R. Dixon, J. Xu, V. Dileep, Y. Zhan, F. Song, V. T. Le, G. G. Yardimci, A. Chakraborty, D. V. Bann, Y. Wang and others, "Integrative detection and analysis of structural variation in cancer genomes," *Nature genetics*, vol. 50, no. 10, pp. 1388--1398, 2018.
- [43] E. C. Jacobson, R. S. Grand, J. K. Perry, M. H. Vickers, A. L. Olins, D. E. Olins and J. M. O'Sullivan, "Hi-C detects novel structural variants in HL-60 and HL-60/S4 cell lines," *Genomics*, vol. 112, no. 1, pp. 151--162, 2020.
- [44] U. S. Melo, R. Schöpflin, R. Acuna-Hidalgo, M. A. Mensah, B. Fischer-Zirnsak, M. Holtgrewe, M.-K. Klever, S. Türkmen, V. Heinrich, I. D. Pluym and others, "Hi-C identifies complex genomic rearrangements and TAD-shuffling in developmental diseases," *The American Journal of Human Genetics*, vol. 106, no. 6, pp. 872--884, 2020.
- [45] A. Himmelbach, A. Ruban, I. Walde, H. Simková, J. Doležel, A. Hastie, N. Stein and M. Mascher, "Discovery of multi-megabase polymorphic inversions by chromosome conformation capture sequencing in large-genome plant species," *The Plant Journal*, vol. 96, no. 6, pp. 1309--1316, 2018.
- [46] S. Walkowiak, L. Gao, C. Monat, G. Haberer, M. T. Kassa, J. Brinton, R. H. Ramirez-Gonzalez, M. C. Kolodziej, E. Delorean, D. Thambugala and others, "Multiple wheat genomes reveal global variation in modern breeding," *Nature*, vol. 588, no. 7837, pp. 277--283, 2020.
- [47] H. Sun, B. A. Rowan, P. J. Flood, R. Brandt, J. Fuss, A. M. Hancock, R. W. Michelmore, B. Huettel and K. Schneeberger, "Linked-read sequencing of gametes allows efficient genome-wide analysis of meiotic recombination," *Nature communications*, vol. 10, no. 1, pp. 1--9, 2019.
- [48] R. Rommel Fuentes, T. Hesselink, R. Nieuwenhuis, L. Bakker, E. Schijlen, W. van Doonijeweert, S. Diaz Trivino, J. R. de Haan, G. Sanchez Perez, X. Zhang and others, "Meiotic recombination profiling of interspecific hybrid F1 tomato pollen by linked read sequencing," *The Plant Journal*, vol. 102, no. 3, pp. 480--492, 2020.
- [49] S. R. Viswanathan, G. Ha, A. M. Hoff, J. A. Wala, J. Carrot-Zhang, C. W. Whelan, N. J. Haradhvala, S. S. Freeman, S. C. Reed, J. Rhoades and others, "Structural alterations driving castration-resistant prostate cancer revealed by linked-read genome sequencing," *Cell*, vol. 174, no. 2, pp. 433--447, 2018.
- [50] M. Kawazu, S. Kojima, T. Ueno, Y. Totoki, H. Nakamura, A. Kunita, W. Qu, J. Yoshimura, M. Soda, T. Yasuda and others, "Integrative analysis of genomic alterations in triple-negative breast cancer in association with homologous recombination deficiency," *PLoS genetics*, vol. 13, no. 6, p. e1006853, 2017.
- [51] D. Lutgen, R. Ritter, R.-A. Olsen, H. Schielzeth, J. Gruselius, P. Ewels, J. T. García, H. Shirihai, M. Schweizer, A. Suh and others, "Linked-read sequencing enables haplotype-resolved resequencing at population scale," *Molecular ecology resources*, vol. 20, no. 5, pp. 1311--1322, 2020.
- [52] K.-T. Tan, H. Kim, J. Carrot-Zhang, Y. Zhang, W. J. Kim, G. Kugener, J. A. Wala, T. P. Howard, Y.-Y. Chi, R. Beroukhi and others, "Haplotype-resolved germline and somatic alterations in renal medullary carcinomas," *Genome medicine*, vol. 13, no. 1, pp. 1--13, 2021.
- [53] K. Srikanth, J.-E. Park, D. Lim, J. Cha, S.-R. Cho, I.-C. Cho and W. Park, "A comparison between hi-C and 10X genomics linked read sequencing for whole genome phasing in Hanwoo cattle," *Genes*, vol. 11, no. 3, p. 332, 2020.
- [54] L. Zhang, X. Fang, H. Liao, Z. Zhang, X. Zhou, L. Han, Y. Chen, Q. Qiu and S. C. Li, "A comprehensive investigation of metagenome assembly by linked-read sequencing," *Microbiome*, vol. 8, no. 1, pp. 1--11, 2020.
- [55] M. Roodgar, B. H. Good, N. R. Garud, S. Martis, M. Avula, W. Zhou, S. M. Lancaster, H. Lee, A. Babveyh, S. Nesamoney and others, "Longitudinal linked-read sequencing reveals ecological and evolutionary responses of a human gut microbiome during antibiotic treatment," *Genome research*, vol. 31, no. 8, pp. 1433--1446, 2021.
- [56] J. Ghurye and M. Pop, "Modern technologies and algorithms for scaffolding assembled genomes," *PLoS computational biology*, vol. 15, no. 6, p. e1006994, 2019.
- [57] X. Liao, M. Li, Y. Zou, F.-X. Wu, J. Wang and others, "Current challenges and solutions of de novo assembly," *Quantitative Biology*, vol. 7, no. 2, pp. 90--109, 2019.
- [58] J. N. Burton, A. Adey, R. P. Patwardhan, R. Qiu, J. O. Kitzman and J. Shendure, "Chromosome-scale scaffolding of de novo genome assemblies based on chromatin interactions," *Nature biotechnology*, vol. 31, no. 12, pp. 1119--1125, 2013.
- [59] N. Kaplan and J. Dekker, "High-throughput genome scaffolding from in vivo DNA interaction frequency," *Nature biotechnology*, vol. 31, no. 12, pp. 1143--1147, 2013.
- [60] X. Zhang, S. Zhang, Q. Zhao, R. Ming and H. Tang, "Assembly of allele-aware, chromosomal-scale autopolyploid genomes based on Hi-C data," *Nature plants*, vol. 5, no. 8, pp. 833--845, 2019.
- [61] O. Dudchenko, S. S. Batra, A. D. Omer, S. K. Nyquist, M. Hoeger, N. C. Durand, M. S. Shamim, I. Machol, E. S. Lander, A. P. Aiden and others, "De novo assembly of the *Aedes aegypti* genome using Hi-C yields chromosome-length scaffolds," *Science*, vol. 356, no. 6333, pp. 92--95, 2017.
- [62] J. Ghurye, M. Pop, S. Koren, D. Bickhart and C.-S. Chin, "Scaffolding of long read assemblies using long range contact information," *BMC genomics*, vol. 18, no. 1, pp. 1--11, 2017.

- [63] H. Marie-Nelly, M. Marbouty, A. Cournac, J.-F. Flot, G. Liti, D. P. Parodi, S. Syan, N. Guillén, A. Margeot, C. Zimmer and others, "High-quality genome (re) assembly using chromosomal contact data," *Nature communications*, vol. 5, no. 1, pp. 1--10, 2014.
- [64] L. Baudry, N. Guiguelmoni, H. Marie-Nelly, A. Cormier, M. Marbouty, K. Avia, Y. L. Mie, O. Godfroy, L. Sterck, J. M. Cock and others, "instaGRAAL: chromosome-level quality scaffolding of genomes using a proximity ligation-based scaffold," *Genome biology*, vol. 21, no. 1, pp. 1--22, 2020.
- [65] N. H. Putnam, B. L. O'Connell, J. C. Stites, B. J. Rice, M. Blanchette, R. Calef, C. J. Troll, A. Fields, P. D. Hartley, C. W. Sugnet and others, "Chromosome-scale shotgun assembly using an in vitro method for long-range linkage," *Genome research*, vol. 26, no. 3, pp. 342--350, 2016.
- [66] R. Nakabayashi and S. Morishita, "HiC-Hiker: a probabilistic model to determine contig orientation in chromosome-length scaffolds with Hi-C," *Bioinformatics*, vol. 36, no. 13, pp. 3966--3974, 2020.
- [67] S. Wang, H. Wang, F. Jiang, A. Wang, H. Liu, H. Zhao, B. Yang, D. Xu, Y. Zhang and W. Fan, "EndHiC: assemble large contigs into chromosomal-level scaffolds using the Hi-C links from contig ends," *arXiv preprint arXiv:2111.15411*, 2021.
- [68] A. Sur, W. S. Noble and P. J. Myler, "A benchmark of Hi-C scaffolders using reference genomes and de novo assemblies," *bioRxiv*, 2022.
- [69] E. Vidal, F. le Dily, J. Quilez, R. Stadhouders, Y. Cuartero, T. Graf, M. A. Marti-Renom, M. Beato and G. J. Filion, "OneD: increasing reproducibility of Hi-C samples with abnormal karyotypes," *Nucleic acids research*, vol. 46, no. 8, pp. e49--e49, 2018.
- [70] S. Wang, S. Lee, C. Chu, D. Jain, P. Kerpedjiev, G. M. Nelson, J. M. Walsh, B. H. Alver and P. J. Park, "HiNT: a computational method for detecting copy number variations and translocations from Hi-C data," *Genome biology*, vol. 21, no. 1, pp. 1--15, 2020.
- [71] D. Jain, C. Chu, B. H. Alver, S. Lee, E. A. Lee and P. J. Park, "HiTea: a computational pipeline to identify non-reference transposable element insertions in Hi-C data," *Bioinformatics*, vol. 37, no. 8, pp. 1045--1051, 2021.
- [72] X. Wang, J. Xu, B. Zhang, Y. Hou, F. Song, H. Lyu and F. Yue, "Genome-wide detection of enhancer-hijacking events from chromatin interaction data in rearranged genomes," *Nature methods*, vol. 18, no. 6, pp. 661--668, 2021.
- [73] A. Adey, J. O. Kitzman, J. N. Burton, R. Daza, A. Kumar, L. Christiansen, M. Ronaghi, S. Amini, K. L. Gunderson, F. J. Steemers and others, "In vitro, long-range sequence information for de novo genome assembly via transposase contiguity," *Genome research*, vol. 24, no. 12, pp. 2041--2049, 2014.
- [74] L. Coombe, J. Zhang, B. P. Vandervalk, J. Chu, S. D. Jackman, I. Birol and R. L. Warren, "ARKS: chromosome-scale scaffolding of human genome drafts with linked read kmers," *BMC bioinformatics*, vol. 19, no. 1, pp. 1--10, 2018.
- [75] S. Yeo, L. Coombe, R. L. Warren, J. Chu and I. Birol, "ARCS: scaffolding genome drafts with linked reads," *Bioinformatics*, vol. 34, no. 5, pp. 725--731, 2018.
- [76] M. Hiltunen, M. Ryberg and H. Johannesson, "ARBitR: an overlap-aware genome assembly scaffold for linked reads," *Bioinformatics*, vol. 37, no. 15, pp. 2203--2205, 2021.
- [77] L. Guo, M. Xu, W. Wang, S. Gu, X. Zhao, F. Chen, O. Wang, X. Xu, I. Seim, G. Fan and others, "SLR-superscaffolder: a de novo scaffolding tool for synthetic long reads using a top-to-bottom scheme," *BMC bioinformatics*, vol. 22, no. 1, pp. 1--16, 2021.
- [78] V. Kuleshov, M. P. Snyder and S. Batzoglu, "Genome assembly from synthetic long read clouds," *Bioinformatics*, vol. 32, no. 12, pp. i216--i224, 2016.
- [79] K. Paszkiewicz and D. J. Studholme, "De novo assembly of short sequence reads," *Briefings in bioinformatics*, vol. 11, no. 5, pp. 457--472, 2010.
- [80] A. R. Khan, M. T. Pervez, M. E. Babar, N. Naveed and M. Shoab, "A comprehensive study of de novo genome assemblers: current challenges and future prospective," *Evolutionary Bioinformatics*, vol. 14, p. 1176934318758650, 2018.
- [81] I. Tolstoganov, A. Bankevich, Z. Chen and P. A. Pevzner, "cloudSPAdes: assembly of synthetic long reads using de Bruijn graphs," *Bioinformatics*, vol. 35, no. 14, pp. i61--i70, 2019.
- [82] L. Mak, D. Meleshko, D. C. Danko, W. N. Barakzai, N. Belchikov and I. Hajirasouliha, "Ariadne: Synthetic Long Read Deconvolution Using Assembly Graphs," *BioRxiv*, pp. 2021--05, 2022.
- [83] A. Bishara, E. L. Moss, M. Kolmogorov, A. E. Parada, Z. Weng, A. Sidow, A. E. Dekas, S. Batzoglu and A. S. Bhatt, "High-quality genome sequences of uncultured microbes by assembly of read clouds," *Nature biotechnology*, vol. 36, no. 11, pp. 1067--1075, 2018.
- [84] M. A. DePristo, E. Banks, R. Poplin, K. V. Garimella, J. R. Maguire, C. Hartl, A. A. Philippakis, G. Del Angel, M. A. Rivas, M. Hanna and others, "A framework for variation discovery and genotyping using next-generation DNA sequencing data," *Nature genetics*, vol. 43, no. 5, pp. 491--498, 2011.
- [85] Z. Wei, W. Wang, P. Hu, G. J. Lyon and H. Hakonarson, "SNVer: a statistical tool for variant calling in analysis of pooled or individual next-generation sequencing data," *Nucleic acids research*, vol. 39, no. 19, pp. e132--e132, 2011.
- [86] D. C. Koboldt, Q. Zhang, D. E. Larson, D. Shen, M. D. McLellan, L. Lin, C. A. Miller, E. R. Mardis, L. Ding and R. K. Wilson, "VarScan 2: somatic mutation and copy number alteration discovery in cancer by exome sequencing," *Genome research*, vol. 22, no. 3, pp. 568--576, 2012.
- [87] Z. Lai, A. Markovets, M. Ahdesmaki, B. Chapman, O. Hofmann, R. McEwen, J. Johnson, B. Dougherty, J. C. Barrett and J. R. Dry, "VarDict: a novel and versatile variant caller for next-generation sequencing in cancer research," *Nucleic acids research*, vol. 44, no. 11, pp. e108--e108, 2016.
- [88] P. Morisse, F. Legeai and C. Lemaitre, "LEVIATHAN: efficient discovery of large structural variants by leveraging long-range information from Linked-Reads data," *bioRxiv*, 2021.
- [89] M. Eslami Rasekh, G. Chiatante, M. Miroballo, J. Tang, M. Ventura, C. T. Amemiya, E. E. Eichler, F. Antonacci and C. Alkan, "Discovery of large genomic inversions using long range information," *BMC genomics*, vol. 18, no. 1, pp. 1--12, 2017.

- [90] F. Karaoğlanoğlu, C. Ricketts, E. Ebrén, M. E. Rasekh, I. Hajirasouliha and C. Alkan, "VALOR2: characterization of large-scale structural variants using linked-reads," *Genome biology*, vol. 21, no. 1, pp. 1--11, 2020.
- [91] L. Fang, C. Kao, M. V. Gonzalez, F. A. Mafra, R. Pellegrino da Silva, M. Li, S.-S. Wenzel, K. Wimmer, H. Hakonarson and K. Wang, "LinkedSV for detection of mosaic structural variants from linked-read exome and genome sequencing data," *Nature communications*, vol. 10, no. 1, pp. 1--15, 2019.
- [92] R. Elyanow, H.-T. Wu and B. J. Raphael, "Identifying structural variants using linked-read sequencing data," *Bioinformatics*, vol. 34, no. 2, pp. 353--360, 2018.
- [93] L. C. Xia, J. M. Bell, C. Wood-Bouwens, J. J. Chen, N. R. Zhang and H. P. Ji, "Identification of large rearrangements in cancer genomes with barcode linked reads," *Nucleic acids research*, vol. 46, no. 4, pp. e19--e19, 2018.

888

889

890

891

## 892 Acknowledgements

893

894 We thank Dr Anthony Schmitt, Arima Genomics for providing data to evaluate Hi-C  
895 technology and its developments. We are also grateful to Professor Yingguang Frank  
896 Chan, Friedrich Miescher Laboratory of the Max Planck Society, for providing much  
897 needed technical support on haplotagging library preparation and sequencing.

898 Funding: L.J. is supported by grant 32070601 from the National Natural Science  
899 Foundation of China. Z.N., M.Q., J.F-G., K.O. and E.M. G. are supported by the  
900 Wellcome Trust (WT206194). Author contributions: L.J., M.Q. and Z.N. proposed  
901 and designed the project. M.Q., K.O. and E.M.G. made the haplotagging sequencing  
902 libraries and produced the data. Z.N. performed data analysis, while J.F-G. developed  
903 the statistic model for coverage assessment and barcode collisions. H.W., X.S. and  
904 L.J. drew Figure 1. L.J., J.F-G. and Z.N. wrote the paper. All the authors read and  
905 approved the final manuscript.

906

## 907 Data availability

908

909 Hi-C reads were sequenced and provided by Arima Genomics and have been submitted  
910 to NCBI under BioProjectID PRJNAxxxxx. The 10x human genome reads were  
911 produced by Genome-in-a-bottle and can be downloaded from [ftp://ftp-  
912 trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/data/NA12878/10Xgenomics\\_Chromi  
913 umGenome\\_LongRanger2.0\\_06202016/NA12878\\_GRCh38.bam](ftp://ftp-trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/data/NA12878/10Xgenomics_Chromi). The 10x  
914 hummingbird dataset has been archived on NCBI/EBI BioProject under

915 accession PRJNA489243. All the haplotagging datasets are part of Darwin Tree of  
916 Life project and have been submitted to NCBI with **BioProjectID PRJNA435xxx**. No  
917 sequencing effort has been paid on the TELL-Seq platform and the contents presented  
918 in this study are based on the TELL-Seq paper, Chen *et al.* [12].

919

920

## 921 List of Figures

922 Figure 1. Work flow of library preparations for four long range platforms. (A) Hi-C;  
923 (B) 10x; (C) haplotagging; (D) Tell-seq.

924 Figure 2. Hi-C contact maps for three human samples. (A) Arima V2 CEU  
925 (NA12878); (B) Arima V2 AJ (NA24385) and (C) Arima V1 CEU (NA12878). Here  
926 the data are shown in the form of two-dimensional symmetric matrices, where x and y  
927 coordinates represent the intensity of the physical interaction between two genome  
928 regions x and y at the DNA level. Each chromosome is seen as a shaded box and  
929 also there are no data points in Chromosome Y as these are female samples.

930 Figure 3. Characteristics of Hi-C reads. (A) length distribution; (B) translocation rate;  
931 (C) base coverage. 3A also contains a fourth dataset, (Arima V2 Oak) which  
932 demonstrates the breakdown of the power-law relationship.

933 Figure 4. Distributions of barcode length for various 10x and haplotagging samples.  
934 Here, for each barcode fragment, the minimum number of read pairs is set to 5 and we  
935 only collect fragments with a length  $\geq 100$ bp.

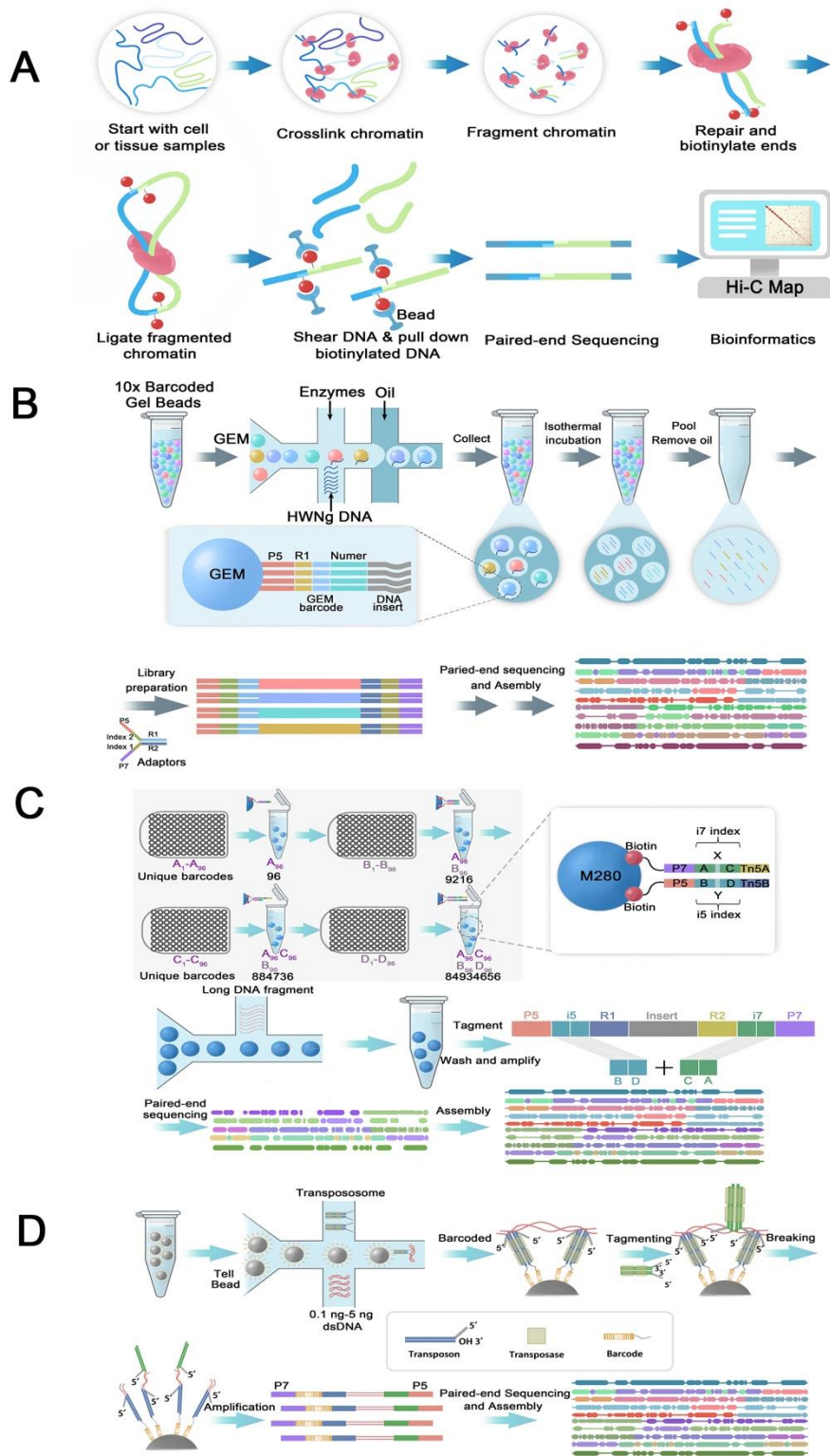
936 Figure 5. Base coverage profiles from various 10x and haplotagging samples with (A)  
937 unnormalized datasets and (B) normalized datasets. Here one Illumina PCR free  
938 dataset is also superimposed for comparison

939 Figure 6. Hi-C contact maps on contigs and scaffolded assemblies. (A) Contigs; (B)  
940 Assembly with V1 data; (C) Assembly with Arima V2 reads. Here contigs are shown  
941 as individual boxes along the diagonal direction in (A). Scaffolds are seen in (B) and  
942 (C), where scaffolding level in (C) is better than that in (B).

943

944





945

946 Figure 1

947

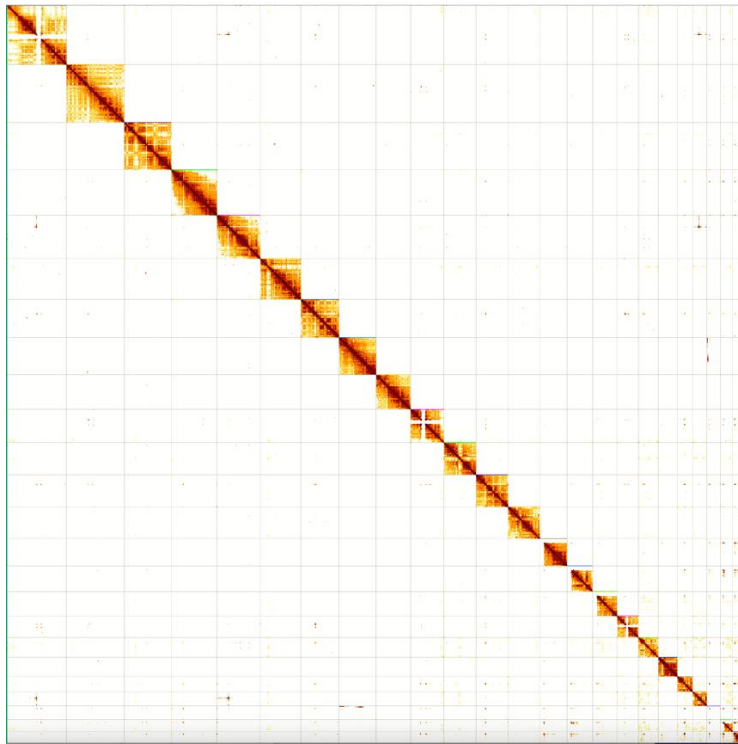
948

949

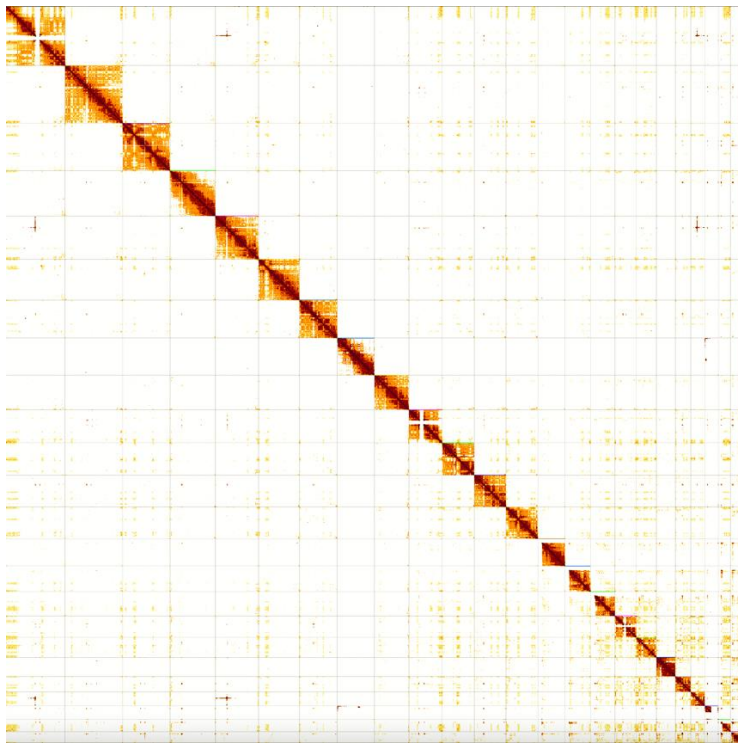
A

B

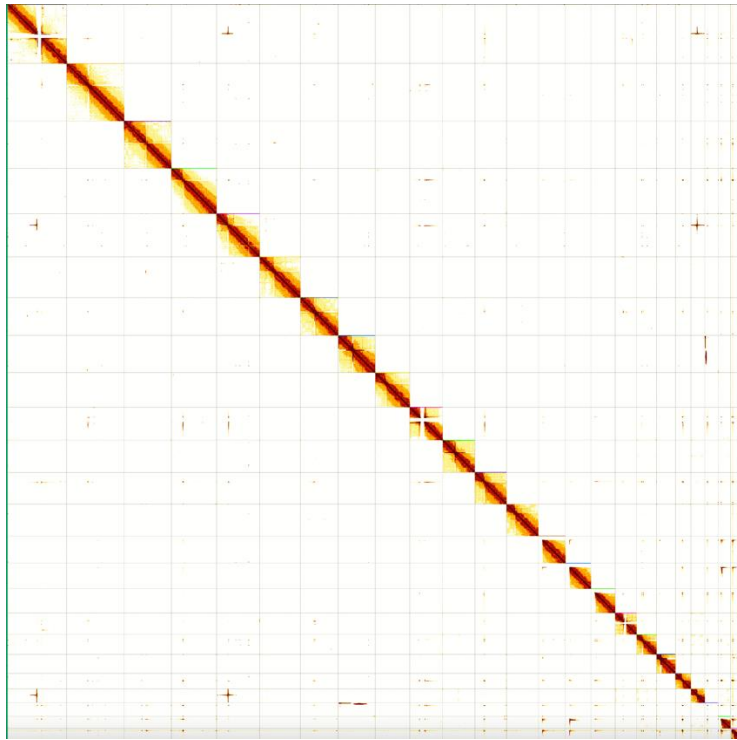
C



950



951



952  
953  
954

955 Figure 2

956  
957  
958  
959

960

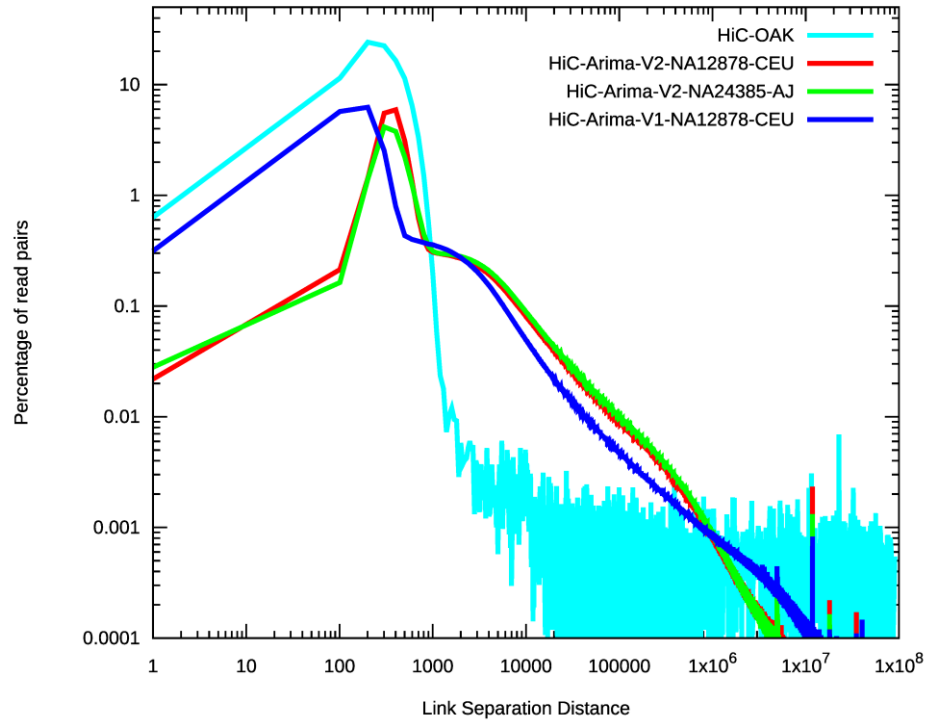
961

962 A

963

B

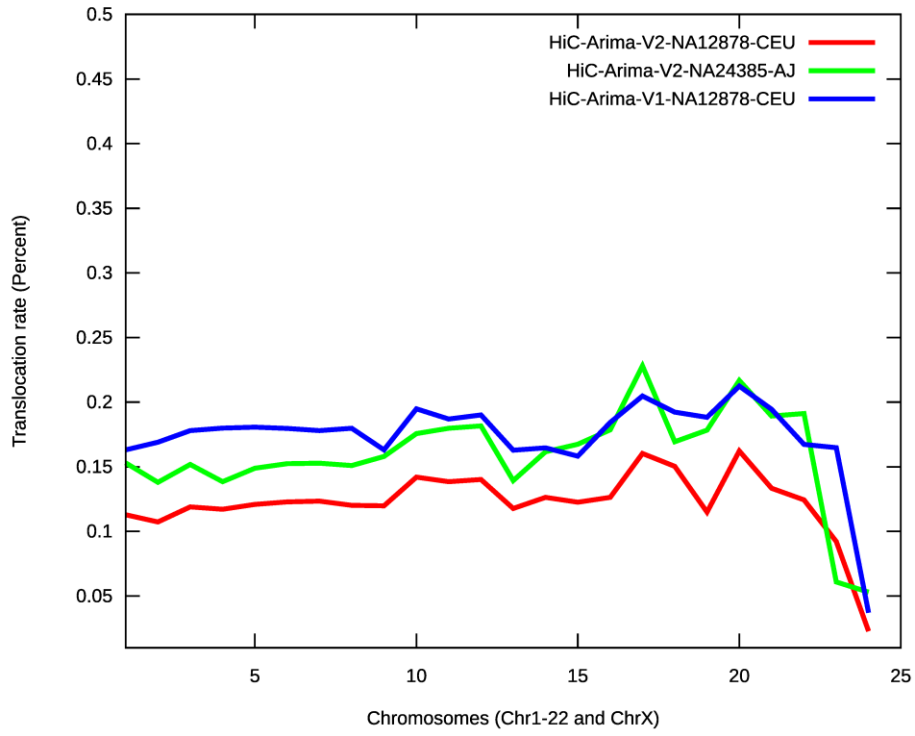
C



964

965

966

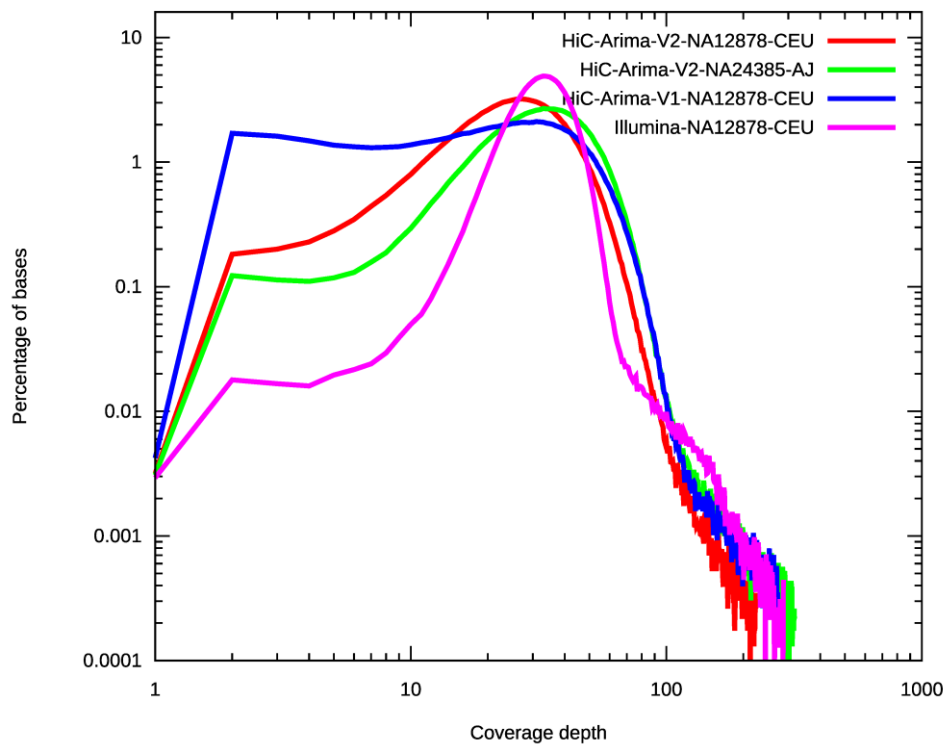


967

968

969

970



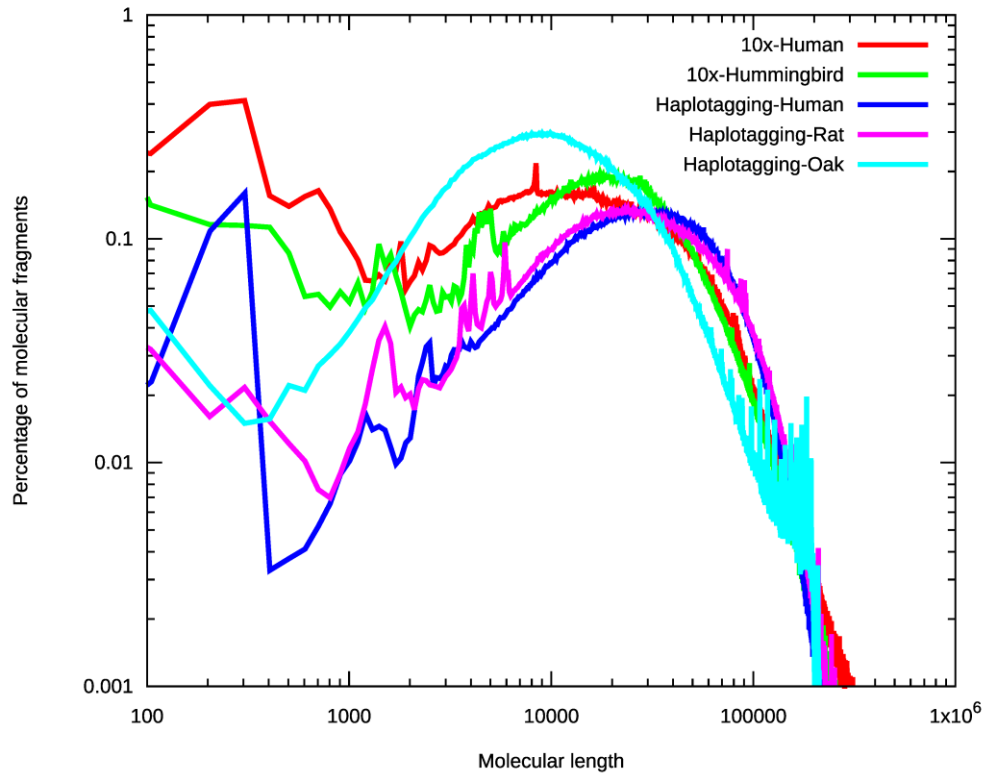
971

972 Figure 3

973

974

975



976

977

978 Figure 4

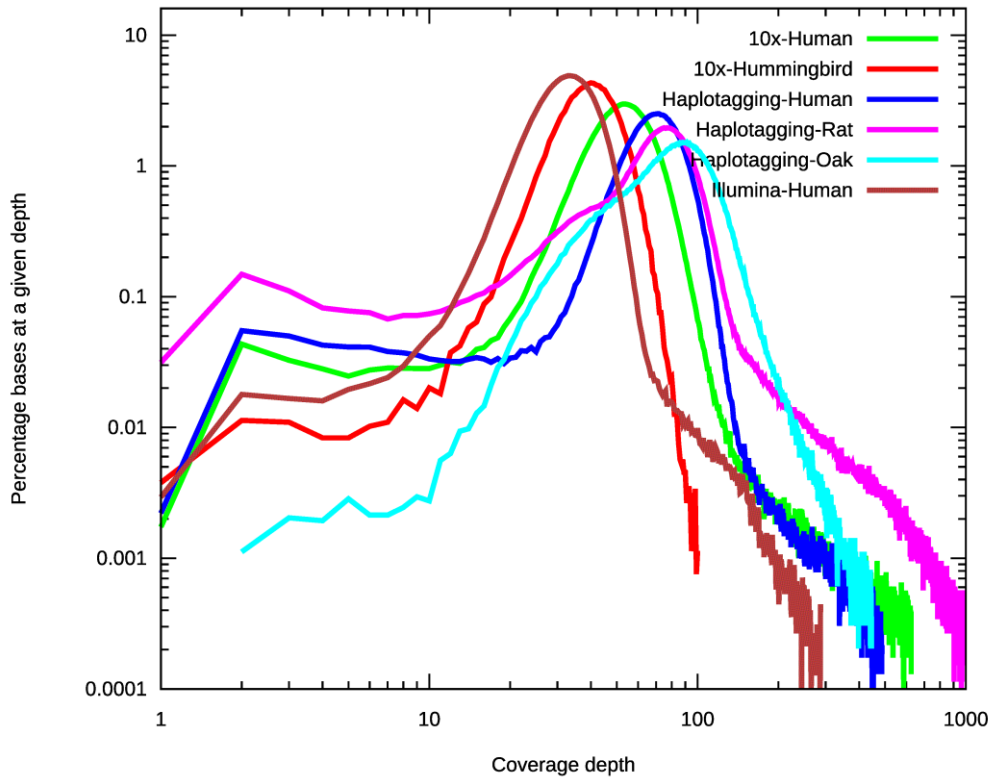
979

980

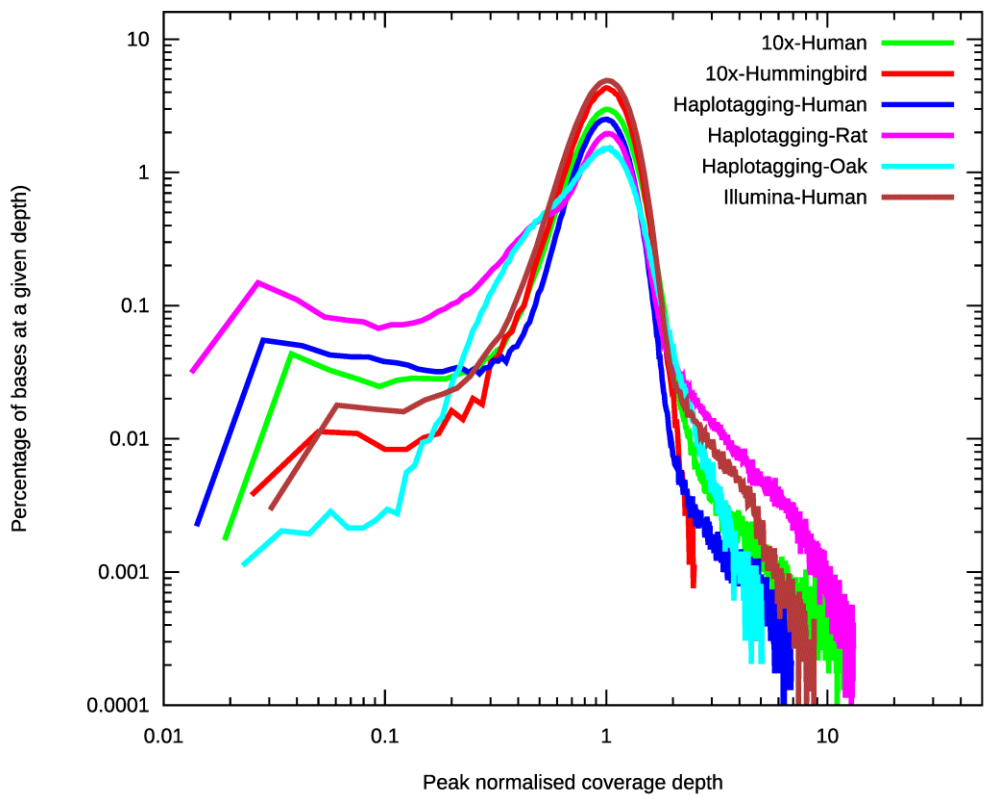
981 A

B

982



983



984

985

Figure 5

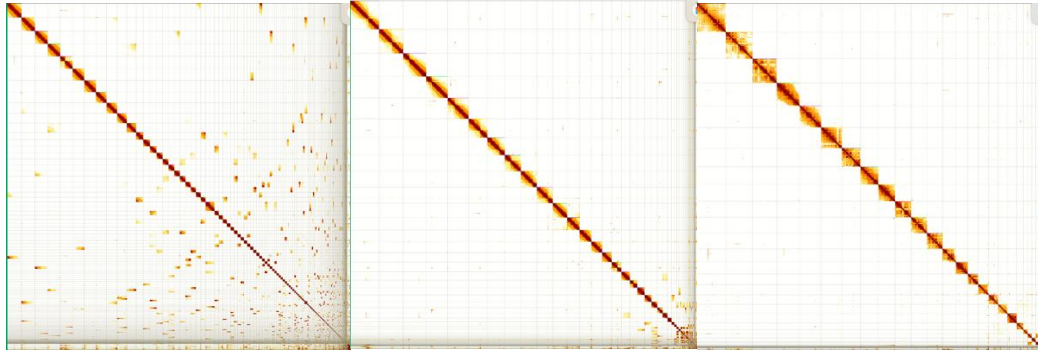
986

987

988 A

B

C



989

990 Figure 6

991



992

993 Table 1: Features of Hi-C reads with three different datasets

Datasets	Read pairs	Unmapped (%)	PCR duplicated (%)	Translocated (%)	Accessibility	N50 (Mb)	N20 (Mb)	N10 (Mb)
Arima V2 NA12878- CEU	352,429,304	20.9	6.8	12.7	0.596	47.9	96.3	130.3
Arima V2 NA24385-AJ	413,162,798	24.8	6.1	16.2	0.529	47.2	104.9	141.9
Arima V1 NA12878- CEU	415,173,112	28.6	10.1	18.5	0.328	28.2	63.1	100.0

994

995

996 Table 2: Features of 10x and Haplotagging linked reads technology with different samples. For the Tools, LR+S

997 means the LongRanger tool and scaff10x, whilst EMA+S means EMA and scaff10x.

Species	Platform	Tools	Read pairs	Unmapped Reads (%)	PCR duplication (%)	Accessability	Molecular length (N50)	Barcode N50 (>=5 reads)	Barcode N50 (>=3 reads)	Reads clustered (%)
Human- NA12878	10x	LR+S	669,583,370	10.1	21.0	0.689	94,611	81	80	88.2
Hummingbird	10x	LR+S	159,605,373	14.9	5.0	0.801	63,292	22	21	87
Human	Haplotagging	EMA+S	678,683,208	2.52	30.7	0.678	73,294	12	11	55
Rat	Haplotagging	EMA+S	742,824,305	2.5	30.0	0.675	78,250	11	11	50
Oak	Haplotagging	EMA+S	208,869,403	4.17	39.5	0.563	54,969	11	8	50

998

999 Table 3: Coverage Evenness statistics

Datasets	Coverage Mean	Coverage Variance	Unevenness
Arima V2 NA12878-CEU	31.2	195.0	5.3
Arima v2 NA24385-AJ	38.8	273.9	6.1
Arima V1 NA12878-CEU	32.1	363.1	10.3
10x NA12878	57.2	565.6	8.9
Haplotagging	73.0	446.0	5.1
Illumina	35.1	131.2	2.7

Nonhuman	10x Hummingbird	41.8	98.6	1.36
	Haplotagging Rat	83.3	3715.1	43.6
	Haplotagging Oak	90.64	1133.9	11.5

1000

1001

1002

1003

1004

1005 Table 4: Assembly stats from different Hi-C datasets

Data and Assembly	Total Bases (Gb)	Number of Sequences	SEQUENCE LENGTH (MB)			
			Mean	N50	N90	Maximum
HIFI READS – HG002	167.76	14,949,433	0.011	0.011	0.009	0.021
CONTIGS WITH HIFIASM	2.866	1,126	2.54	45.1	6.93	116
CONTIGS AFTER PURGE_DUPS	2.83	434	6.52	45.1	8.46	116
SCAFFOLDS WITH ARIMA V1	2.83	162	17.5	152	78.4	324
SCAFFOLDS WITH ARIMA V2	2.83	152	18.6	144	75.12	235

1006

1007

<i>Application</i>	<i>Software Tool</i>	<i>Year</i>	<i>Properties</i>	<i>URL</i>
<i>Hi-C</i>	LACHESIS	2013		<a href="https://github.com/shendurelab/LACHESIS">https://github.com/shendurelab/LACHESIS</a>
	dnaTri	2013	Deterministic, agglomerative hierarchical clustering	<a href="https://github.com/NoamKaplan/dna-triangulation">https://github.com/NoamKaplan/dna-triangulation</a>
	GRAAL	2014	Probabilistic MCMC	<a href="https://github.com/koszullab/GRAAL">https://github.com/koszullab/GRAAL</a>
	instaGRAAL	2020	Probabilistic MCMC, refined for large genomes	<a href="https://github.com/koszullab/instaGRAAL">https://github.com/koszullab/instaGRAAL</a>
	SALSA2	2019	Novel iterative scaffolding method	<a href="https://github.com/marbl/SALSA">https://github.com/marbl/SALSA</a>
	3D-DNA	2017	Deterministic best-neighbour, megascaffold approach	<a href="https://github.com/aidenlab/3d-dna">https://github.com/aidenlab/3d-dna</a>
	HiRise	2016	Maximum Likelihood algorithm, official Dovetail product	<a href="https://github.com/DovetailGenomics/HiRise">https://github.com/DovetailGenomics/HiRise</a>
	ALLHiC	2019	Deterministic hierarchical clustering on autopolyploid or heterozygous genomes	<a href="https://github.com/tangerzhang/ALLHiC">https://github.com/tangerzhang/ALLHiC</a>
	HIC-Hiker	2020	Probabilistic, dynamic programming approach to improve quality of already-scaffolded data	<a href="https://github.com/ryought/hic_hiker">https://github.com/ryought/hic_hiker</a>
	EndHic	2021	Improves quality of already-scaffolded data	<a href="https://github.com/fanagislab/EndHic">https://github.com/fanagislab/EndHic</a>
YaHS	2022	Probabilistic, novel inference algorithm	<a href="https://github.com/c-zhou/yahs">https://github.com/c-zhou/yahs</a>	
<i>Variation detection</i>	HiCnv	2018	Detects Copy Number Variations using Hidden Markov Models	<a href="https://github.com/ay-lab/HiCnv">https://github.com/ay-lab/HiCnv</a>
	OneD	2018		<a href="https://github.com/qenvio/dryhic">https://github.com/qenvio/dryhic</a>
	HiCtrans	2018	Detects Translocations using Hidden Markov Models	<a href="https://github.com/ay-lab/HiCtrans">https://github.com/ay-lab/HiCtrans</a>
	HiNT	2020	Detects both CNV and Translocations	<a href="https://github.com/parklab/HiNT">https://github.com/parklab/HiNT</a>
	HiTea	2021	Identifies mobile transposable element insertions	<a href="https://github.com/parklab/HiTea">https://github.com/parklab/HiTea</a>
	NeoLoopFinder	2021	Finds SV-induced chromatin loops	<a href="https://github.com/XiaoTaoWang/NeoLoopFinder">https://github.com/XiaoTaoWang/NeoLoopFinder</a>
EagleC	2022	Deep-Learning method for full-spectrum SV detection	<a href="https://github.com/XiaoTaoWang/EagleC">https://github.com/XiaoTaoWang/EagleC</a>	
<i>Chain-Linked Reads</i>	fragScaf	2014	Agglomerative hierarchical clustering	<a href="https://github.com/adeylab/fragScaff">https://github.com/adeylab/fragScaff</a>
	Architect	2016		<a href="https://github.com/kuleshov/architect">https://github.com/kuleshov/architect</a>
	ARCS	2018	Designed specifically for 10x	<a href="https://github.com/bcgsc/arcs">https://github.com/bcgsc/arcs</a>
	ARKS	2018	k-mer mapping for improved efficiency in ARCS	
	ARBitR	2021	Explicitly designed for multiple CLR platforms	<a href="https://github.com/markhil/ARBitR">https://github.com/markhil/ARBitR</a>
	SLR-superscaffolder	2021	Divisive hierarchical clustering	<a href="https://github.com/BGI-Qingdao/SLR-superscaffolder">https://github.com/BGI-Qingdao/SLR-superscaffolder</a>
	Supernova	2017	Official 10x assembly product	<a href="https://support.10xgenomics.com/de-novo-assembly/software/overview/latest/welcome">https://support.10xgenomics.com/de-novo-assembly/software/overview/latest/welcome</a>
	cloudSPAdes	2019	De Bruijn assembler, extensible to metagenomic or hybrid data	<a href="https://github.com/ablab/spades/releases/tag/cloudspades-paper">https://github.com/ablab/spades/releases/tag/cloudspades-paper</a>
	Ariadne	2021	cloudSPAdes module, deconvolves barcodes accurately	<a href="https://github.com/lauren-mak/ariadne">https://github.com/lauren-mak/ariadne</a>
	Athena	2018	Improves metagenomic assembly	<a href="https://github.com/abishara/athena_meta">https://github.com/abishara/athena_meta</a>
TuringAssembler	2020			
TELL-seq data analysis pipeline	2020	Introduced explicitly for TELL-Seq data	<a href="https://universalsequencing.com/software/">https://universalsequencing.com/software/</a>	
Long Ranger	2019	Official 10x variation detection tool, uses augmented GATK approach	<a href="https://support.10xgenomics.com/genome-exome/software/downloads/latest">https://support.10xgenomics.com/genome-exome/software/downloads/latest</a>	
GROC-SVs	2017	Similar to Long Ranger, uses local assembly to improve resolution	<a href="https://github.com/grocsvs/grocsvs">https://github.com/grocsvs/grocsvs</a>	
NAIBR	2018	Probabilistic model using “split molecule” approach	<a href="https://github.com/raphael-group/NAIBR">https://github.com/raphael-group/NAIBR</a>	
VALOR	2017	Detects genomic inversion from “split molecule” signature	<a href="https://github.com/BilkentCompGen/valor">https://github.com/BilkentCompGen/valor</a>	
VALOR2	2020	Expanded from VALOR to detect more types of SV		
ZoomX	2018	Novel probabilistic approach to detect large rearrangements	<a href="https://bitbucket.org/charade/zoomx/src">https://bitbucket.org/charade/zoomx/src</a>	
LEVIATHAN	2021	Can detect SVs in highly fragmented and heterozygosity genomes	<a href="https://github.com/morispi/LEVIATHAN">https://github.com/morispi/LEVIATHAN</a>	

# 1011 Supplementary

## 1012 Evenness Metric

1013 In the ideal case, the sampling of the genome would be perfectly uniform, such that every base  
1014 was covered exactly the same number of times. Since this is practically impossible, we would instead  
1015 prefer that every base had the same chance of being covered, and allowing for some statistical noise. If  
1016 we model the sequencing process as one which samples each base of the genome at a mean rate  $\lambda$ ,  
1017 which is independent of the sampling rate of other bases, then the probability that a given base enters  
1018 the library  $k$  times (i.e. has a coverage of  $k$ ) is:

$$1019 \quad p(k|\lambda) = \frac{\lambda^k e^{-\lambda}}{k!} = \mathcal{P}(k|\lambda)$$

1020 This is the standard Poisson distribution, and would be the result of a perfectly even sampling of the  
1021 genome. If, however, there is not a single value of  $\lambda$ , but multiple different values, such that the  
1022 probability of a given value of  $\lambda$  is given by the distribution function  $f(\lambda)$ , then the probability of  
1023 finding a coverage value of  $k$  is given by a Polypoison distribution,  $p(k|f)$  such that:

$$1024 \quad p(k|f) = \int_0^\infty f(\lambda) \frac{\lambda^k e^{-\lambda}}{k!} d\lambda.$$

1025 The integral is carried out over the full support of the parameter  $\lambda$ , that being the half-infinite interval.  
1026 We note that if  $f$  is a normalized distribution function on this interval, then the total probability still  
1027 obeys:

$$1028 \quad \sum_{k=0}^{\infty} p(k|f) = \int_0^\infty f(\lambda) \left( \sum_{k=0}^{\infty} \frac{\lambda^k e^{-\lambda}}{k!} \right) d\lambda \\ 1029 \quad = \int_0^\infty f(\lambda) d\lambda = 1$$

1030 The mean and the variance of a Polypoison distribution are found from:

$$1031 \quad \langle k \rangle = \sum_{k=0}^{\infty} k p(k|f) = \int_0^\infty \lambda f(\lambda) d\lambda \\ 1032 \quad \text{Var}(k) = \left( \sum_{k=0}^{\infty} k^2 p(k|f) \right) - \langle k \rangle^2 = \int_0^\infty \lambda^2 f(\lambda) d\lambda + \langle k \rangle - \langle k \rangle^2$$

1033 We note that the dimensional conflict of  $\langle k \rangle$  and  $\langle k \rangle^2$  appearing in linear combinations is not a  
1034 problem since the Poisson distribution inherently only deals with dimensionless ‘counts’. Writing the  
1035 results above in terms of the variance and mean of  $f$ , we find that:

$$1036 \quad \langle \text{coverage} \rangle = \langle f \rangle \\ 1037 \quad \text{Var}(\text{coverage}) = \text{Var}(f) + \langle f \rangle$$

1038 Previous works have set  $f(\lambda)$  equal to the Gamma distribution, in which case  $p(k|f)$  is equal to  
1039 the Negative Binomial Distribution. However, we note that there is no particular need to assign a  
1040 functional form to  $f$ , since all we are interested in is the dispersion of this relationship around the  
1041 mean. The index of dispersion is given by:

1042

$$D(f) = \frac{\text{Var}(f)}{\langle f \rangle}$$

1043

Hence:

1044

$$D = \frac{\text{Var}(\text{coverage}) - \langle \text{coverage} \rangle}{\text{Var}(\text{coverage})}$$

1045

We reiterate that this is a measure of the dispersion of  $f$  around its mean, and is hence a measure of how Poisson-like the data is: zero indicates the data is perfectly Poisson like, whilst larger values

1046

indicate that there are a significant number of processes altering how the genome is sampled. We

1047

therefore use this quantity as a metric of the unevenness of the coverage of the genome,  $\mathcal{U} = D$ , with

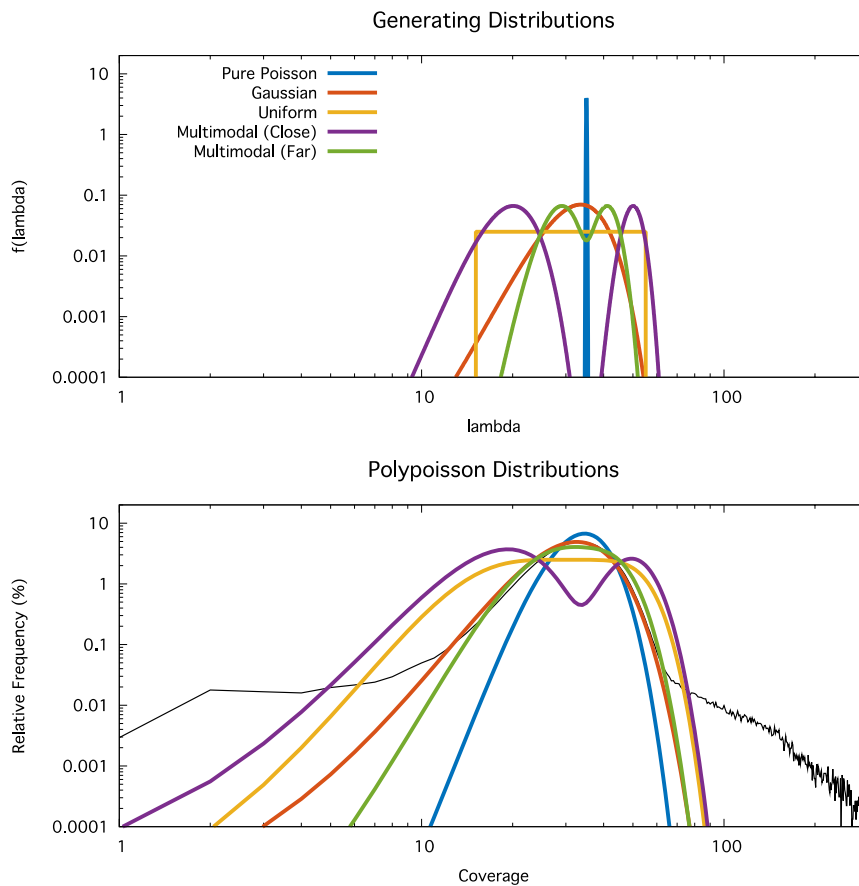
1048

smaller values being indicative of a more even coverage.

1049

1050

1051



1052

1053

Supplementary Figure 1: An example of how generating distributions  $f(\lambda)$  (top) result in different Polypoison

1054

distributions (bottom). All distributions are chosen to have the same mean as the black curve (the Illumina human

1055

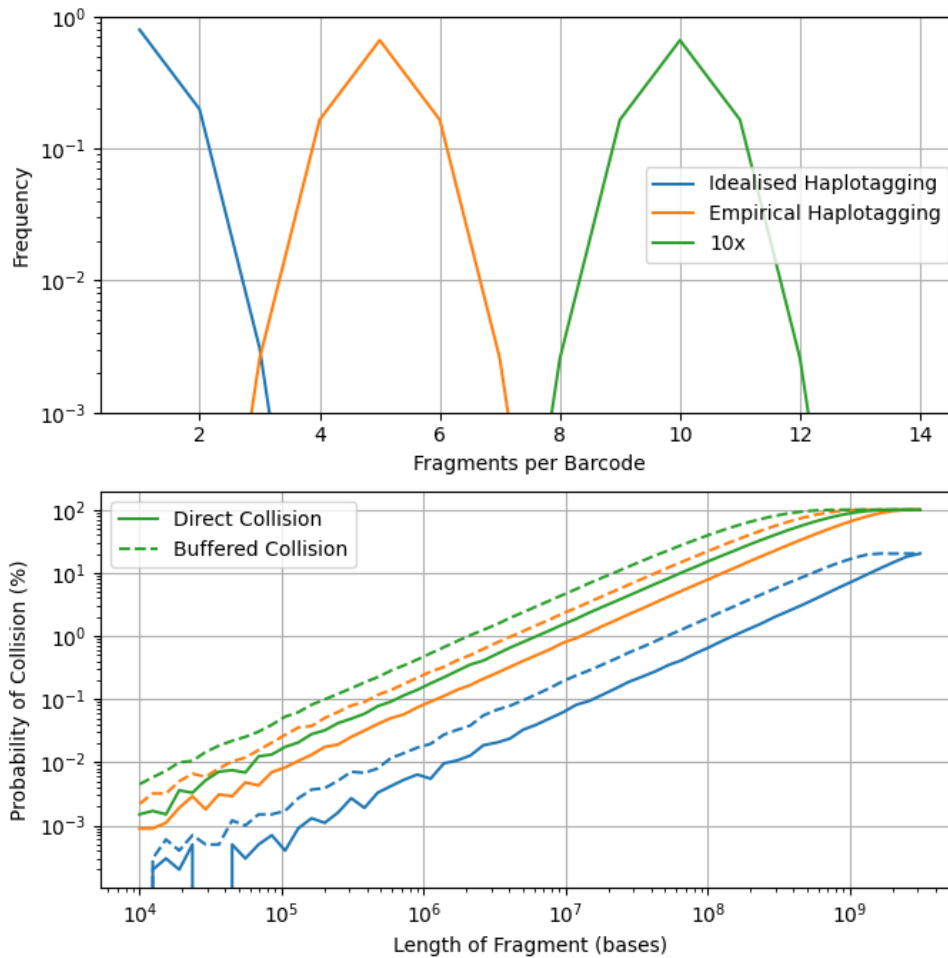
data from Fig. 3), but with other parameters chosen for demonstration purposes rather than to provide a good fit to

1056

the data. The multimodal models demonstrate that even though the Poisson distribution is monomodal, suitable

1057

generating functions can generate multimodel Polypoison distributions.



1058  
 1059  
 1060  
 1061  
 1062  
 1063  
 1064  
 1065  
 1066  
 1067

Supplementary Figure 2: (Top) various distributions of the number of fragments of HMW-DNA which share a barcode, the blue and orange curves are designed to approximate haplotagging, whilst the purple and brown demonstrate 10x. (Bottom) the probability of barcode ‘collisions’ which result, as a function of the length of the fragment, assuming a diploid genome length of 6.3Gb. Solid lines demonstrate direct collisions: overlapping fragments which share the same barcode, whilst the dashed line shows ‘buffered collisions’, where the shared-barcode fragments do not overlap, but are too close together for reads to be unambiguously assigned to one or the other.

Supplementary Table 1: Collision-Frequency Analysis of the Chain-Linked Read platforms

Datasets	Platform	Genome Length	Fragments-per-barcode	1% Collision Size (kb)	Mean Fragment Length (kbp)	Collision Frequency (%)
Human-NA12878	10x	6.3Gb	10	2,100	59.2	0.03
Hummingbird	10x	1.8Gb	10	580	44.6	0.08
Human	Haplotagging	6.3Gb	5	4,200	56.2	0.01
Rat	Haplotagging	5.5Gb	5	3,700	57.2	0.02
Oak	Haplotagging	1.4Gb	5	970	38.5	0.04

## Instructions on running assembly pipelines

1069

1070 Software packages

1071

1072 scaffHiC

1073 <https://github.com/wtsi-hpag/scaffHiC>

1074 Note: scaffHiC contains PretextView and we here use scaffHiC to process data and generate Hi-C maps as well as length distributions. We did not use it for scaffolding as yahs is noticeably better in genome scaffolding.

1075

1076

1077 PretextView

1078 <https://github.com/wtsi-hpag/PretextView>

1079

1080 purge\_dups

1081 [https://github.com/dfguan/purge\\_dups](https://github.com/dfguan/purge_dups)

1082

1083 yahs

1084 <https://github.com/c-zhou/yahs>

1085

1086 samtools

1087 <https://github.com/samtools/>

1088

1089 Produce sorted bam file - AJ.bam

1090 `/nfs/users/nfs_z/zn1/src/scaffHiC/src/scaff-bin/bwa-mem2 mem -t 54 -5SPM GRCH38.fasta`1091 `/lustre/scratch117/sciops/team117/hpag/zn1/project/HiC/arima/human/QC/GM24385.AJ.R1.fastq.gz`1092 `/lustre/scratch117/sciops/team117/hpag/zn1/project/HiC/arima/human/QC/GM24385.AJ.R2.fastq.gz > align-AJ.sam`1093 `samtools view -@ 50 -bS align-AJ.sam > Sorted_names.bam`1094 `samtools fixmate -@ 50 -m Sorted_names.bam Fixmate.bam > try.out`1095 `samtools sort -@ 50 -o Sorted.bam Fixmate.bam > try.out`1096 `rm -rf align-AJ.sam Sorted_names.bam Fixmate.bam`1097 `samtools markdup -@ 50 -r -s Sorted.bam Dupmarked.bam > try.out`1098 `mv Dupmarked.bam AJ.bam`

1099

1100 Coverage analysis

1101 `samtools depth Sorted.bam | egrep _0 | awk '($2%100==0){print $0}' > depth.dat`1102 `sort -n -k 3 depth.dat | awk '{print $1,$3}' > depth-raw.dat`1103 `/nfs/users/nfs_z/zn1/src/scaffHiC/src/scaff-bin/distribution_hic-coverage depth-raw.dat | awk '{print $2,$3}' > depth-freq.dat`

1104

1105 Hi-C contact map

1106 `/nfs/users/nfs_z/zn1/src/scaffHiC/src/scaffhic -nodes 54 -depth 50 -score 200 -map arima-AJ.map -plot arima-AJ.png -length`1107 `500000 -file 0 -fq1 /lustre/scratch117/sciops/team117/hpag/zn1/project/HiC/arima/human/QC/GM24385.AJ.R1.fastq.gz -fq2`1108 `/lustre/scratch117/sciops/team117/hpag/zn1/project/HiC/arima/human/QC/GM24385.AJ.R2.fastq.gz GRCH38.fasta aj-`1109 `arima.fasta > try.out`

1110

1111 Here we obtained arima-AJ.map and arima-AJ.png

1112 You may use PretextView to view the Hi-C map:

1113 <https://github.com/wtsi-hpag/PretextView>

1114

1115 Genome assembly

1116

1117 Contigs

1118 `~/zn1/src/hifiasm/hifiasm -o hg002-hifiasm -t 80 HG002-HiFi-all.fastq.gz > try.out`1119 `egrep "^S" hg002-hifiasm.p_ctg.gfa | awk '{print ">"$2"\n"$3}' > hg002-hifiasm.fasta`

1120

1121 Purge\_dups

1122 `/nfs/users/nfs_z/zn1/src/minimap2/minimap2-2.17_x64-linux/minimap2 -t 30 -xmap-pb hg002-hifiasm.fasta HG002-HiFi-`1123 `all.fastq.gz | gzip -c - > align.paf.gz`1124 `/nfs/users/nfs_z/zn1/src/purge_dups/bin/pbcstat align.paf.gz`1125 `/nfs/users/nfs_z/zn1/src/purge_dups/bin/calcuts PB.stat > cutoffs`1126 `/nfs/users/nfs_z/zn1/src/purge_dups/bin/split_fa hg002-hifiasm.fasta > Human.split`1127 `/nfs/users/nfs_z/zn1/src/minimap2/minimap2-2.17_x64-linux/minimap2 -t 20 -xasm5 -DP Human.split Human.split | gzip -c - >`1128 `split.self.paf.gz`1129 `/nfs/users/nfs_z/zn1/src/purge_dups/bin/purge_dups -2 -T cutoffs -c PB.base.cov split.self.paf.gz > dups.bed`1130 `/nfs/users/nfs_z/zn1/src/purge_dups/bin/get_seqs dups.bed hg002-hifiasm.fasta > purged.fa 2> hap.fa`

1131

1132 Scaffolding

1133 `/nfs/users/nfs_z/zn1/src/scaffHiC/src/scaff-bin/bwa-mem2 mem -t 54 -5SPM purged.fa`1134 `/lustre/scratch117/sciops/team117/hpag/zn1/project/HiC/arima/human/QC/GM24385.AJ.R1.fastq.gz`1135 `/lustre/scratch117/sciops/team117/hpag/zn1/project/HiC/arima/human/QC/GM24385.AJ.R2.fastq.gz > align-purge.sam`1136 `samtools view -@ 50 -bS align-purge.sam > Sorted_names.bam`

```
1137 samtools fixmate -@ 50 -m Sorted_names.bam Fixmate.bam > try.out
1138 samtools sort -@ 50 -o Sorted.bam Fixmate.bam > try.out
1139 rm -rf align-AJ.sam Sorted_names.bam Fixmate.bam
1140 samtools markdup -@ 50 -r -s Sorted.bam Dupmarked.bam > try.out
1141 mv Dupmarked.bam AJ-scaff.bam
1142
1143 ~zn1/src/yahs/yahs -o HG002-yahs.fa purged.fa AJ-scaff.bam > try.out
1144
1145 Hi-C map for scaffolded assembly
1146 /nfs/users/nfs_z/zn1/src/scaffHiC/src/scaffhic -nodes 54 -depth 50 -score 200 -map yahs-final-AJ.map -plot yahs-final-AJ.png -
1147 length 500000 -file 0 -fq1
1148 /lustre/scratch117/sciops/team117/hpag/zn1/project/HiC/arima/human/QC/jz1/GM24385.AJ.R1.fastq.gz -fq2
1149 /lustre/scratch117/sciops/team117/hpag/zn1/project/HiC/arima/human/QC/jz1/GM24385.AJ.R2.fastq.gz HG002-yahs.fa arima-
1150 AJ.fasta > try.out
1151
1152 Here we have yahs-final-AJ.map and yahs-final-AJ.png.
1153
```