

GigaScience

Multifaceted Hi-C benchmarking: what makes a difference in chromosome-scale genome scaffolding? --Manuscript Draft--

Manuscript Number:	GIGA-D-19-00211	
Full Title:	Multifaceted Hi-C benchmarking: what makes a difference in chromosome-scale genome scaffolding?	
Article Type:	Research	
Funding Information:	RIKEN	Dr. Ichiro Hiratani Dr. Shigehiro Kuraku
	Ministry of Education, Culture, Sports, Science and Technology (18H05530)	Dr. Ichiro Hiratani
Abstract:	<p>Background: Hi-C, a derivative of chromosome conformation capture (3C) targeting the whole genome, was originally developed as a means for characterizing chromatin conformation. More recently, this method has also been frequently employed in elongating nucleotide sequences obtained by de novo genome sequencing and assembly, in which the number of resultant sequences rarely converge into the chromosome number. Despite the prevailing and irreplaceable use, sample preparation methods for Hi-C have not been intensively discussed, especially from the standpoint of genome scaffolding.</p> <p>Results: To gain insights into the best practice of Hi-C scaffolding, we performed a multifaceted methodological comparison using vertebrate samples and optimized various factors during sample preparation, sequencing, and computation. As a result, we have identified some key factors that help improve Hi-C scaffolding including the choice and preparation of tissues, library preparation conditions, and restriction enzyme(s), as well as the choice of scaffolding program and its usage.</p> <p>Conclusions: This study provides the first comparison of multiple sample preparation kits/protocols and computational programs for Hi-C scaffolding, by an academic third party. We introduce a customized protocol designated the 'inexpensive and controllable Hi-C (iconHi-C) protocol', in which the optimal conditions revealed by this study have been incorporated, and release the resultant chromosome-scale genome assembly of the Chinese softshell turtle <i>Pelodiscus sinensis</i>.</p>	
Corresponding Author:	Shigehiro Kuraku JAPAN	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Mitsutaka Kadota	
First Author Secondary Information:		
Order of Authors:	Mitsutaka Kadota	
	Osamu Nishimura	
	Hisashi Miura	
	Kaori Tanaka	
	Ichiro Hiratani	
	Shigehiro Kuraku	
Order of Authors Secondary Information:		

Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p>	Yes

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)?

1 **Multifaceted Hi-C benchmarking: what makes a difference in**
2 **chromosome-scale genome scaffolding?**

3

4 Mitsutaka Kadota^{1*}, Osamu Nishimura^{1*}, Hisashi Miura², Kaori Tanaka^{1,3}, Ichiro
5 Hiratani², and Shigehiro Kuraku¹

6

7 ¹Laboratory for Phyloinformatics, RIKEN Center for Biosystems Dynamics Research
8 (BDR), Kobe, 650-0047, Japan, ²Laboratory for Developmental Epigenetics, RIKEN
9 BDR, Kobe, 650-0047, Japan, ³Present address: Division of Transcriptomics, Medical
10 Institute of Bioregulation, Kyushu University, Fukuoka, 812-0054, Japan

11

12 *These authors contributed equally to this study.

13

14 Correspondence address. Shigehiro Kuraku, Laboratory for Phyloinformatics, RIKEN
15 BDR, Japan. Tel: +81 78 306 3048; Fax: +81 78 306 3048; E-mail:
16 shigehiro.kuraku@riken.jp

17

18

19 **Abstract**

20 **Background:** Hi-C, a derivative of chromosome conformation capture (3C) targeting
21 the whole genome, was originally developed as a means for characterizing chromatin
22 conformation. More recently, this method has also been frequently employed in
23 elongating nucleotide sequences obtained by *de novo* genome sequencing and assembly,
24 in which the number of resultant sequences rarely converge into the chromosome
25 number. Despite the prevailing and irreplaceable use, sample preparation methods for
26 Hi-C have not been intensively discussed, especially from the standpoint of genome
27 scaffolding.

28 **Results:** To gain insights into the best practice of Hi-C scaffolding, we performed a
29 multifaceted methodological comparison using vertebrate samples and optimized
30 various factors during sample preparation, sequencing, and computation. As a result, we
31 have identified some key factors that help improve Hi-C scaffolding including the
32 choice and preparation of tissues, library preparation conditions, and restriction
33 enzyme(s), as well as the choice of scaffolding program and its usage.

34 **Conclusions:** This study provides the first comparison of multiple sample preparation
35 kits/protocols and computational programs for Hi-C scaffolding, by an academic third
36 party. We introduce a customized protocol designated the ‘inexpensive and controllable
37 Hi-C (iconHi-C) protocol’, in which the optimal conditions revealed by this study have
38 been incorporated, and release the resultant chromosome-scale genome assembly of the
39 Chinese softshell turtle *Pelodiscus sinensis*.

40

41 **Keywords:** Hi-C, genome scaffolding, chromosomes, proximity-guided assembly,
42 softshell turtle

43

44 **Background**

45 Chromatin, a complex of nucleic acids (DNA and RNA) and proteins, exhibits a
46 complex three-dimensional organization in the nucleus, which enables intricate
47 regulation of genome information expression through spatiotemporal controls (reviewed
48 in [1]). In order to characterize chromatin conformation on a genomic scale, the Hi-C
49 method was introduced as a derivative of chromosome conformation capture (3C) (Fig.
50 1A; [2]). This method detects chromatin contacts on a genomic scale through digestion
51 of crosslinked DNA molecules with restriction enzymes, followed by proximity ligation
52 of the digested DNA molecules. Massively parallel sequencing of the library harboring
53 ligated DNA molecules enables comprehensive quantification of contacts between
54 different genomic regions inside and between chromosomes, which is presented in a
55 heatmap conventionally called the ‘contact map’ [3].

56 Analyses of chromatin conformation with Hi-C have revealed more frequent
57 contacts between more closely linked genomic regions, which has prompted this
58 method to be employed in elongating *de novo* genome sequences, more recently [4]. In
59 *de novo* genome sequencing, the number of assembled sequences is usually far larger
60 than the number of chromosomes in the karyotype of the species of interest, irrespective
61 of the sequencing platform chosen [5]. The application of Hi-C scaffolding enabled
62 remarkable enhancement of sequence continuity to reach a chromosome scale and
63 integration of fragmentary sequences into longer sequences, which are similar in
64 number to that of chromosomes in the karyotype. In early 2018, commercial Hi-C
65 library preparation kits were introduced to the market (Fig. 1B), and *de novo* genome
66 assembly was revolutionized by the release of versatile computational programs for Hi-
67 C scaffolding (Table 1), namely LACHESIS [6], HiRise [7], SALSA [8, 9], and 3d-dna

68 [10]. These movements assisted the rise of mass sequencing projects targeting a number
69 of species, such as Earth BioGenome Project (EBP) [11], Genome 10K
70 (G10K)/Vertebrate Genome Project (VGP) [12, 13], and DNA Zoo Project [14].
71 Optimization of Hi-C sample preparation, however, has been limitedly attempted [15].
72 Thus, it remains unexplored which factor in particular makes a difference in the results
73 of Hi-C scaffolding, mainly because of its costly and resource-demanding nature.

74 Together with performing protocol optimization using human culture cells, we
75 focused on the softshell turtle *Pelodiscus sinensis* (Fig. 2). This species has been
76 adopted as a study system for evolutionary developmental biology (Evo-Devo),
77 including the study on the formation of the dorsal shell (carapace) (reviewed in [16]). It
78 is anticipated that relevant research communities have access to genome sequences of
79 optimal quality. In Japan, live materials (adults and embryos) of this species are
80 available through local farms mainly between May and August, which allows its high
81 utility for sustainable research. Based on a previous cytogenetic report, the karyotype of
82 this species consists of 33 chromosome pairs including Z and W ($2n = 66$) that show a
83 wide variety of sizes (conventionally categorized into macrochromosomes and
84 microchromosomes) [17]. Despite its moderate global GC-content in its whole genome
85 at around 44%, an earlier study suggested the intragenomic heterogeneity of GC-content
86 between and within the chromosomes, along with their sizes [18]. A wealth of
87 cytogenetic efforts on this species accumulated fluorescence *in situ* hybridization
88 (FISH)-based mapping data for 162 protein-coding genes covering almost all
89 chromosomes [17-19], which serves as structural landmarks for validating genome
90 assembly sequences.

91 A draft sequence assembly of the softshell turtle genome was built with short

92 reads and released already in 2013 [20]. This sequence assembly achieved the N50
93 scaffold length of >3.3 Mb but remains fragmented into approximately 20,000
94 sequences (see Supplementary Table S1). The longest sequence in this assembly is only
95 slightly larger than 16 Mb, which is much shorter than the largest chromosome size
96 estimated from the karyotype report [17]. The total size of the assembly is
97 approximately 2.2 Gb, which is a moderate size for a vertebrate species. Because of its
98 affordable genome size, sufficiently complex structure, and availability of validation
99 methods, we reasoned that the genome of this species is a suitable target for our
100 methodological comparison, and its improved genome assembly is expected to assist a
101 wide range of genome-based studies employing this species.

102

103

104 **Results**

105

106 **Stepwise QC before large-scale sequencing**

107 It would be ideal to judge the quality of prepared libraries before costly sequencing.

108 Following existing literature [15, 21], we routinely control the quality of Hi-C DNAs

109 and Hi-C libraries by observing DNA size shifts with digestion targeting the restriction

110 sites in properly prepared samples (Fig. 3). More concretely, a successfully ligated Hi-C

111 DNA sample should exhibit a slight length recovery of restricted DNA fragments after

112 ligation (QC1), which serves as an indicator of qualified samples (e.g., Sample 1 in Fig.

113 3B). In contrast, an unsuccessfully prepared Hi-C DNA does not exhibit this length

114 recovery (e.g., Sample 2 in Fig. 3B). In a later step, DNA molecules in a successfully

115 prepared HindIII-digested Hi-C library should contain the NheI restriction site at a high

116 probability. Thus, the length distribution after the NheI digestion of the prepared library
117 serves as an indicator of qualified or disqualified products (QC2; Fig. 3C). This series
118 of QCs is incorporated into our protocol by default (Supplementary Protocol S1) and
119 can also be performed along with sample preparation using commercial kits provided
120 that it employs a single restriction enzyme.

121 Some of the libraries we have prepared passed the QC steps before sequencing
122 but yielded an unpreferably large proportion of unusable read pairs. To identify such
123 libraries, we routinely performed small-scale sequencing with the purpose of quick and
124 inexpensive QC using the HiC-Pro program [22] (see Fig. 4 for the read pair categories
125 assigned by HiC-Pro). Our test with variable input data sizes (500 K–200 M read pairs)
126 resulted in highly similar breakdowns into different categories of read pair properties
127 (Supplementary Table S2) and guaranteed the QC with an extremely small data size of 1
128 M or fewer reads. These post-sequencing QC steps that do not incur a large cost are
129 expected to help avoid large-scale sequencing of unsuccessful libraries that have
130 somehow passed through QC1 and QC2 steps. Importantly, libraries that have passed
131 this QC can be further sequenced in more depth as necessary.

132

133 **Optimization of sample preparation conditions**

134 We identified overt differences between sample preparation protocols of already
135 published studies and those of commercial kits (Fig. 1B). Therefore, we first sought to
136 optimize the conditions of several preparation steps using human culture cells.

137 To evaluate the effect of the degree of cell fixation, we prepared Hi-C libraries
138 from GM12878 cells fixed for 10 and 30 minutes. Our comparison did not detect any
139 marked difference in the quality of Hi-C DNA (QC1; Fig. 5A) and Hi-C library (QC2;

140 Fig. 5B). However, libraries with longer fixation showed larger proportions of dangling
141 end read pairs and re-ligation read pairs, as well as a smaller proportion of valid
142 interaction reads (Fig. 5C). Increased duration of cell fixation reduces the proportion of
143 long-range (>1 Mb) interactions among the overall captured interactions (Fig. 5D).

144 The reduced preparation time with commercial Hi-C kits (up to two days
145 according to their advertisement) is attributable mainly to shortened duration of
146 restriction and ligation (Fig. 1B). To monitor the effect of shortening these enzymatic
147 reactions, we analyzed the progression of restriction and ligation in a time course
148 experiment using human GM12878 cells. The results show persistent progression of
149 restriction until 16 hours and of ligation until 6 hours (Fig. 6).

150

151 **Multifaceted comparison using softshell turtle samples**

152 On the basis of the detailed optimization of sample preparation conditions described
153 above, we built an original protocol, designated the ‘iconHi-C protocol’, with 10 min-
154 long cell fixation, 16 hour-long restriction, 6 hour-long ligation, and successive QC
155 steps (Methods; also see Supplementary Protocol S1; Fig. 1B).

156 We performed Hi-C sample preparation and scaffolding using tissues from a
157 female Chinese softshell turtle which is known to have both Z and W chromosomes
158 [17]. For this purpose, we prepared Hi-C libraries with variable tissues (liver or blood
159 cells), restriction enzymes (HindIII or DpnII), and protocols (our iconHi-C protocol, the
160 Arima Genomics kit in conjunction with the KAPA Hyper Prep Kit, or the Phase
161 Genomics kit) as outlined in Fig. 7A (see Supplementary Table S3; Supplementary Fig.
162 S1). As in some existing protocols (e.g., [23]), we performed T4 DNA polymerase
163 treatment in our iconHi-C protocol (Library a–d), expecting reduced proportions of

164 ‘dangling end’ read pairs that contain no ligated junction and thus do not contribute to
165 Hi-C scaffolding. We also incorporated this T4 DNA polymerase treatment in the
166 workflow of the Arima kit (Library e vs. Library f without this additional treatment).
167 We also tested a lesser degree of PCR amplification (11 cycles) along with the use of
168 the Phase Genomics kit which compels as many as 15 cycles by default (Library h vs.
169 Library g; Fig. 7A).

170 The samples prepared with the iconHi-C protocol, which is compatible with the
171 abovementioned QC1 and QC2, were all judged as qualified, by these QCs (Fig. 7B).
172 The prepared Hi-C libraries were sequenced to obtain one million 127nt-long read pairs
173 and subjected to post-sequencing QC with the HiC-Pro program (Fig. 8). As a result of
174 this QC, the largest proportion of ‘valid interaction’ pairs was observed for Arima
175 libraries (Library e and f). As for the iconHi-C libraries (Library a–d), fewer
176 ‘unmapped’ and ‘religation’ pairs were detected with the DpnII libraries than with
177 HindIII libraries. It should be noted that the QC results for the softshell turtle libraries
178 generally produced lower proportions of the ‘valid interaction’ category and larger
179 proportions of ‘unmapped pairs’ and ‘pairs with singleton’ than those for human
180 libraries. This cross-species difference is accounted for by possibly incomplete genome
181 sequences used as a reference for Hi-C read mapping (Supplementary Table S1). This
182 evokes a caution in comparing QC results across species.

183

184 **Scaffolding with variable inputs and computational conditions**

185 In this study, only well-maintained, open-source programs, namely 3d-dna and
186 SALSA2, were used in conjunction with variable combinations of an input library, an
187 input read amount, an input sequence cutoff length, and a number of iterative misjoin

188 correction rounds (Fig. 9A). As a result of scaffolding, we observed a wide spectrum of
189 basic metrics, including the N50 scaffold length (0.6–303 Mb), the largest scaffold
190 length (8.7–703 Mb), and the number of chromosome-sized (>10 Mb) sequences (0–65)
191 (Fig. 9; Supplementary Table S4).

192 First of all, with the default parameters, 3d-dna consistently produced more
193 continuous assemblies than SALSA2 (see Assembly 1 vs. 5, 3 vs. 6, 9 vs. 10, and 11 vs
194 12 in Fig. 9). Second, increasing the number of iterative corrections ('-r' option with 3d-
195 dna) resulted in relatively large N50 lengths but with more missing orthologs (see
196 Assembly 13–15). Third, a smaller input sequence cutoff length ('-i' option with 3d-
197 dna) resulted in a smaller number of resultant scaffolds but again, with more missing
198 orthologs (see Assembly 13, 16–18). Fourth, using the liver libraries consistently
199 resulted in a higher continuity than using the blood cell libraries (see Assembly 1 vs. 2
200 as well as 3 vs. 4 in Fig. 9).

201 Of those, Assembly 8, employing input Hi-C reads derived from both liver and
202 blood, exhibited an outstandingly large N50 scaffold length (303 Mb) but a larger
203 number of undetected reference ortholog (141 orthologs) than most of the other
204 assemblies. The largest scaffold (scaffold 5) in this assembly is approximately 703 Mb
205 long, causing the large N50 length, and accounts for approximately one-third of the
206 whole genome in length, as a result of possible overassembly bridging 14 putative
207 chromosomes (see Supplementary Fig. S2).

208 The choice of restriction enzymes has not yet been discussed in depth, in the
209 context of genome scaffolding. In the present study, we separately prepared Hi-C
210 libraries with HindIII and DpnII. We did not mix multiple enzymes in a reaction (apart
211 from using the Arima kit originally employing two enzymes) and instead performed a

212 single scaffolding run with both HindIII-based and DpnII-based reads (see Assembly 7
213 in Fig. 9). Our comparison of multiple metrics expectedly highlights a more successful
214 result with DpnII than with HindIII (see Assembly 1 vs. 3 as well as 2 vs. 4; Fig. 9).
215 However, the mixed input of HindIII-based and DpnII-based reads did not necessarily
216 yield a better scaffolding result (see Assembly 3 vs. 7).

217

218 **Validation of scaffolding results with transcriptome and FISH data**

219 In addition to the above-mentioned evaluation of the scaffolding results based on
220 sequence length and gene space completeness, we attempted to evaluate the sequence
221 continuity with independently obtained data. First, we mapped assembled transcript
222 sequences onto our Hi-C scaffold sequences (see Methods). This did not reveal any
223 substantial differences between the assemblies (Supplementary Table S5), probably
224 because the sequence continuity after Hi-C scaffolding already exceeded that of RNA-
225 seq library inserts even when the lengths of intervening introns in the genome are taken
226 into consideration. The present analysis with RNA-seq data did not provide an effective
227 resort of continuity validation.

228 Second, we referred to the fluorescence *in situ* hybridization (FISH) mapping
229 data for 162 protein-coding genes from published cytogenetic studies [17-19], which
230 allowed us to check the locations of those genes with our resultant Hi-C assemblies. In
231 this analysis, we evaluated Assembly 3, 7, and 9 (see Fig. 9A) that showed better
232 scaffolding results in terms of sequence length distribution and gene space completeness
233 (Fig. 9B). As a result, we confirmed the positioning of almost all genes and their
234 continuity over the centromeres, which encompassed not only large but also small
235 chromosomes (conventionally called ‘macro-’ and ‘micro-chromosomes’; Fig. 10). Two

236 genes that were not confirmed by Assembly 7 (*UCHL1* and *COX15*; Fig. 10) were
237 found in separate scaffold sequences shorter than 1 Mb, which indicates insufficient
238 scaffolding. On the other hand, the gene array including *RBM5*, *TKT*, *WNT7A*, and
239 *WNT5A*, previously shown by FISH, was consistently unconfirmed by all the three
240 assemblies (Fig. 10), which did not provide any clue for among-assembly evaluation or
241 even indicated an erroneous interpretation of FISH data in a previous study.

242

243

244 **Discussion**

245

246 **Starting materials: not genomic DNA extraction but *in situ* cell fixation**

247 In genome sequencing, best practices for high molecular weight DNA extraction have
248 often been discussed (e.g., [24]). This factor is fundamental to building longer contigs,
249 whether employing short-read or long-read sequencing platforms. Also, the proximity
250 ligation method using Chicago libraries provided by Dovetail Genomics which is based
251 on *in vitro* chromatin reconstruction [7], uses genomic DNA as starting materials.
252 Instead, proximity guided assembly enabled by Hi-C employs cellular nuclei preserving
253 chromatin conformation, which brings a new technical challenge for appropriate
254 sampling and sample preservation in genomics.

255 In preparing the starting materials, it seems important to optimize the degree of
256 cell fixation depending on your sample choice, to obtain an optimal result in Hi-C
257 scaffolding (Fig. 5). Another practical lesson about tissue choice was obtained by
258 examining Assembly 8 (Fig. 9A). This assembly was produced by 3d-dna scaffolding
259 with both liver and blood libraries (Library b and d), which led to an unacceptable result

260 possibly caused by overassembly (Fig. 9B–D; also see Results). It is likely that
261 enhanced cellular heterogeneity, possibly introducing excessive conflicting chromatin
262 contacts, did not allow the scaffolding program to properly group and order the input
263 genome sequences. In brief, we recommend the use of samples with modest cell-type
264 heterogeneity amenable to thorough fixation.

265

266 **Considerations in sample preparation**

267 In this study, we could not test all commercial Hi-C kits available in the market. This is
268 partly because the Dovetail Hi-C kit specifies a non-open source program HiRise as the
269 only supported downstream computation solution and does not allow a direct
270 comparison with other kits, namely those from Phase Genomics and Arima Genomics.

271 According to our calculation, it would be at least three times more economical to
272 prepare a Hi-C library with the iconHi-C protocol than with a commercial kit.

273 Practically, the cost difference would be even larger, either when one cannot fully
274 consume the purchased kit or when one cannot undertake post-sequencing computation
275 steps and thus cover additional outsourcing cost for this.

276 Genomic regions targeted by Hi-C are determined by the choice of restriction
277 enzymes. Theoretically, 4-base cutters (e.g., DpnII), potentially with more frequent
278 restriction sites on the genome, are expected to provide a higher resolution than 6-base
279 cutters (e.g., HindIII) [15]. However, it might not be so straightforward when the
280 species-by-species variation of GC-content, as well as its intra-genomic heterogeneity,
281 are taken into consideration. The use of multiple enzymes in a single reaction could be
282 promising, but not all scaffolding programs are compatible with multiple enzymes from
283 a computational viewpoint (see Table 1 for a comparison of scaffolding program

284 specifications). Another technical downside is the incompatibility of DNA ends
285 restricted by multiple enzymes, with restriction-based QCs, such as QC2 in our iconHi-
286 C protocol (Fig. 3). Therefore, in this study, DpnII and HindIII were separately
287 employed in conjunction with the iconHi-C protocol, which resulted in higher
288 scaffolding performance with the DpnII library (Figs. 8 and 9), as expected. In addition,
289 we input the separately prepared DpnII and HindIII libraries together in scaffolding
290 (Assembly 7), but this attempt did not lead to higher scaffolding performance (Figs.
291 9B–D and 10). The Arima Hi-C kit employs two different enzymes that can produce
292 much more combinations of restriction sites, because one of the two enzymes
293 recognizes the nucleotide stretch GATTC. Scaffolding with the libraries prepared using
294 this kit resulted in one of the most acceptable assemblies (Assembly 9). However, this
295 result did not explicitly exceed the performance of scaffolding with the iconHi-C
296 libraries including the one employing only a single enzyme DpnII (Library d).

297 One concern about the use of commercial kits (except the Arima Hi-C kit used
298 with the Arima-QC2) is overamplification by PCR, as their manuals specify certain
299 numbers of PCR cycles *a priori* (15 cycles for the Phase Genomics Proximo Hi-C kit
300 and 11 cycles for the Dovetail Hi-C kit). In our iconHi-C protocol, an optimal number
301 of PCR cycles is estimated by means of a preliminary real-time PCR using a small
302 aliquot (Step11.25–29 in Supplementary Protocol S1) as traditionally performed for
303 other library types (e.g., [25]). This procedure allowed us to minimize the PCR cycles
304 down to five cycles (Supplementary Table S3). The Dovetail Hi-C kit recommends that
305 one consumes larger amounts of kit components than specified for a single sample,
306 depending on the genome size, as well as the degree of genomic heterozygosity and
307 repetitiveness, of the species of interest. However, with our iconHi-C protocol, we

308 always performed a single library preparation, irrespective of those species-specific
309 factors, which we understand suffices in all the cases we have tested.

310 Commercial Hi-C kits, usually advertised for easiness and quickness, have
311 largely shortened the protocol down to two days, in comparison with existing non-
312 commercial protocols (e.g., [15, 23]). Such time-saving protocols are achieved mainly
313 by shortened durations of restriction enzyme digestion and ligation (Fig. 1B). Our
314 assessment, however, showed unsaturated reaction within such shortened time frames
315 employed in the commercial kits (Fig. 6). Also, our attempt to insert a step for T4 DNA
316 polymerase treatment in sample preparation with the Arima Hi-C kit resulted in reduced
317 ‘dangling end’ reads (Library e vs. Library f in Fig. 8). As for the Phase Genomics
318 Proximo Hi-C kit, transposase-based library preparation contributes largely to
319 shortening its protocol, but this decreases the operability of library insert lengths.
320 Especially if Hi-C sample preparation is performed for a limited number of samples, as
321 practiced typically for genome scaffolding, one would opt to consider these points, even
322 in using commercial kits, in order to further improve the quality of prepared libraries
323 and scaffolding products.

324

325 **Considerations in sequencing**

326 The quantity of Hi-C read pairs to be input for scaffolding is critical because it accounts
327 for the majority of the cost of Hi-C scaffolding. Our protocol introduces a thorough
328 safety system to prevent sequencing unsuccessful libraries, firstly with pre-sequencing
329 QCs for size shift analysis (Fig. 3) and secondly with small-scale (down to 500 K read
330 pairs) sequencing (see Results; also see Supplementary Table S2, S6).

331 Our comparison shows a dramatic decrease in assembly quality when less than

332 100 M read pairs were used (see the comparison among Assembly 19–23 above in Fig.
333 9). Still, we obtained optimal results with a smaller number of reads (ca. 160 M per 2.2
334 Gb genome) than recommended by commercial kits (e.g., 100 M per 1 Gb genome for
335 the Dovetail Hi-C kit and 200 M per Gb genome for the Arima Hi-C kit). As generally
336 and repeatedly discussed, the proportion of informative reads and their diversity, rather
337 than just the number of all obtained reads, are critical.

338 In terms of read length, we did not perform any comparison in this study.
339 Longer reads may enhance the fidelity in characterizing the read pair property and
340 allows precise QC. Still, the existing Illumina sequencing platform has enabled
341 economical acquisition of 150 nt-long paired-end reads, which did not prompt us to
342 vary the read length.

343

344 **Considerations in computation**

345 In this study, 3d-dna produced a more reliable scaffolding output than SALSA2,
346 whether sample preparation employed a single or multiple enzyme(s) (Fig. 9B–D). On
347 the other hand, 3d-dna needed more time to complete scaffolding than SALSA2. Apart
348 from the choice of the program, there are quite a few points to consider, in order to
349 achieve successful scaffolding for a smaller investment. In general, it is advised not to
350 take Hi-C scaffolding results for granted, and it is necessary to improve them by
351 referring to contact maps, using an interactive tool such as Juicebox [14]. In this study,
352 however, we compared raw scaffolding outputs to evaluate sample preparation and
353 reproducible computational steps.

354 Our study employed variable parameters of the scaffolding programs (Fig. 9A).
355 First, available Hi-C scaffolding programs have different default length cut-off values

356 for input sequences (e.g., 15000 bp for the parameter ‘-i’ with 3d-dna and 1000 bp for
357 the parameter ‘-c’ with SALSA2). Only sequences longer than the cut-off length value
358 contribute to sequence elongation towards the chromosome sizes, and those shorter than
359 that are implicitly excluded from the scaffolding process and remain unchanged.
360 Typically with the Illumina sequencing platform, genomic regions with unusually high
361 frequencies of GC-content and repetitive elements are not assembled into sequences
362 with sufficient lengths (see [26]). Such genomic regions tend to be excluded from
363 chromosome-scale Hi-C scaffolds because their length is smaller than the threshold. It is
364 also possible that such regions are excluded because few Hi-C read pairs are mapped to
365 such regions, even if they exceed the cutoff length. One needs to deliberately set the
366 length cutoff in accordance with the overall continuity of the input assembly and
367 possible interest into particular, fragmentary sequences expected to be elongated. It
368 should be warned that lowering the length threshold can result in frequent misjoins in
369 the scaffolding output (Fig. 9B–D) or too much computational time. Regarding the
370 number of iterative misjoin correction rounds (the parameter ‘-r’ with 3d-dna and ‘i’
371 with SALSA2), our attempts with increased values did not necessarily yield favorable
372 results (Fig. 9B–D), which did not provide a consistent optimal range of values but
373 rather suggests the importance of performing multiple scaffolding runs with varied
374 parameters.

375

376 **Considerations in assessing chromosome-scale genome sequences**

377 Our assessment with cytogenetic data confirmed the continuity of gene linkage over the
378 obtained chromosome-scale sequences (Fig. 10). This validation was necessitated by
379 almost saturated scores of typical gene space completeness assessment such as BUSCO

380 (Supplementary Table S4) as well as transcript contig mapping (Supplementary Table
381 S5), both of which did not provide an effective metric for evaluation.

382 For further evaluation of our scaffolding results, we referred to sequence length
383 distribution of the genome assemblies of other turtle species that are regarded as
384 chromosome-scale. This showed comparable values for the basic metrics to our Hi-C
385 scaffolding results on the softshell turtle, that is, a N50 length of 127.5 Mb and the
386 maximum sequence length of 344.5 Mb for the green sea turtle (*Chelonia mydas*)
387 genome assembly released by the DNA Zoo Project and a N50 length of 131.6 Mb and
388 the maximum length of 370.3 Mb for the Goode's thornscrub tortoise (*Gopherus*
389 *evgoodei*) genome assembly released by the Vertebrate Genome Project (VGP).
390 Scaffolding results should be evaluated by referring to an estimate N50 length and the
391 maximum length based on the actual number and the length distribution of
392 chromosomes in the intrinsic karyotype of the species in question or its close relative.
393 Turtles tend to have the N50 length of approximately 130 Mb and the maximum length
394 of 350 Mb, while many teleost fish genomes exhibit an N50 length of as low as 20–30
395 Mb and the maximum length of <100 Mb [27]. If these metrics show excessive values,
396 scaffolded sequences harbor overassembly that erroneously boosts length-based metrics.
397 Larger values that researchers conventionally regard as signs for successful sequence
398 assembly do not necessarily indicate higher precision.

399 The total length of assembly sequences is expected to increase after Hi-C
400 scaffolding, because scaffolding programs simply insert a stretch of the unassigned base
401 'N' with a uniform length between input sequences in most cases (500 bp as default
402 with both 3d-dna and SALSA2). However, this has a minor impact on the total
403 assembly sequence length. In fact, inserting the 'N' stretches of arbitrary lengths has

404 been an implicit, rampant practice even before Hi-C scaffolding prevailed—for
405 example, the most and second most frequent lengths of the ‘N’ stretch in the publicly
406 available zebrafish genome assembly Zv10 are 100 and 10 bp, respectively.

407

408 **Conclusions**

409 In this study, we introduced the iconHi-C protocol in which successive QC steps are
410 implemented, and assessed possible keys for improving Hi-C scaffolding. Overall, our
411 study shows that a small variation in sample preparation or computation for scaffolding
412 can have a large impact on scaffolding output, and any scaffolding output should ideally
413 be validated by independent information, such as cytogenetic data, long reads, or
414 genetic linkage maps. Our present study aimed to evaluate the output of reproducible
415 computational steps, which in practice should be followed by modifying the raw
416 scaffolding output by referring to independent information or by analyzing chromatin
417 contact maps. The study employed only limited combinations of species, sample prep
418 methods, scaffolding programs, and its parameters, and we will continue testing
419 different conditions for kits/programs that did not necessarily perform well here with
420 our specific materials.

421

422 **Methods**

423

424 **Initial genome assembly sequences**

425 The softshell turtle (*Pelodiscus sinensis*) assembly published previously [20] was
426 downloaded from NCBI GenBank (GCA_000230535.1), whose gene space
427 completeness and length statistics were assessed by gVolante [28] (see Supplementary

428 Table S1 for the assessment results). Although it could be suggested to remove
429 haplotigs before Hi-C scaffolding [29], we omitted this step because of the low
430 frequency of the reference orthologs with multiple copies (0.72 %; Supplementary
431 Table S1), indicating a minimal degree of haplotig contamination.

432

433 **Animals and cells**

434 We sampled tissues (liver and blood cells) from a female purchased from a local farmer
435 in Japan, because the previous whole genome sequencing used the whole blood of a
436 female [20]. All the experiments were conducted in accordance with the Guideline of
437 the Institutional Animal Care and Use Committee of RIKEN Kobe Branch (Approval
438 ID: A2017-12).

439 Human lymphoblastoid cell line GM12878 was purchased from the Coriell Cell
440 Repositories and cultured in RPMI-1640 media (Thermo Fisher Scientific)
441 supplemented with 15% FBS, 2 mM L-glutamine, and 1x antibiotic-antimycotic
442 solution (Thermo Fisher Scientific), at 37 °C, 5 % CO₂, as described previously [30].

443

444 **Hi-C sample preparation using the original protocol**

445 We have made modifications to a protocol introduced in previous literature [23, 31]
446 (Fig. 1B). The full version of the modified ‘inexpensive and controllable Hi-C (iconHi-
447 C)’ protocol is described in Supplementary Protocol S1.

448

449 **Hi-C sample preparation using commercial kits**

450 The Proximo Hi-C kit (Phase Genomics) which employs the restriction enzyme Sau3A1
451 and transposase-based library preparation [32] (Fig. 1B) was used for preparing a

452 library from the 50 mg softshell turtle liver following its official ver. 1.0 animal
453 protocol (Library g in Fig. 7A) and a library from the 10 mg liver amplified with a
454 reduced number of PCR cycles based on a preliminary real-time qPCR using an aliquot
455 (Library h; see [25] for the detail of the pre-determination of optimal PCR cycles). The
456 Arima Hi-C kit (Arima Genomics) which employs a restriction enzyme cocktail (Fig.
457 1B) was used in conjunction with the KAPA Hyper Prep Kit (KAPA Biosystems),
458 protocol ver. A160108 v00, to prepare a library using the softshell turtle liver, following
459 its official animal vertebrate tissue protocol (ver. A160107 v00) (Library f) and a library
460 with an additional step of T4 DNA polymerase treatment for reducing ‘dangling end’
461 reads (Library e). This additional treatment is detailed in Step 8.2 (for DpnII-digested
462 samples) in Supplementary Protocol S1.

463

464 **DNA sequencing**

465 Small-scale sequencing for library QC was performed in-house to obtain 127 nt-long
466 paired-end reads on an Illumina HiSeq 1500 in the Rapid Run Mode. Large-scale
467 sequencing for Hi-C scaffolding was performed to obtain 151 nt-long paired-end reads
468 on an Illumina HiSeq X. The obtained reads were subjected to quality control with
469 FastQC ver. 0.11.5 (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>), and
470 low-quality regions and adapter sequences in the reads were removed using Trim Galore
471 ver. 0.4.5 (https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/) with the
472 parameters ‘-e 0.1 -q 30’.

473

474 **Post-sequencing quality control of Hi-C libraries**

475 For post-sequencing library QC, one million trimmed read pairs for each Hi-C library

476 were sampled using the ‘subseq’ function of the program seqtk ver. 1.2-r94
477 (<https://github.com/lh3/seqtk>). The resultant sets of read pairs were processed using
478 HiC-Pro ver. 2.11.1 [22] with bowtie2 ver. 2.3.4.1 [33] to evaluate the insert structure
479 and mapping status onto the softshell turtle genome assembly PelSin_1.0
480 (GCF_000230535.1) or human genome assembly hg19. This resulted in the
481 categorization between valid interaction pairs and invalid pairs, and the latter is divided
482 into ‘dangling end’, ‘religation’, ‘self circle’, and ‘single-end’ (Fig. 4). To process the
483 read pairs derived from the libraries prepared using either HindIII or DpnII (Sau3AI)
484 with the iconHi-C protocol (Library a–d) and the Phase Genomics Proximo Hi-C kit
485 (Library g and h), the restriction fragment file required by HiC-Pro was prepared
486 according to the script ‘digest_genome.py’ provided with HiC-Pro. To process the reads
487 derived from the Arima Hi-C kit (Library e and f), all restriction sites (‘GATC’ and
488 ‘GANTC’) were inserted into the script. In addition, the nucleotide sequences of all
489 possible ligated sites generated by restriction enzymes were included in a configuration
490 file of HiC-Pro. The details and the sample code are included in Supplementary
491 Protocol S2.

492

493 **Computation for Hi-C scaffolding**

494 In order to control our comparison with intended input data sizes, certain numbers of
495 trimmed read pairs were sampled for each library with seqtk as described above.

496 Scaffolding was processed with the following methods employing two program
497 pipelines, 3d-dna and SALSA2.

498 Scaffolding with the program 3d-dna was preceded by Hi-C read mapping onto
499 the genome with Juicer ver. 20180805 [34] using the default parameters with BWA

500 ver.0.7.17-r1188 [35]. The restriction fragment file required by Juicer was prepared by
501 the script ‘generate_site_positions.py’ provided with Juicer or our original script
502 compatible with multiple restriction enzymes to convert the restriction fragment file of
503 HiC-Pro to the format required by Juicer (Supplementary Protocol S2). Scaffolding with
504 3d-dna ver. 20180929 was performed with variable parameters (see Fig. 9A).

505 Scaffolding with the program SALSA2 using Hi-C reads was preceded by Hi-C
506 read pair processing with the Arima mapping pipeline ver. 20181207
507 (https://github.com/ArimaGenomics/mapping_pipeline) together with BWA, SAMtools
508 ver. 1.8-21-gf6f50ac [36] and Picard ver. 2.18.12
509 (<https://github.com/broadinstitute/picard>). The mapping result in the binary alignment
510 map (bam) format was converted into a BED file by bamToBed of Bedtools ver. 2.26.0
511 [37], whose output was used as an input of scaffolding using SALSA2 ver. 20181212
512 with the default parameters.

513

514 **Completeness assessment of Hi-C scaffolds**

515 gVolante ver. 1.2.1 [28] was used to perform an assessment of sequence length
516 distribution and gene space completeness based on the coverage of one-to-one reference
517 orthologs with BUSCO v2/v3 employing the one-to-one ortholog set ‘Tetrapoda’
518 supplied with BUSCO [38]. For the assessment, no threshold of cut-off length was set.

519

520 **Continuity assessment with RNA-seq read mapping**

521 Paired-end reads obtained by RNA-seq of softshell turtle embryos at multiple stages
522 were downloaded from NCBI SRA (DRX001576) and were assembled with the
523 program Trinity ver. 2.7.0 [39] with the default parameters. The assembled transcript

524 sequences were mapped with pblat [40] to the Hi-C scaffold sequences, and the output
525 was assessed with isoblat ver. 0.31 [41].

526

527 **Comparison with chromosome FISH results**

528 Cytogenetic validation of Hi-C scaffolding results was performed by comparing the
529 gene locations on the scaffold sequences with those in preexisting chromosome FISH
530 data for 162 protein-coding genes [17-19]. The nucleotide exonic sequences for those
531 162 genes retrieved from GenBank were aligned with Hi-C scaffold sequences using
532 BLAT ver. 36x2 [42], and their positions and orientation along the Hi-C scaffold
533 sequences were analyzed.

534

535 **Availability of supporting data**

536 All sequence data generated from this study have been submitted to the DDBJ Sequence
537 Read Archive (DRA) under accession IDs DRA008313. The datasets supporting the
538 results of this article are available in the FigShare
539 (<https://figshare.com/s/6ea495a65fc231a74458>).

540

541 **Additional files**

542 Supplementary Figure S1. Quality control of the Hi-C libraries.

543

544 Supplementary Figure S2. Structural analysis of the possibly overassembled scaffold in
545 Assembly #8

546

547 Supplementary Figure S3. Results of quality controls before sequencing.

548

549 Supplementary Table S1. Statistics of Chinese softshell turtle draft genome assembly

550 before Hi-C.

551

552 Supplementary Table S2. HiC-Pro results of the human GM12878 HindIII Hi-C library

553 with reduced reads

554

555 Supplementary Table S3. Quality control of Hi-C libraries.

556

557 Supplementary Table S4. Scaffolding results with variable input data and computational

558 parameters

559

560 Supplementary Table S5. Mapping results of assembled transcript sequences onto Hi-C

561 scaffolds

562

563 Supplementary Table S6. HiC-Pro results of the softshell turtle liver DpnII library

564 (Library d) with reduced reads

565

566 Supplementary Table S7. Quality control of the human GM12878 Hi-C libraries

567

568 Supplementary Protocol S1. Protocol of iconHi-C

569

570 Supplementary Protocol S2. Computational protocol to support multiple enzymes

571

572

573

574 **Abbreviations**

575 PCR: polymerase chain reaction; FISH, fluorescence *in situ* hybridization; BUSCO,

576 benchmarking universal single-copy orthologs; NCBI, National Center for

577 Biotechnology Information; NGS, next generation DNA sequencing

578

579 **Funding**

580 This work was supported by intramural grants within RIKEN to S.K. and I.H. and by a

581 Grant-in-Aid for Scientific Research on Innovative Areas to I.H. (18H05530) from the

582 Ministry of Education, Culture, Sports, Science, and Technology (MEXT).

583

584 **Competing interests**

585 The authors declare that they have no competing interests

586

587 **Acknowledgements**

588 The authors acknowledge Naoki Irie, Juan Pascual Anaya and Tatsuya Hirasawa in

589 Laboratory for Evolutionary Morphology, RIKEN BDR for suggestions for sampling,

590 Rawin Poonperm for comments and discussion on the iconHi-C protocol, Olga

591 Dudchenko, Erez Lieberman-Aiden, Arang Rhie, Sergey Koren, and Jay Ghurye for

592 their technical suggestions for sample preparation and computation, Yoshinobu Uno for

593 guidance to cytogenetic data interpretation, and Anthony Schmitt of Arima Genomics

594 and Stephen Eacker of Phase Genomics for providing information about the Hi-C kits.

595 They also thank the other members of Laboratory for Phyloinformatics and Laboratory

596 for Developmental Epigenetics in RIKEN BDR for technical support and discussion.

597

598 **Author contributions**

599 S.K., I.H., H.M., and M.K. conceived the study. M.K. and K.T. performed laboratory
600 works, and O.N. performed bioinformatic analysis. M.K., O.N., and H.M. analyzed the
601 data. S.K., M.K., and O.N. drafted the manuscript. All authors contributed to the
602 finalization of the manuscript.

603

604 **References**

- 605 1. Rowley MJ and Corces VG. Organizational principles of 3D genome
606 architecture. *Nature Reviews Genetics*. 2018;19 12:789-800.
607 doi:10.1038/s41576-018-0060-8.
- 608 2. Lieberman-Aiden E, van Berkum NL, Williams L, Imakaev M, Ragozy T,
609 Telling A, et al. Comprehensive Mapping of Long-Range Interactions Reveals
610 Folding Principles of the Human Genome. *Science*. 2009;326 5950:289-93.
611 doi:10.1126/science.1181369.
- 612 3. Rao Suhas SP, Huntley Miriam H, Durand Neva C, Stamenova Elena K,
613 Bochkov Ivan D, Robinson James T, et al. A 3D Map of the Human Genome at
614 Kilobase Resolution Reveals Principles of Chromatin Looping. *Cell*. 2014;159
615 7:1665-80. doi:10.1016/j.cell.2014.11.021.
- 616 4. Burton JN, Adey A, Patwardhan RP, Qiu R, Kitzman JO and Shendure J.
617 Chromosome-scale scaffolding of de novo genome assemblies based on
618 chromatin interactions. *Nature Biotechnology*. 2013;31:1119.
619 doi:10.1038/nbt.2727.

- 620 5. Sedlazeck FJ, Lee H, Darby CA and Schatz MC. Piercing the dark matter:
621 bioinformatics of long-range sequencing and mapping. *Nature Reviews*
622 *Genetics*. 2018;19 6:329-46. doi:10.1038/s41576-018-0003-4.
- 623 6. Bickhart DM, Rosen BD, Koren S, Sayre BL, Hastie AR, Chan S, et al. Single-
624 molecule sequencing and chromatin conformation capture enable de novo
625 reference assembly of the domestic goat genome. *Nature Genetics*. 2017;49:643.
626 doi:10.1038/ng.3802.
- 627 7. Putnam NH, O'Connell BL, Stites JC, Rice BJ, Blanchette M, Calef R, et al.
628 Chromosome-scale shotgun assembly using an in vitro method for long-range
629 linkage. *Genome Research*. 2016; doi:10.1101/gr.193474.115.
- 630 8. Ghurye J, Pop M, Koren S, Bickhart D and Chin C-S. Scaffolding of long read
631 assemblies using long range contact information. *BMC Genomics*. 2017;18
632 1:527. doi:10.1186/s12864-017-3879-z.
- 633 9. Ghurye J, Rhie A, Walenz BP, Schmitt A, Selvaraj S, Pop M, et al. Integrating
634 Hi-C links with assembly graphs for chromosome-scale assembly. *bioRxiv*.
635 2018:261149. doi:10.1101/261149.
- 636 10. Dudchenko O, Batra SS, Omer AD, Nyquist SK, Hoeger M, Durand NC, et al.
637 De novo assembly of the *Aedes aegypti* genome using Hi-C yields chromosome-
638 length scaffolds. *Science*. 2017;356 6333:92-5. doi:10.1126/science.aal3327.
- 639 11. Lewin HA, Robinson GE, Kress WJ, Baker WJ, Coddington J, Crandall KA, et
640 al. Earth BioGenome Project: Sequencing life for the future of life. *Proceedings*
641 *of the National Academy of Sciences of the United States of America*. 2018;115
642 17:4325-33. doi:10.1073/pnas.1720115115.
- 643 12. Koepfli KP, Paten B and O'Brien SJ. The Genome 10K Project: a way forward.

- 644 Annual review of animal biosciences. 2015;3:57-111. doi:10.1146/annurev-
645 animal-090414-014900.
- 646 13. Editorial. A reference standard for genome biology. *Nature Biotechnology*.
647 2018;36:1121. doi:10.1038/nbt.4318.
- 648 14. Dudchenko O, Shamim MS, Batra SS, Durand NC, Musial NT, Mostofa R, et al.
649 The Juicebox Assembly Tools module facilitates de novo assembly of
650 mammalian genomes with chromosome-length scaffolds for under \$1000.
651 bioRxiv. 2018:254797. doi:10.1101/254797.
- 652 15. Belaghzal H, Dekker J and Gibcus JH. Hi-C 2.0: An optimized Hi-C procedure
653 for high-resolution genome-wide mapping of chromosome conformation.
654 *Methods (San Diego, Calif)*. 2017;123:56-65. doi:10.1016/j.ymeth.2017.04.004.
- 655 16. Kuratani S, Kuraku S and Nagashima H. Evolutionary developmental
656 perspective for the origin of turtles: the folding theory for the shell based on the
657 developmental nature of the carapacial ridge. *Evolution & Development*.
658 2011;13 1:1-14. doi:10.1111/j.1525-142X.2010.00451.x.
- 659 17. Matsuda Y, Nishida-Umehara C, Tarui H, Kuroiwa A, Yamada K, Isobe T, et al.
660 Highly conserved linkage homology between birds and turtles: bird and turtle
661 chromosomes are precise counterparts of each other. *Chromosome research : an
662 international journal on the molecular, supramolecular and evolutionary aspects
663 of chromosome biology*. 2005;13 6:601-15. doi:10.1007/s10577-005-0986-5.
- 664 18. Kuraku S, Ishijima J, Nishida-Umehara C, Agata K, Kuratani S and Matsuda Y.
665 cDNA-based gene mapping and GC3 profiling in the soft-shelled turtle suggest a
666 chromosomal size-dependent GC bias shared by sauropsids. *Chromosome
667 research : an international journal on the molecular, supramolecular and*

- 668 evolutionary aspects of chromosome biology. 2006;14 2:187-202.
669 doi:10.1007/s10577-006-1035-8.
- 670 19. Uno Y, Nishida C, Tarui H, Ishishita S, Takagi C, Nishimura O, et al. Inference
671 of the protokaryotypes of amniotes and tetrapods and the evolutionary processes
672 of microchromosomes from comparative gene mapping. PloS one. 2012;7
673 12:e53027. doi:10.1371/journal.pone.0053027.
- 674 20. Wang Z, Pascual-Anaya J, Zadissa A, Li W, Niimura Y, Huang Z, et al. The
675 draft genomes of soft-shell turtle and green sea turtle yield insights into the
676 development and evolution of the turtle-specific body plan. Nature Genetics.
677 2013;45:701. doi:10.1038/ng.2615.
- 678 21. Belton JM, McCord RP, Gibcus JH, Naumova N, Zhan Y and Dekker J. Hi-C: a
679 comprehensive technique to capture the conformation of genomes. Methods.
680 2012;58 3:268-76. doi:10.1016/j.ymeth.2012.05.001.
- 681 22. Servant N, Varoquaux N, Lajoie BR, Viara E, Chen CJ, Vert JP, et al. HiC-Pro:
682 an optimized and flexible pipeline for Hi-C data processing. Genome Biol.
683 2015;16:259. doi:10.1186/s13059-015-0831-x.
- 684 23. Sofueva S, Yaffe E, Chan WC, Georgopoulou D, Vietri Rudan M, Mira-
685 Bontenbal H, et al. Cohesin-mediated interactions organize chromosomal
686 domain architecture. The EMBO journal. 2013;32 24:3119-29.
687 doi:10.1038/emboj.2013.237.
- 688 24. Mayjonade B, Gouzy J, Donnadiou C, Pouilly N, Marande W, Callot C, et al.
689 Extraction of high-molecular-weight genomic DNA for long-read sequencing of
690 single molecules. BioTechniques. 2016;61 4:203-5. doi:10.2144/000114460.
- 691 25. Tanegashima C, Nishimura O, Motone F, Tatsumi K, Kadota M and Kuraku S.

- 692 Embryonic transcriptome sequencing of the ocellate spot skate *Okamejei*
693 *kenojei*. *Scientific data*. 2018;5:180200. doi:10.1038/sdata.2018.200.
- 694 26. Botero-Castro F, Figuet E, Tilak MK, Nabholz B and Galtier N. Avian Genomes
695 Revisited: Hidden Genes Uncovered and the Rates versus Traits Paradox in
696 Birds. *Molecular biology and evolution*. 2017;34 12:3123-31.
697 doi:10.1093/molbev/msx236.
- 698 27. Hotaling S and Kelley JL. The rising tide of high-quality genomic resources.
699 *Molecular Ecology Resources*. 2019;19 3:567-9. doi:10.1111/1755-0998.12964.
- 700 28. Nishimura O, Hara Y and Kuraku S. gVolante for standardizing completeness
701 assessment of genome and transcriptome assemblies. *Bioinformatics* (Oxford,
702 England). 2017;33 22:3635-7. doi:10.1093/bioinformatics/btx445.
- 703 29. Roach MJ, Schmidt SA and Borneman AR. Purge Haplotigs: allelic contig
704 reassignment for third-gen diploid genome assemblies. *BMC Bioinformatics*.
705 2018;19 1:460. doi:10.1186/s12859-018-2485-7.
- 706 30. Kadota M, Hara Y, Tanaka K, Takagi W, Tanegashima C, Nishimura O, et al.
707 CTCF binding landscape in jawless fish with reference to Hox cluster evolution.
708 *Scientific Reports*. 2017;7 1:4957. doi:10.1038/s41598-017-04506-x.
- 709 31. Miura H, Takahashi S, Poonperm R, Tanigawa A, Takebayashi S and Hiratani I.
710 Spatiotemporal developmental dynamics of chromosome organization revealed
711 by single-cell DNA replication profiling. in press.
- 712 32. Adey A, Morrison HG, Asan, Xun X, Kitzman JO, Turner EH, et al. Rapid, low-
713 input, low-bias construction of shotgun fragment libraries by high-density in
714 vitro transposition. *Genome Biology*. 2010;11 12:R119. doi:10.1186/gb-2010-
715 11-12-r119.

- 716 33. Langmead B and Salzberg SL. Fast gapped-read alignment with Bowtie 2.
717 Nature Methods. 2012;9:357. doi:10.1038/nmeth.1923.
- 718 34. Durand NC, Shamim MS, Machol I, Rao SS, Huntley MH, Lander ES, et al.
719 Juicer Provides a One-Click System for Analyzing Loop-Resolution Hi-C
720 Experiments. Cell systems. 2016;3 1:95-8. doi:10.1016/j.cels.2016.07.002.
- 721 35. Li H and Durbin R. Fast and accurate short read alignment with Burrows-
722 Wheeler transform. Bioinformatics (Oxford, England). 2009;25 14:1754-60.
723 doi:10.1093/bioinformatics/btp324.
- 724 36. Li H. A statistical framework for SNP calling, mutation discovery, association
725 mapping and population genetical parameter estimation from sequencing data.
726 Bioinformatics (Oxford, England). 2011;27 21:2987-93.
727 doi:10.1093/bioinformatics/btr509.
- 728 37. Quinlan AR and Hall IM. BEDTools: a flexible suite of utilities for comparing
729 genomic features. Bioinformatics (Oxford, England). 2010;26 6:841-2.
730 doi:10.1093/bioinformatics/btq033.
- 731 38. Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM.
732 BUSCO: assessing genome assembly and annotation completeness with single-
733 copy orthologs. Bioinformatics (Oxford, England). 2015;31 19:3210-2.
734 doi:10.1093/bioinformatics/btv351.
- 735 39. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al.
736 Full-length transcriptome assembly from RNA-Seq data without a reference
737 genome. Nat Biotechnol. 2011;29 7:644-52. doi:10.1038/nbt.1883.
- 738 40. Wang M and Kong L. pblat: a multithread blat algorithm speeding up aligning
739 sequences to genomes. BMC Bioinformatics. 2019;20 1:28.

740 doi:10.1186/s12859-019-2597-8.

741 41. Ryan JF. Baa.pl: A tool to evaluate de novo genome assemblies with RNA
742 transcripts. arXiv e-prints. 2013.

743 42. Kent WJ. BLAT--the BLAST-like alignment tool. *Genome Res.* 2002;12 4:656-
744 64. doi:10.1101/gr.229202.

745 43. Imakaev M, Fudenberg G, McCord RP, Naumova N, Goloborodko A, Lajoie
746 BR, et al. Iterative correction of Hi-C data reveals hallmarks of chromosome
747 organization. *Nature Methods.* 2012;9:999. doi:10.1038/nmeth.2148.

748

749

750

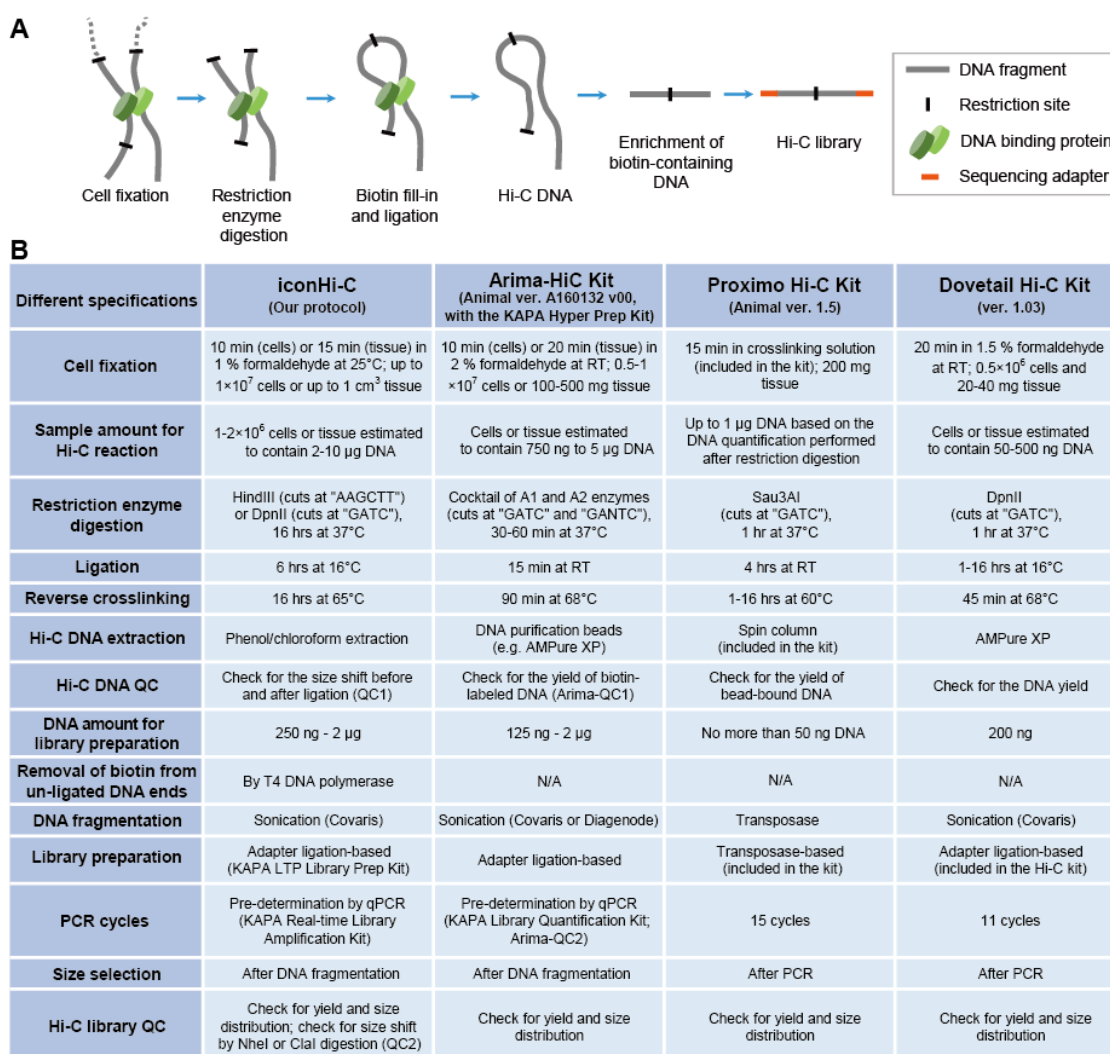
751 **Table 1:** Overview of the specification of the scaffolding programs released to date.

Program	Support and availability	Input data requirement	Other information	Literature
LACHESIS	Developer's support discontinued; intricate installation	Generic bam format	No function to correct scaffold misjoins	[4]
HiRise	Open source version at GitHub not updated since 2015	Generic bam format	Employed in Dovetail Chicago/Hi-C service. Default input sequence length cutoff=1000 bp	[7]
3d-dna	Actively maintained and supported by the developer	Not compatible with multiple enzymes; Accept only Juicer mapper format	Default parameters: -t 15000 (input sequence length cutoff), -r 2 (no. of iterations for misjoin correction)	[10, 34]
SALSA2	Actively maintained and supported by the developer	Compatible with multiple enzymes; generic bam (bed) file, assembly graph, unitig, 10x link files	Default parameters: -c 1000 (input sequence length cutoff), -i 3 (no. of iterations for misjoin correction)	[8, 9]

752

753

754 **Figures**



755

756 **Figure 1: Hi-C library preparation. (A) Basic procedure. (B) Comparison of Hi-C**
 757 **library preparation methods. Included here are only the major differences between the**
 758 **methods. The KAPA Hyper Prep Kit (KAPA Biosystems) is assumed to be conjunctly**
 759 **used with Arima Hi-C Kit, among the several specified kits. See Supplementary**
 760 **Protocol S1 for the full version of the iconHi-C protocol which was derived from the**
 761 **protocol previously introduced [23].**

762

763

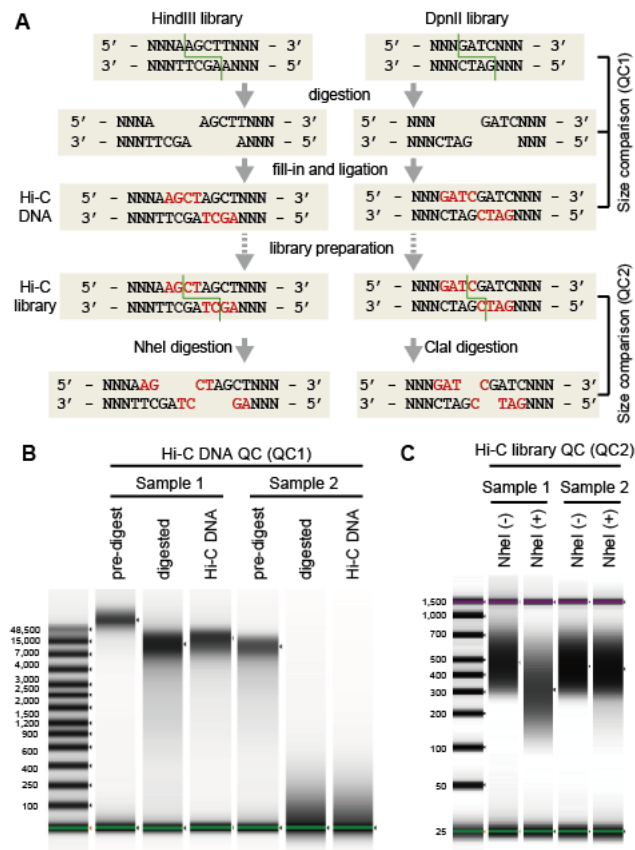


764

765 **Figure 2:** A juvenile softshell turtle *Pelodiscus sinensis*.

766

767

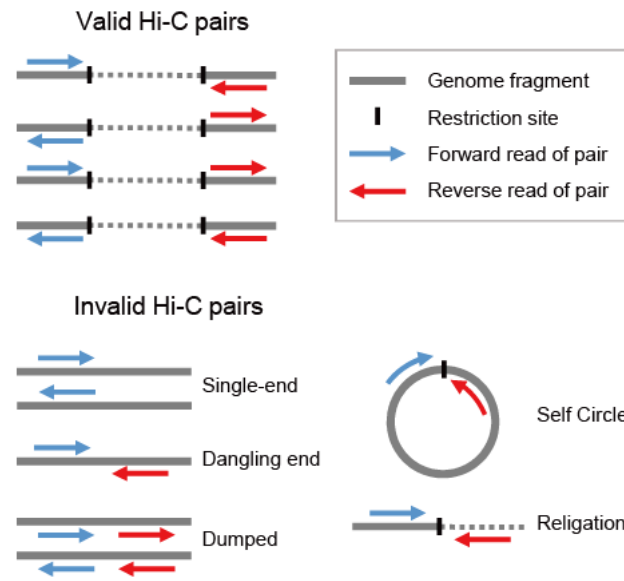


769

770 **Figure 3:** Structure of Hi-C DNA and principle of quality controls. (A) Schematic
 771 representation of the library preparation workflow based on HindIII or DpnII digestion.
 772 Patterns of restriction are indicated by the green lines. Nucleotides that were filled in are
 773 indicated by the letters in red. (B) Size shift analysis of HindIII-digested Hi-C DNA
 774 (QC1). Shown are the representative images of qualified (Sample 1) and disqualified
 775 samples (Sample 2). (C) Size shift analysis of the HindIII-digested Hi-C library (QC2).
 776 Shown are the representative images of the qualified (Sample 1) and disqualified
 777 (Sample 2) samples. Size distributions were measured with Agilent 4200 TapeStation.

778

779

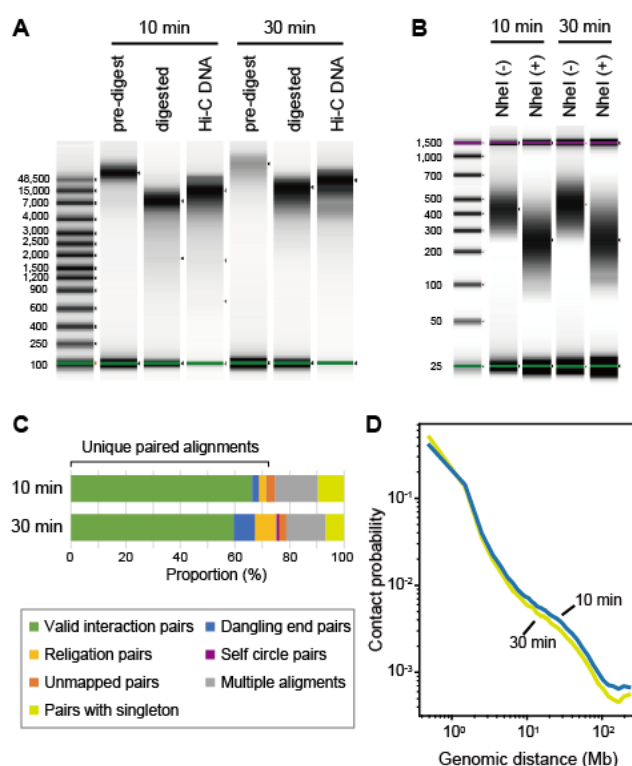


780

781 **Figure 4:** Post-sequencing quality control of Hi-C reads. Read pairs were categorized
 782 into valid and invalid pairs by HiC-Pro, based on their status in the mapping to the
 783 reference genome (see Methods). This figure was adapted from the literature originally
 784 introducing HiC-Pro [22].

785

786



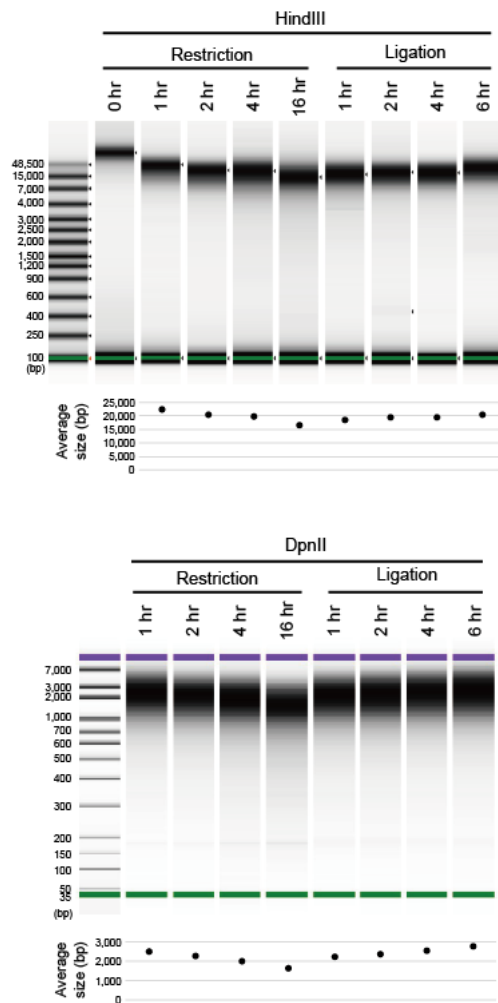
788

789 **Figure 5:** Effect of cell fixation duration. (A) QC1 of the HindIII-digested Hi-C DNA
 790 of human GM12878 cells fixed for 10 or 30 minutes in 1% formaldehyde. (B) QC2 of
 791 the HindIII-digested library of human GM12878 cells. (C) Quality control of the
 792 sequence reads by HiC-Pro using 1M read pairs. See Fig. 4 for the details of the read
 793 pair categorization. See Supplementary Table S7 for the actual proportion of the reads
 794 in each category. (D) Contact probability measured by the ratio of observed and
 795 expected frequencies of Hi-C read pairs mapped along the same chromosome [43].

796

797

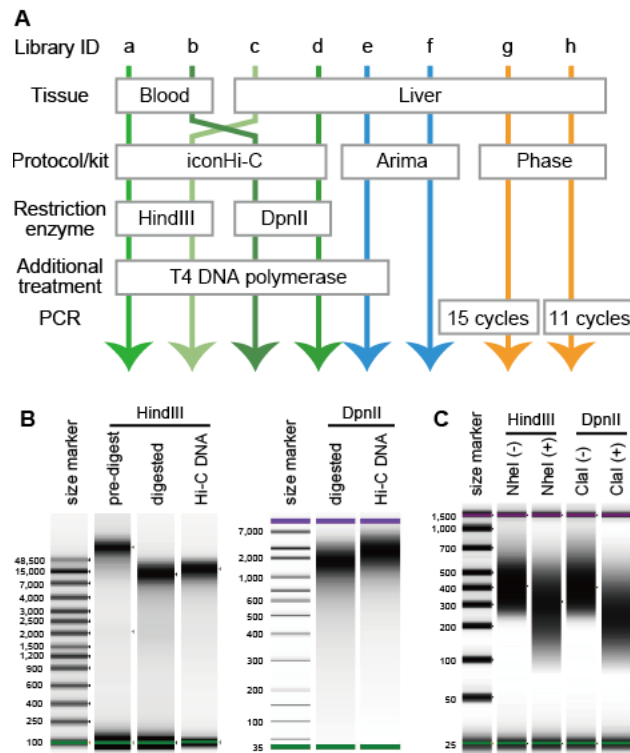
798



799

800 **Figure 6:** Testing variable durations of restriction and ligation of Hi-C DNA. Length
 801 distributions of the DNA molecules prepared from human GM12878 cells after variable
 802 durations of restriction and ligation are shown. Size distribution for the HindIII-digested
 803 samples (top) and DpnII-digested samples (bottom) were measured by Agilent 4200
 804 TapeStation and Agilent Bioanalyzer, respectively.

805



807

808 **Figure 7:** Softshell turtle Hi-C libraries prepared for our methodological comparison.

809 (A) Lineup of the prepared libraries. This chart includes only the conditions that varied

810 preparation methods between these libraries, and the rest of the preparation workflows

811 are described in Supplementary Protocol S1 for the non-commercial ('iconHi-C')

812 protocol and the manuals of the commercial kits. (B) Quality control of Hi-C DNA

813 (QC1) for Library c and d. The prepared Hi-C DNA for the Chinese softshell turtle liver

814 samples were digested with either HindIII or DpnII. (C) Quality control of Hi-C

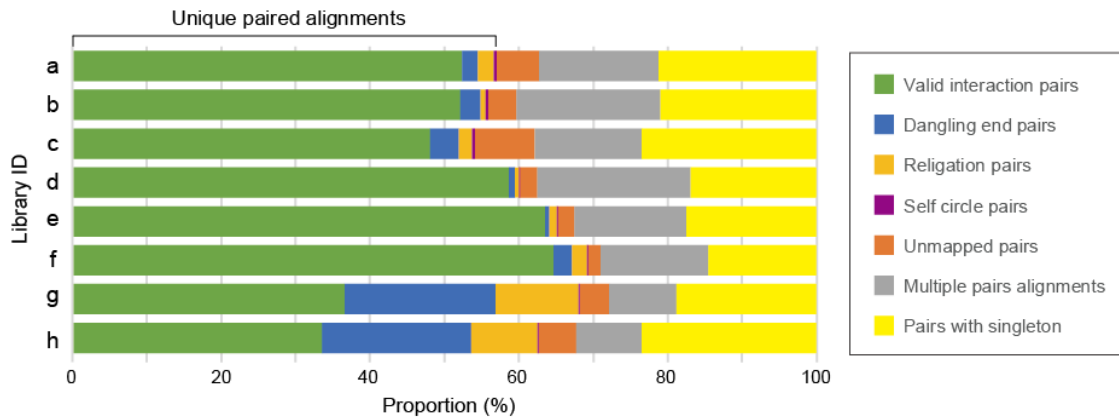
815 libraries (QC2). The prepared softshell turtle liver HindIII library was digested by NheI,

816 and the DpnII library was digested by ClaI (see Fig. 3 for the technical principle). See

817 Supplementary Fig. S3 for the QC1 and QC2 results for the samples prepared from the

818 blood of this species.

819



820

821 **Figure 8:** Results of the post-sequencing quality control with HiC-Pro. One million

822 read pairs were used for computation with HiC-Pro. See Fig. 7A for the preparation

823 conditions of Library a-h, Fig. 4 for the categorization, and Supplementary Table S3 for

824 the actual proportion of the reads in each category. Post-sequencing quality control

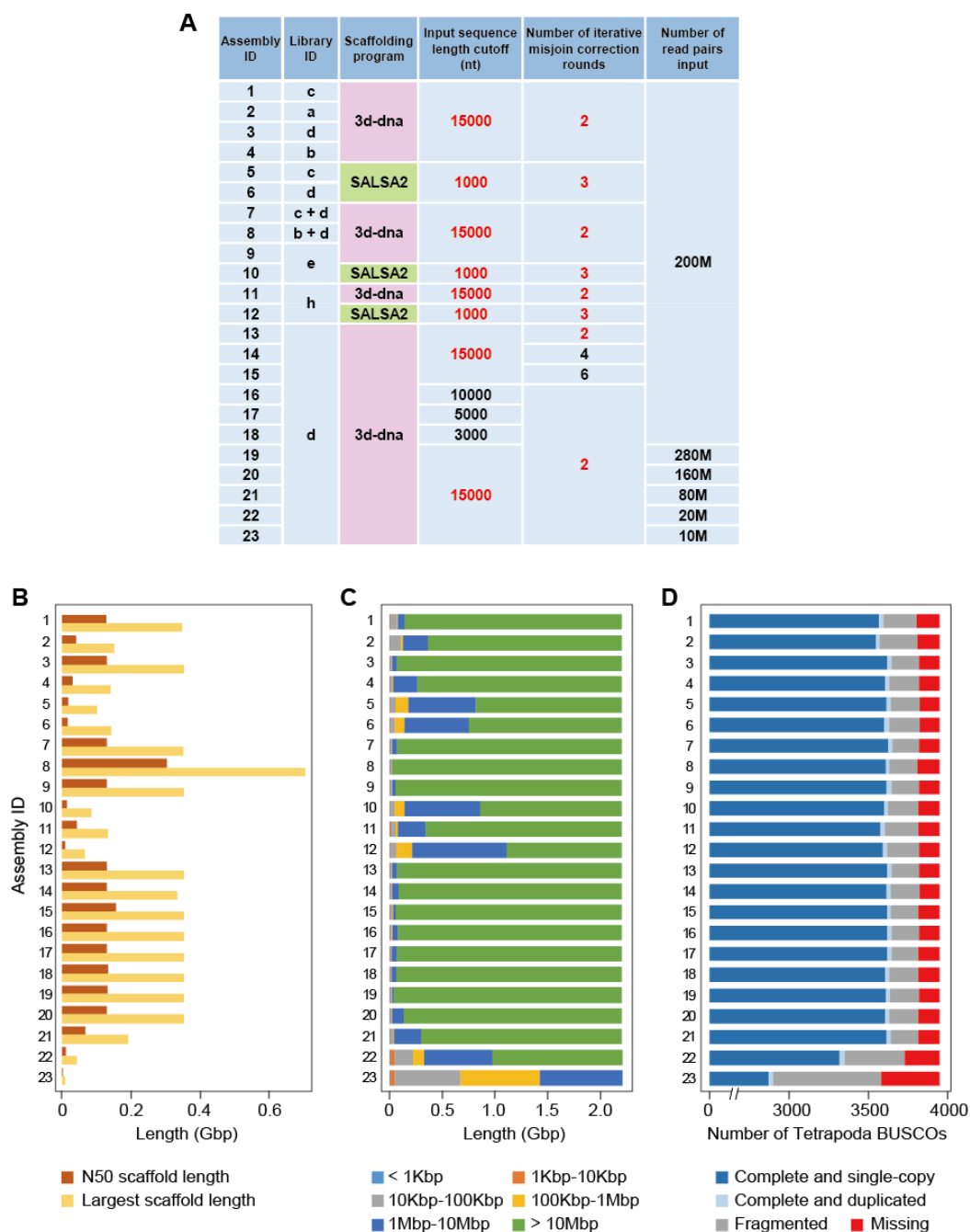
825 using variable read amounts (500 K–200 M pairs) for one of these softshell turtle

826 libraries (Supplementary Table S6) and human GM12878 libraries (Supplementary

827 Table S2) shows the validity of this quality control with as few as 500 K read pairs.

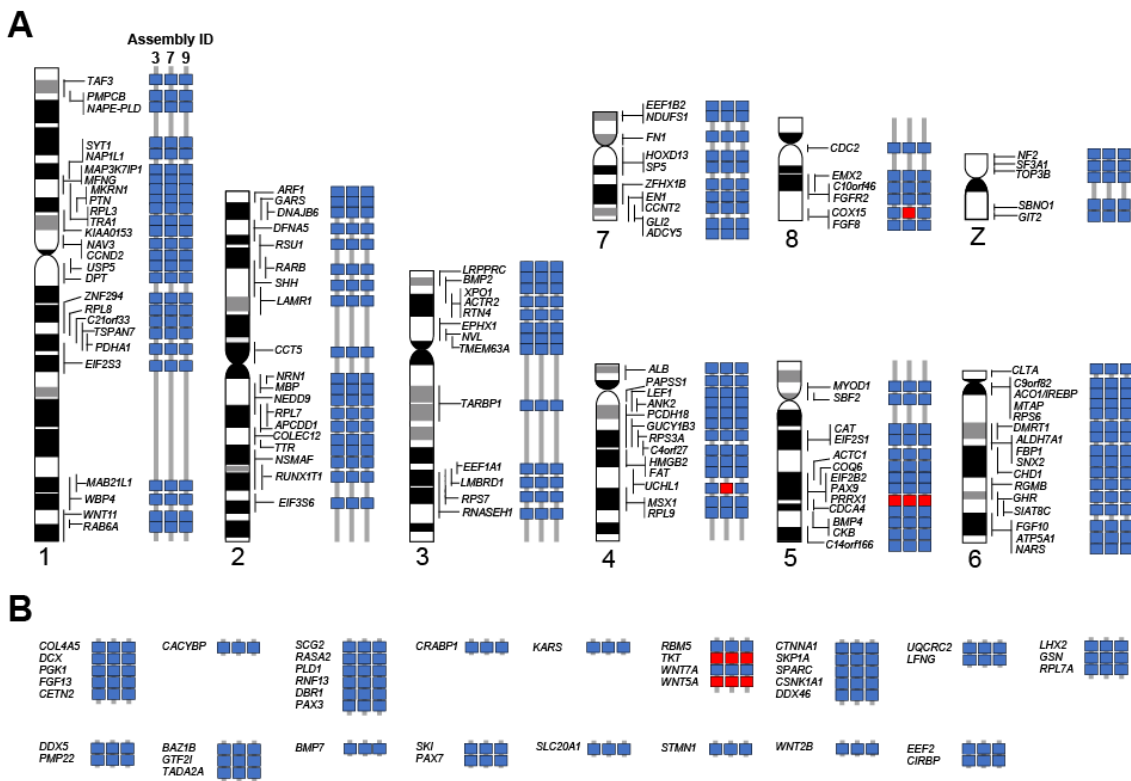
828

829



830

831 **Figure 9:** Comparison of Hi-C scaffolding products. (A) Scaffolding conditions to
 832 produce Assembly 1 to 23. Default parameters are shown with red letters. (B) Total and
 833 N50 scaffold lengths. (C) Scaffold length distributions. (D) Gene space completeness.
 834 See the panel A for Library IDs and Supplementary Table S4 for raw values of the
 835 metrics in B–D.



837

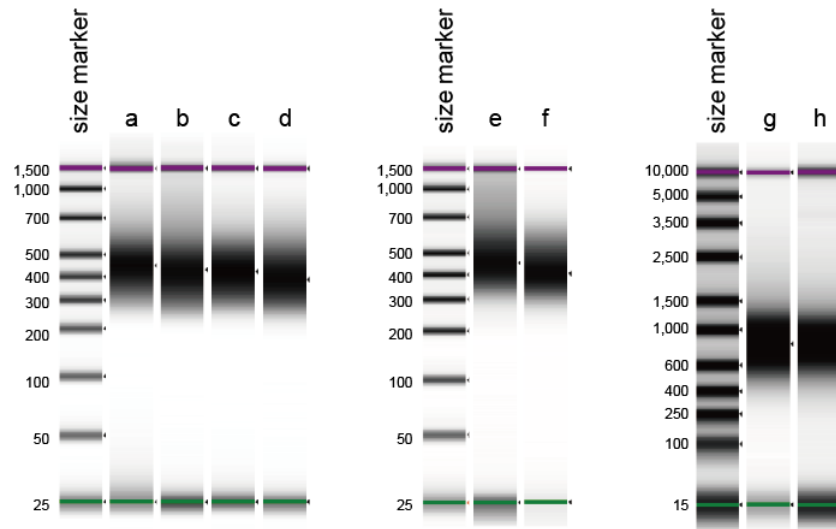
838 **Figure 10:** Cytogenetic validation of Hi-C scaffolding results. On the scaffolded
 839 sequences of Assembly 3, 7, and 9, we evaluated the consistency of the positions of the
 840 selected genes that were previously localized on 8 macrochromosomes and Z
 841 chromosome (A) and microchromosomes (B) by chromosome FISH [17-19] (see
 842 Results). Concordant and discordant gene locations on individual assemblies are
 843 indicated with blue and red boxes, respectively. The arrays of genes without idiograms
 844 in B were identified on chromosomes that are cytogenetically indistinguishable from
 845 each other.

846

847

848

849



850

851 **Supplementary Figure S1:** DNA size distribution of the softshell turtle Hi-C libraries.

852 Size distribution of the libraries was analyzed by Agilent 4200 TapeStation using the

853 High Sensitivity D1000 kit for Library a-f and the High Sensitivity D5000 kit for Library

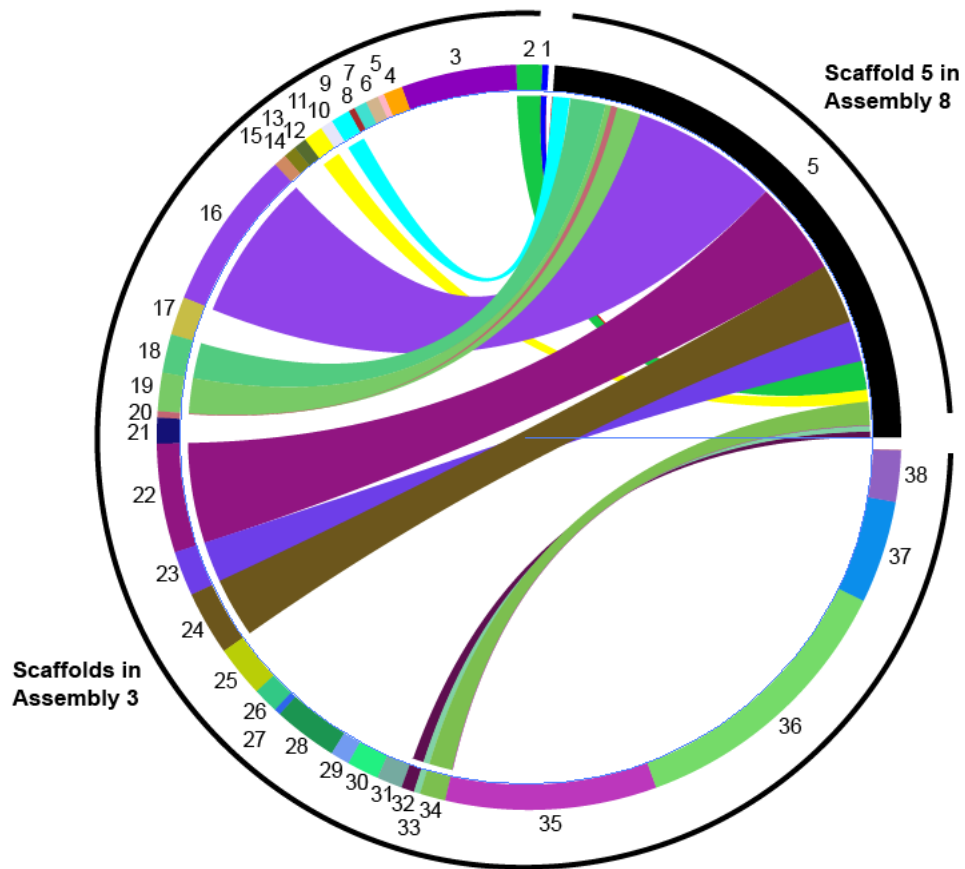
854 g and h.

855

856

857

858



859

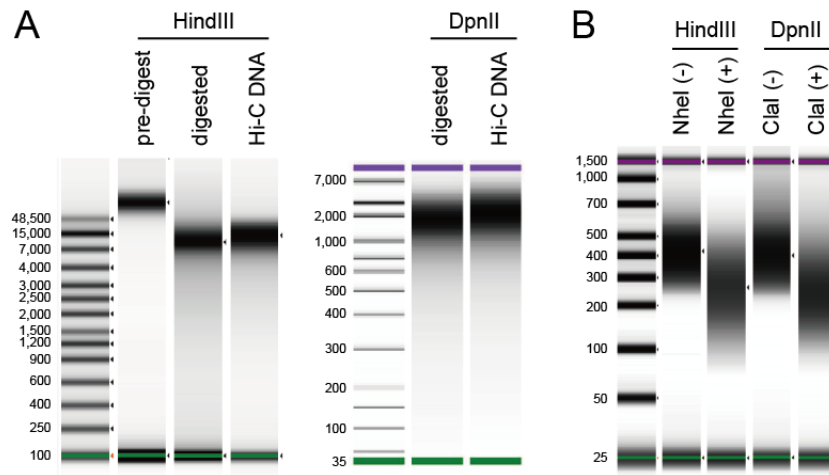
860 **Supplementary Figure S2:** Structural analysis of the possibly overassembled scaffold
861 in Assembly 8. This figure shows the nucleotide sequence-level correspondence of the
862 whole sequence of the scaffold 5 of Assembly 8 to 14 scaffolds of Assembly 3. Note
863 that the scaffold 5 of Assembly 8 accounts for approximately one-third of the estimated
864 genome size, and that some of the scaffolds of Assembly 3 in the figure have multiple
865 high-similarity regions in the scaffold 5 of Assembly 8.

866

867

868

869



870

871 **Supplementary Figure S3:** Pre-sequencing quality control of softshell turtle blood Hi-

872 C libraries (Library a and b). (A) Quality control of Hi-C DNAs (QC1). Hi-C DNA was

873 prepared from the Chinese softshell turtle blood by HindIII or DpnII digestion (see Fig.

874 7A for the detail). (B) Quality control of Hi-C libraries (QC2). The prepared softshell

875 turtle blood library employing HindIII was digested by NheI, and the one employing

876 DpnII was digested by ClaI (see Fig. 3 for the technical principle).

877

878



[Click here to access/download](#)

Supplementary Material

[Supplementary_TableS1_draft-genome.pdf](#)





[Click here to access/download](#)

Supplementary Material

[Supplementary_TableS2_GM12878-reduced-reads.pdf](#)





[Click here to access/download](#)

Supplementary Material

[Supplementary_TableS3_Ps_Lib_QC_1M-Mod.pdf](#)





[Click here to access/download](#)

Supplementary Material

[Supplementary_TableS4_all_scaffolding.pdf](#)

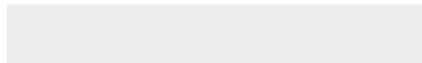




[Click here to access/download](#)

Supplementary Material

[Supplementary_TableS5_RNAassembly_mapping.pdf](#)





[Click here to access/download](#)

Supplementary Material

[Supplementary_TableS6_Ps-reduced-reads.pdf](#)





Click here to access/download

Supplementary Material

Supplementary_TableS7_10-and-30-minutes_QC.pdf





Click here to access/download
Supplementary Material
Supplementary_Protocol_S1.pdf





Click here to access/download

Supplementary Material

Supplementary_Protocol_S2_to_support_multiple_enzymes.pdf



Shigehiro Kuraku, Ph.D.
 Team Leader
 Laboratory for Phyloinformatics
 RIKEN Center for Biosystems Dynamics Research (BDR)

Tel: +81-(0)-78-306-3331
 Email: shigehiro.kuraku@riken.jp
 URL: <https://www.bdr.riken.jp/en/research/labs/kuraku-s/>

June 5, 2019

Editorial Board Member, *GigaScience*

Dear Dr. Takashi Gojobori,

Accompanying this letter is our manuscript entitled, '**Multifaceted Hi-C benchmarking: what makes a difference in chromosome-scale genome scaffolding?**' by **Kadota, Nishimura, et al.** to be considered for publication in the journal *GigaScience*. It is accompanied by the Supplementary Information file that is covering the detailed statistics and results of individual analyses.

Chromosome-scale scaffolding using Hi-C has increasingly been employed in *de novo* genome assembly, but its best practice has not been discussed in depth from methodological viewpoints. In the submitted manuscript, we report a benchmarking for evaluating various factors in sample preparation, sequencing, and computation. As a result, we have identified some key factors that help improve Hi-C scaffolding, such as the choice of tissues and restriction enzymes, duration of enzymatic reactions, and the choice of scaffolding programs and parameters. To our knowledge, this is the first-ever comparison of multiple sample preparation kits/protocols and computational programs for Hi-C scaffolding, by a third party in academia. The largest product of our study is the release of an original Hi-C protocol that incorporates the lessons learned from our benchmarking. We understand that *GigaScience* has created an active forum of readers interested in both practical aspects of genome sequencing and technical aspects in the computation for genome assembly. Therefore, we think that *GigaScience* is the most suitable journal to publish our study reported in the present manuscript.

The submitted manuscript has been shared among all the authors and approved by them. It has not been published and even submitted to any other journal. We have no conflict of interest regarding this manuscript. As preferred reviewers of this manuscript, we would nominate the researchers below for the reasons included.

Chris Amemiya	University of California, Merced, US	camemiya@ucmerced.edu
Expertise in vertebrate genome structure evolution		
Jessica Alföldi	Broad Institute, US	jalfoldi@broadinstitute.org
Expertise in genome sequence data production and analysis		
Nicolas Servant	Institut Curie, France	nicolas.servant@curie.fr
Expertise in technical evaluation of Hi-C data		
Rob Waterhouse	Lausanne University, Switzerland	robert.waterhouse@unil.ch
Expertise in technical evaluation of genome assemblies		



We expect that our study will provide a technical baseline for Hi-C scaffolding for building chromosome-scale genome sequences, which influences a wide spectrum of genomic studies across taxonomic divisions of diverse organisms. We hope that you will find our manuscript reporting an unprecedented suite of technical resources worthy of publication in *GigaScience*.

Sincerely yours,

A handwritten signature in black ink, reading 'Shigehiro Kuraku' in Japanese characters (工樂樹洋).

Shigehiro Kuraku, Ph.D.