

1 How best to identify chromosomal interactions: a comparison of approaches

2

3 James O. J. Davies¹, A. Marieke Oudelaar¹, Douglas R. Higgs¹ and Jim R. Hughes¹

4

5 ¹Medical Research Council (MRC) Molecular Haematology Unit, Weatherall

6 Institute of Molecular Medicine, Oxford University, Oxford, United Kingdom.

7

8

9

10 **Summary (100 words)**

11

12 Chromosome conformation capture (3C) methods are central to understanding
13 the link between nuclear structure and function and the physical interactions
14 between distal regulatory elements and promoters. However, no one method is
15 appropriate for all biological questions as each variant differs markedly in
16 resolution, reproducibility, throughput and biases. A thorough appreciation of
17 the strengths and weaknesses of each technique is critical when choosing the
18 correct method for a specific application or for gauging how best to interpret
19 different sources of data. In addition, the method of analysis can have a profound
20 effect on the output and must be carefully considered.

21

22

23 **Preface (100 words)**

24

25 The location of potential regulatory elements can now be routinely and rapidly
26 mapped genome-wide using chromatin accessibility and ChIP-seq based assays<sup>1-
27 4</sup>, however an outstanding challenge is to determine which regulatory elements
28 control which genes⁵. This is problematic because regulatory elements are
29 scattered over large distances (up to millions of base pairs) around the genes
30 they regulate⁶. In order to influence transcription, the protein complexes at distal
31 regulatory elements, such as enhancers, make physical contact with gene
32 promoters. These interactions can be detected using chromosome conformation
33 capture (3C) technology and form the basis for identifying a gene's regulatory
34 elements.

35

36 There is a growing body of evidence that genomic organization within the
37 nucleus has an important role in determining gene expression patterns⁷.
38 Although a large amount has been learnt about nuclear organization from
39 microscopic study of the nucleus, particularly with fluorescence in situ
40 hybridization (FISH), the majority of recent discoveries regarding sequence
41 specific chromatin structure have been made using 3C technology. These include
42 the description of large scale (>100kb) structures such as topologically
43 associated domains (TADs)^{8,9} and interactions on a much smaller scale such as
44 those with regulatory elements¹⁰ and within the gene body itself^{11,12}.

45 The first assays that used a restriction enzyme digest followed by a ligation step
46 to detect physical interactions between pieces of DNA were performed in the late
47 1980s to show that DNA-DNA interactions occur in plasmids¹³. However, it was
48 not until the seminal paper by Dekker et al.,¹⁴ that it became possible to define
49 interactions between two specific DNA sequences in eukaryotic cells, and studies
50 demonstrating interactions between regulatory elements in mammalian
51 genomes quickly followed¹⁰. The assays are all based on the principle that
52 chromatin interactions can be crosslinked, cut and then re-ligated so that
53 sequences in close physical proximity within the nucleus become linked to each
54 other. The ligation junctions that result thus reflect the three-dimensional
55 organization of the genome at time of fixation and can therefore be used to infer
56 chromatin structure (**Figure 1a**).

57 The library of re-ligated fragments (3C library) could in principle be interrogated
58 by simply sequencing the rearranged genomic DNA using next-generation
59 sequencing technologies. This has been done for small genomes (Sexton et al.
60 2012), but is not feasible for high-complexity mammalian genomes. The different
61 3C methods can be seen as approaches to sample the ligation events from
62 specific regions, or sequence only the informative ligation junctions to overcome
63 the complexity of large genomes.

64 When the 3C library is being manufactured each restriction fragment in every
65 individual cell has many potential ligation partners (**Figure 1b**)¹⁵. Live cell

imaging studies show that chromatin is in a state of constant flux within the nucleus^{15,16} and single-cell 3C assays, therefore, show huge variability in the ligation junctions between cells¹⁷. Thus the ligation junctions present within the 3C library represent the probability distribution of all possible ligation junctions from a group of cells, in a particular state. It is therefore vital to accurately quantify the ligation junctions present in the 3C library and with significant depth of sequencing in order to build up a profile of the probability distribution underlying the ligation junctions (**Figure 1c**). Simply stating that a junction can or cannot be found is not meaningful because if one samples the library at sufficient depth then it is likely that a ligation junction could be found with most other fragments in the genome. It is therefore important to understand if a dataset has sampled the interactions to sufficient depth to be reproducible or whether it is sporadically detecting interactions from within the regulatory landscape. This depends on the sensitivity of the assay, which is determined by the number of unique ligation junctions that can be enriched for a given region of interest in the genome.

3C library considerations

The first step of most 3C techniques is to generate a 3C library (**Figure 2**). Generally cells are initially fixed using formaldehyde crosslinking. Recently protocols have been published that allow the digestion and ligation reactions to be performed without the disruption of the nucleus ('in situ' 3C library manufacture)¹⁸, which has considerable advantages for determining accurate interaction profiles^{18,19}. It is also possible to manufacture 3C libraries without this fixation step by embedding the cells in agarose plugs, in order to maintain the nuclear shape and structure. Importantly, this results in similar interaction profiles to those generated from formaldehyde fixed material (although there is greater background noise), showing that the interactions detected are not simply an artifact of the fixation process itself¹⁸.

In most 3C techniques the chromatin is cut using a restriction enzyme. The choice of restriction enzyme determines the maximum resolution of the data because interactions can only be detected at restriction enzyme cut sites (**Figure 3**). 3C libraries were initially manufactured using six-cutter restriction enzymes, which limited the potential resolution since these enzymes cut at approximately 4 kb intervals (95% of fragments <12.3 kb). Subsequently four-cutter restriction enzymes have been used, which cut on average every 256bp (95% of fragments <800bp), potentially generating data with 16 times greater resolution.

The complexity of 3C libraries is related to the square of the number of restriction fragments in the genome²⁰. Thus increasing the resolution (by using an enzyme that cuts more frequently) or the genome size has dramatic effects on the sequencing requirements. 3C library complexity is often a limiting factor, particularly when investigating the ligation events from specific restriction fragments of interest, as only four interactions can be determined from each

108 fragment in each cell (one from the end of the fragment on each allele). In
109 practical terms the critical factors that determine library complexity are the
110 initial numbers of cells used, the digestion and ligation efficiency and the
111 cumulative loss of material from each step before sequencing.

112 When library complexity and / or sequencing depth are insufficient to analyze
113 the data at the level of individual restriction fragments, the resolution of 3C-
114 based experiments is determined by the bin size used to view the data.
115 Windowing or binning improves the signal strength and reduces biases by
116 combining all of the data in the window together, which allows a meaningful
117 interaction profile to be determined when the raw signal is weak. However,
118 windowing has a number of disadvantages. The profile becomes skewed by the
119 restriction enzyme cut site density and it smoothens the profile of the original
120 signal, obscuring the quality of the underlying data and reducing resolution.
121 Ideally, there should be sufficient signal strength and reproducibility at the level
122 of individual restriction fragments, to allow the data to be reported without
123 windowing.

124 **Chromosome conformation capture (3C)**

125 Initially ligation junctions were detected in 3C libraries by polymerase chain
126 reaction (PCR) followed by gel electrophoresis, which was subsequently
127 replaced by real-time PCR (**Figure 2**).

128 Many important discoveries were made using 3C. The technique was initially
129 used to define the spatial organization of chromosomes in yeast¹⁴. Interactions at
130 the murine β -globin locus were determined soon after¹⁰. Subsequently 3C has
131 been used to analyze many other loci, including the T-helper type 2 cytokine
132 locus²¹, the immunoglobulin light chain locus²² and the alpha globin locus²³⁻²⁵.
133 However, conventional 3C has largely been superseded by other methods
134 because it is laborious and as it is only able to detect interactions between a
135 small number of fixed restriction fragments (**Figure 3**). This approach focuses
136 researchers on confirming suspected interactions rather than identifying new
137 interactions in an unbiased way. In addition, it can be difficult to obtain
138 reproducible results, and correction for biases due to differences in amplification
139 with different primer pairs requires great care.

140 **Circular chromosome conformation capture (4C)**

141 The next major advance in the field was the development of 4C, which allows all
142 of the potential interacting partners to be identified from any specific point of
143 interest in the genome^{26,27}. 4C uses standard 3C library preparation, often with a
144 4-cutter restriction enzyme, after which inverse PCR is used from the 'viewpoint'
145 fragment to amplify any interacting partners (**Figure 2**). Initially the interacting
146 partners were read out using a microarray, but this has now been replaced by
147 high-throughput sequencing.

148 For inverse PCR to work efficiently, small circularized fragments must be
149 generated either by optimizing the initial ligation reaction²⁶ or by the more
150 widely used method of performing a second restriction enzyme digest and
151 ligation of the extracted DNA with a different restriction enzyme²⁷. The inverse
152 PCR in 4C data is affected by the GC content and size of the interacting fragments
153 to the extent that some large GC-rich fragments simply fail to amplify. This can
154 result in serious problems with bias in the interaction profile²⁸ and make the
155 interactions appear more discrete than they really are (**Figure 4**). 4C provides
156 good resolution and sensitivity, when carefully performed, allowing very long-
157 range interactions to be detected. In order to generate high-quality 4C profiles, it
158 is necessary to perform multiple PCRs to generate material from the DNA of a
159 million or more cells²⁹. However, it is difficult to determine the sensitivity of
160 conventional 4C accurately, as it is not possible to differentiate between PCR
161 duplicates and unique ligation junctions.

162 In a recently developed one vs all protocol (described as UMI-4C³⁰), the 3C
163 library is sonicated and sequencing adapters are ligated to one end of each
164 sonicated 3C fragment. Using a primer in the target sequence and a universal
165 adapter primer, interacting fragments can be amplified and sequenced. Such
166 interactions can be quantified by using the unique fragment ends produced by
167 the sonication step and this shows that about 5000-10,000 unique read pairs can
168 be detected per µg of input material. UMI-4C also allows for multiplexing up to
169 20-50 baits, though the individual nested PCR reactions still need to be
170 optimized for each viewpoint³⁰.

171 A large number of seminal observations have been made using 4C techniques. It
172 has been widely used to describe interactions between promoters and enhancers
173 and to show how these change during differentiation and development³¹. Very
174 long-range interactions between active genes have been demonstrated in *cis*
175 using 4C³². 4C has also been used to link potential regulatory single nucleotide
176 polymorphisms (SNPs) with genes³³. It can be used to determine structural
177 changes such as chromosomal rearrangements³⁴ although this is probably more
178 easily done using conventional karyotype analysis and FISH. It has also been
179 used to define disease mechanisms, including the demonstration that
180 chromosomal translocations can result in distal enhancers becoming juxtaposed
181 to an oncogene leading to malignancy³⁵.

182 **Chromosome Conformation Capture Carbon Copy (5C)**

183 5C can be used to study all interactions within a particular locus and is based on
184 the use of highly multiplexed ligation-mediated amplification in the region of
185 interest. It uses primer pairs that anneal on either side of the ligation junctions in
186 a 3C library. These primer pairs are ligated and can then be amplified using T7
187 and T3 sequences incorporated in them in a single PCR reaction, which can be
188 analyzed by microarray or sequencing³⁶ (**Figure 2**).

189 Differences in the hybridization efficiency of the probes can cause bias in 5C. In
190 addition it is only possible to determine interactions between forward and
191 reverse probes and it is impossible to design probes to the ends of some
192 fragments, so the interaction profiles can have large areas, which are not
193 interrogated (**Figure 4**).

194 5C can define functional interactions for all the genes in a locus simultaneously,
195 but it is difficult to be precise about the sensitivity of the method as the levels of
196 PCR duplication cannot be determined. We find that 5C occasionally misses
197 weak, long-range interactions, which are detectable by Hi-C, 4C and Capture-C
198 (**Figure 5**).

199 5C has been used to determine interaction profiles at the pilot regions of the
200 ENCODE project⁶. In addition analysis of the X-chromosome by 5C provided
201 some of the first evidence for the existence of topologically associating domains
202 (TADs)⁹. Massively multiplexed 5C can generate all vs all interaction maps of a
203 region of interest but it requires significant financial resources for somewhat
204 limited resolution and scale. For example, 5826 oligonucleotides were required
205 to cover 1% of the genome (for the ENCODE pilot regions) and only two
206 replicates of three different cell lines were analyzed⁶.

207 **Capture-C**

208 Capture-C is able to generate genome-wide interaction profiles from hundreds of
209 viewpoints in a single assay³⁷. It combines 3C library preparation with a four-
210 cutter restriction enzyme and oligonucleotide capture technology. The 3C
211 libraries are enriched for fragments of interest using biotinylated capture probes
212 designed for each viewpoint, after which these fragments are amplified and
213 sequenced.

214 Although Capture-C involves several rounds of PCR amplification this does not
215 cause much bias, in contrast to 4C. This is because the 3C library is sonicated
216 prior to the PCR, resulting in uniform small fragments with random unique ends
217 that allow PCR duplicates to be collapsed during data analysis.

218 While capable of producing hundreds of informative, genome-wide tracks in a
219 single experiment, initially the individual tracks themselves did not have the
220 depth of data of a good single 4C experiment from the same region³⁷ (**Figure 2**).
221 Next Generation (NG) Capture-C was subsequently developed¹⁵. This uses a new,
222 more flexible and efficient oligonucleotide capture process that markedly
223 increases the sensitivity of the assay. In high-quality datasets, more than 100,000
224 read pairs commonly contribute to the interaction profile. This makes NG
225 Capture-C the most sensitive and highest resolution assay currently available for
226 mammalian genomes, allowing the data to be expressed at maximum resolution
227 (per restriction fragment) (**Figure 3**). This high sensitivity combined with the
228 reliability of the technique allow weak long-range *cis* and even *trans* interactions
229 to be quantified, which is important because it provides confidence that
230 stronger, functional interactions are not being missed. It also makes it possible to

231 adapt the technique to small cell numbers. In addition, several independent 3C
232 libraries (e.g. from different cell types or different stages of development) can be
233 processed in a single tube. This greatly increases throughput and allows
234 meaningful subtractive analysis of chromosome conformation of multiple
235 replicates in different cell types (**Figure 4&5**).

236 **Hi-C**

237 Hi-C generates maps of interaction between all parts of the genome (all against
238 all). This technique uses a modified method for 3C library preparation. A
239 biotinylated nucleotide fill-in is performed after the restriction enzyme digestion
240 and followed by blunt end ligation. After DNA extraction and sonication of the
241 material, a streptavidin bead pull-down is performed to concentrate ligation
242 junctions, which can then be analyzed by high throughput sequencing (**Figure**
243 **2**).

244 Hi-C has great sensitivity for determining megabase scale interactions and is
245 unparalleled in its ability to determine large-scale chromatin structure. This is
246 because it combines data from many restriction fragments to define robust
247 interactions on a chromosomal scale. However, the number of interactions
248 determined from any individual restriction fragment are around 100-fold lower
249 than in 4C or Capture-C, even in the recent Hi-C datasets at kilobase resolution¹⁸.
250 Unless Hi-C is done at very high resolution it is a relatively insensitive method to
251 determine fine-scale (< 40 kb) interactions between regulatory elements present
252 within TADs.

253 Hi-C has relatively few biases, but is still affected by systematic biases relating to
254 the distance between restriction sites, GC content and sequence uniqueness and
255 several methods have been developed to attempt to correct these³⁸⁻⁴⁰.

256 Since it was initially reported²⁰, Hi-C experiments have produced data at
257 increasing resolution predominantly through massive increases in sequencing
258 depth. Initially 40kb bins were used²⁰ but this was improved to 5-10kb⁴¹ and
259 more recently 1kb bins have been reported in human cell lines¹⁸. Producing Hi-C
260 data of this quality for large mammalian genomes requires a gargantuan effort
261 involving the analysis of several billion read pairs per sample.

262 Hi-C has been used extensively to define TADs genome wide⁸ and to determine
263 the structure of the chromosome during mitosis⁴². Hi-C has also been used to link
264 trans interactions with sites associated with chromosomal translocations⁴³. At
265 present this is the best available technique for determining genome wide (all vs
266 all) maps of interactions.

267 Hi-C interactions have also been determined from single cells¹⁷ by picking and
268 sequencing single intact nuclei during Hi-C library preparation. This showed that
269 there is considerable variability between individual cells, but the data are not
270 informative at the level of regulatory features.

271 **Hi-C variants**

272 A further increase in resolution has recently been achieved in *S. Cerevisiae* by
273 substituting restriction enzymes with micrococcal nuclease (MNase)¹². This
274 shows similar but not identical maps of interaction compared to Hi-C. Long-
275 range interactions are poorly captured by Micro-C. It shows that ~90% of all
276 interactions are within 1 kb, which is a blind spot of current approaches.
277 However, it is not clear at the moment how applicable this will be in larger
278 genomes as smaller fragment sizes could increase the need for sequence depth
279 beyond that of the current high-resolution Hi-C protocols. Similarly DNase I has
280 been used in place of a restriction enzyme⁴⁴. Enzymes that can cut at any point in
281 the genome could theoretically improve the resolution down to a single base
282 pair. This means that the resolution of Hi-C type experiments in mammalian
283 genomes becomes constrained by sequence depth rather than cut site density.
284 Interestingly, using DNase I did not significantly improve the resolution over
285 four-cutter restriction enzymes in mammalian cells (5 kb windows were used)
286 but it is likely that enzymes with higher cut site density will be key to improving
287 resolution in the future.

288 The Capture-C approach can be used to enrich both 3C and Hi-C libraries. The
289 combination of Capture-C and Hi-C libraries (Capture Hi-C) can exclude further
290 uninformative background from captures within the target fragment rather than
291 spanning a ligation junction. Compared to the high levels of enrichment of NG
292 Capture-C this ~2 fold benefit seems negligible for small to moderate size
293 designs (hundreds of viewpoints) and must be balanced against the more
294 extended protocol and extra losses in complexity of the library, particularly for
295 lower cell numbers. However, it is potentially beneficial when attempting
296 ambitious designs. A low-resolution design using six cutter enzymes has been
297 attempted for every annotated promoter in the human genome, however the
298 very low levels of enrichment (a mean of 10-fold) combined with the ambitious
299 design reduced the number of unique read pairs from each viewpoint, so that
300 that individual profiles are weakly sampled⁴⁵.

301 A tiled oligonucleotide capture approach has been attempted to define an all vs
302 all map of interaction at the beta globin locus⁴⁶. A potential pitfall of this
303 approach is that the efficiency of the oligonucleotide capture will partially
304 determine the interaction profile and robust methods of correcting this have not
305 been established.

306 **Chromatin Interaction Analysis by Paired-End Tag Sequencing (ChIA-PET)**

307 ChIA-PET combines chromatin immunoprecipitation (ChIP) with 3C, potentially
308 allowing all of the interactions to be identified from all of the sites bound by a
309 protein of interest. In this technique the material is initially crosslinked and then
310 sonicated, after which ChIP is performed to enrich for DNA-protein complexes.
311 DNA linkers are then ligated to the immunoprecipitated material and these are
312 used to join sequences held in close proximity. The linkers are biotinylated and

313 contain MmeI restriction sites allowing short fragments of material to be
314 extracted with a streptavidin bead pull-down. Paired-end tags (PETs) can then
315 be identified by paired-end sequencing. ChIA-PET studies have been performed
316 using ChIP for the oestrogen-receptor- α ⁴⁷, CTCF⁴⁸ and PolII⁴⁸.
317 The technique has limitations because the relatively low levels of enrichment of
318 ChIP reduce the library complexity and so the number of reads used to identify
319 individual interactions is usually extremely low (**Figure 4**). It is thought that if
320 three read pairs are present, this corresponds to a false discovery rate of <0.05 ⁴⁷
321 and this is often used to define significant interactions. By comparison the other
322 3C-based techniques define interactions based on thousands of read pairs⁴⁷. In
323 addition, ChIA-PET is considered vulnerable to bias towards detecting
324 interactions between sequences that bind the protein targeted by the antibody
325 used for the immunoprecipitation.

326

327 **Data analysis**

328 Data analysis of sequencing-based 3C approaches is complex (**Figure 6**). It is
329 difficult to compare the quality of the different methods, because customized
330 data processing is often used, which makes it challenging to relate the output
331 back to the quality of the raw data generated by the experiment.

332 Several approaches attempt to determine interactions that are significantly
333 different from a baseline interaction profile. This is often done using distance
334 modeling approaches. The background interaction profile is inversely related to
335 the distance from the viewpoint. These methods allow significant interactions to
336 stand out from the baseline profile by increasing the weight of more distal
337 interaction counts. This can be done either using a theoretical model of expected
338 baseline interactions (common in 4C) or a statistical distribution based on all of
339 the distribution of counts with distance across the whole dataset (e.g LOWESS in
340 5C). The potential problem with this approach is that very strong functional
341 interactions commonly occur between neighboring sequences in the chromatin
342 fiber and identification of these can be highly dependent on the tool or
343 parameters used. Measurement of diffusion within chromatin in live cells,
344 suggests that the contact probability decreases with the fourth power of the
345 distance along the chromatin fiber i.e. a 2-fold increase in distance results in a
346 16-fold reduction in contact probability¹⁶. Thus it seems inappropriate to
347 disregard interactions with regulatory elements close to the promoter (where
348 they are optimally placed to interact with it) relative to more distal elements that
349 have much lower contact frequencies.

350 A complementary approach that avoids distance modeling uses the tissue
351 specificity of regulatory elements to identify their activity dependent interaction
352 profile. At its most simple this relies on the comparison of the interaction profile
353 in a tissue where the target gene is inactive with a tissue in which it is active to
354 determine which interactions only appear in the active state. This approach has
355 been shown to be effective at identifying interactions with very proximal

356 regulatory elements as well as distal interactions, but would obviously ignore
357 interactions that are common between the two tissues.
358 Overall, careful thought needs to be given to the common basic steps as well as to
359 the most appropriate downstream analysis for the biological question and the
360 use of different combinations of approaches is advisable to understand how
361 robust any specific interpretation may be. It would seem that, in addition to
362 highly processed data, routine reporting of the raw read counts used to
363 substantiate any given interaction and a clear description of understood sources
364 of bias would greatly aid transparency in the field.

365 **Choice of method**

366 Considering the number of variations of the 3C-based methods, there is often
367 confusion over which method is appropriate to which question. Hi-C is unique in
368 its ability to determine genome wide (all against all) interaction profiles and
369 define the landscape of whole genomes to globally determine basic rules of
370 genome organization^{8,18,20,42}. However, in order to define the details of small
371 scale interactions that dictate regulation of individual genes Hi-C needs to be
372 performed at the highest resolution currently available, which requires several
373 billion reads per sample. Thus for investigating the detailed interactions for
374 small numbers of genes in multiple samples it would be more appropriate to use
375 4C or NG Capture-C, which can generate interaction profiles for single loci at
376 higher depth using only ~1 million reads per viewpoint. Of these two methods,
377 the more recently developed NG Capture-C approach would seem the most
378 general solution even for simple designs, considering the control of PCR
379 duplicates and ability to multiplex both samples and viewpoints simultaneously.
380 With the challenges of genome-wide association studies (GWAS) analysis in
381 mind, NG Capture-C was specifically designed to determine the regulatory
382 landscapes of hundreds of genes or regulatory elements simultaneously from the
383 small numbers of cells available from primary tissues.

384 Approaches, such as 5C, can be used for determining the all vs all chromatin
385 structure of targeted sections of the genome. Although several seminal
386 observations have been made with 5C (which predates Hi-C), the need to
387 manufacture large numbers of probes with 5C and the falling costs of sequencing
388 means that Hi-C, which was developed by the same group, is increasingly used in
389 its place.

390

391 **Future directions**

392 3C-based technologies have the potential to overcome the barriers caused by our
393 limited understanding of how genome structure relates to its function. The “all vs
394 all” technologies are driving research into the general mechanisms that dictate
395 chromatin structure, while the targeted technologies provide high-resolution
396 approaches to dissect the regulation of specific loci. Increasingly these
397 technologies are being adopted by non-specialist laboratories. It is therefore

important that the appropriate applications and limitations of each method are clear, as well as the sources of technological and bioinformatic biases that might confound correct data interpretation. A general concern for data interpretation of the 3C technologies are potential biases due to the common formaldehyde crosslinking step, which could result in a systematic skew. Though it has been shown that crosslinking does not affect the gross patterns detected¹⁸, the effects at a finer scale are unclear and this potential problem is being actively addressed. Another general source of bias is the uneven degree of digestion across the genome when restriction enzymes are used to fragment the chromatin. To date it has been impractical to effectively map these efficiencies genome-wide, particularly at high resolution, due to the depth of sequencing required. However, as sequencing costs drop it would be inevitable that such maps will be generated and the effect on 3C signal determined. As these potential biases are systematically determined, they can be incorporated for normalization into increasingly sophisticated tools for data analysis.

The other crucial development needed for accurate data interpretation is the ability to interpret 3C data in its three-dimensional context in the nucleus. At the moment 3C data are mostly represented in two dimensions, but expertise from fields such as polymer physics is increasingly being used to try to understand what the patterns seen in 3C assays are reflecting in terms of three-dimensional structures within the nucleus^{49,50}. This will be important to better understand the processes that shape chromatin structures and how they in turn may shape the activity of the genome. This drive will undoubtedly bring in other methodologies and expertise such as super-resolution imaging and live cell chromatin tagging, to view these structures dynamically at the single-cell level. These approaches will help refine and inform the models derived from the 3C data.

Irrespective of the future developments required to fully understand the relationship between nuclear structure and function, it is clear that chromatin conformation assays (especially the 4C and Capture-C approaches) have reached a level of maturity and approachability that has turned them into general tools, similar to ChIP-seq or ATAC-seq, to analyze regulatory landscapes. Although impossible to predict the full impact of these methods, investigation of mutations in non-coding regulatory or structural elements, particularly those identified in GWAS studies, will clearly benefit from such analyses.

It is now crucial that the field develops and agrees on standardized data formats and standards for quality control to promote unbiased and comparable interpretation.

Figure Legends

Figure 1. Common principles in 3C-based techniques.

- a. The chromatin fiber is initially digested into short restriction fragments, after which a ligation reaction is performed to create large DNA concatemers in which the order of the fragments reflects the three dimensional structure of the chromatin at the time of fixation.
- b. Only approximately 1 in 20 (~5%) of the restriction fragments ligate back to their original partner in a DpnII digested library¹⁵. If one assumes that the cut restriction site has the highest probability of ligating back to its original partner in the fixed nucleus, then there are at least 20 potential ligation partners for every cut site in every cell.
- c. The determined interaction frequencies from 3C-based experiments such as 4C and Capture-C can be represented as an interaction profile from a single viewpoint. All vs all interaction data generated by Hi-C like approaches are usually represented in heat maps.

Figure 2. Comparison of different 3C-based methodologies.

- 3C libraries are generated by initially crosslinking the chromatin using formaldehyde, after which a restriction enzyme digest is performed. This is followed by a ligation reaction, after which the crosslinking is reversed and the DNA extracted.
- 3C libraries can be interrogated by several different techniques. In the initial 3C experiments ligation junctions were detected by PCR amplification combined with gel electrophoresis or real time PCR.
- In 4C, 3C libraries are circularized by either optimizing the ligation reaction or performing a second restriction enzyme and ligase reaction. Inverse PCR is then performed and the product of this reaction is sequenced.
- 5C uses a highly multiplexed ligation-mediated PCR to amplify ligation junctions in 3C libraries, which are then measured by sequencing.
- Capture-C involves sonication of the 3C library, followed by ligation of barcoded sequencing adaptors, which allows multiple samples to be mixed at this stage. A hybridization reaction with biotinylated oligonucleotides targeted against the viewpoints of interest and biotin pull-down are subsequently performed (this is repeated in NG Capture-C to increase enrichment), after which the material can be sequenced.
- In the 3C library preparation in Hi-C experiments, overhanging restriction cut sites are filled in using biotinylated nucleotides, followed by a blunt-end ligation reaction. The material is sonicated and a biotin/streptavidin pull-down is performed to concentrate ligation junctions, which are then identified by sequencing.
- The ChIA-PET protocol also starts with a crosslinked and digested 3C library. Next, this material undergoes a chromatin immunoprecipitation step using an

antibody against a transcription factor or chromatin associated protein. Biotinylated linkers are ligated to the immunoprecipitated chromatin, after which these linkers are ligated to one another. An MmeI digest is then performed and the material containing the linkers is pulled down and sequenced.

Figure 3. Comparison of different 3C-based techniques.

a. The maximum resolution is determined by the frequency of cutting of the chromatin when an enzyme that cuts at fixed intervals is used (restriction enzyme with 6 bp or 4 bp recognition or MNase). Since DNase I and sonication (ChIA-PET) do not have specific cut sites they could theoretically generate very high resolution data. Techniques that report data for each ligation junction (e.g. 3C, 5C and NG Capture-C) can be analysed at the maximum resolution determined by the cut site density. However, the majority of 3C-based techniques require windowing or binning of data to generate sufficient data for meaningful profiles because this improves the sensitivity and reduces bias. This reduces the resolution and the window/bin size becomes the major determinant of resolution.

b. Comparison of the sensitivity of the different techniques. The sensitivity is largely determined by the number of reads from each individual viewpoint although techniques such as Hi-C, which undertake an all vs all approach, can improve the sensitivity markedly by combining data from multiple view points.

c. Comparison of the multiplexing ability of the different techniques. 4C, Capture-C and Hi-C are all capable of reporting genome-wide interactions. However the number of viewpoints differs radically between the methods. 4C generates genome-wide profiles from 1 viewpoint, whereas Hi-C uses an all vs all approach. The highest number of viewpoints currently used in Capture-C experiments is 450, but there is no theoretical limit and more viewpoints could be used if desired. Importantly, NG Capture-C can also analyze multiple samples simultaneously.

Figure 4. Case study 1: regulation of the alpha globin locus in erythroid cells.

The human alpha globin genes (*HBA1* and *HBA2*) are regulated by four enhancers (R1-R4) located ~10-40kb upstream of the gene within the introns of the neighboring gene *NPRL3*. These elements are characterized by DNase I hypersensitivity (see **Supplementary Fig. 1** for further chromatin data).

a. Comparison of ChIA-PET, 5C and Hi-C data in the human alpha globin locus in K562 cells. The top track shows ChIA-PET data derived from a ChIP for RNA polymerase II⁵¹. All of the reads within a 20kb region around the two alpha globin genes are included. The numbers on the reads denote the number of paired end tags (PETs) contributing to the reported interaction. Note that despite the generous inclusion criteria, the interactions with the enhancers are defined by only 16 PETs. The next profile shows 5C data from the two alpha globin promoters⁶. Note that there are limited numbers of reported fragments

because it is only possible to determine interactions between forward and reverse probes and the limited resolution of data generated by a 6-cutter restriction enzyme (HindIII). In addition the distance modeling approach downgrades the raw interaction count with the R4 regulatory element, so that it is not statistically significant. The heat map at the bottom displays Hi-C data at 5kb resolution¹⁸ (Yue Lab Hi-C data browser <http://promoter.bx.psu.edu/hi-c/>), in which the interactions with the regulatory elements are highlighted by black boxes.

b. Comparison of 4C and NG Capture-C data in the mouse alpha globin locus (*Hba-a1* and *Hba-a2*) in primary erythroid cells. The alpha globin promoters are used as viewpoints and highlighted in red. The top track shows the interactions with the regulatory elements as defined by 4C⁵². As some larger fragments fail to amplify, the profile shows gaps, despite the use of a windowing approach. The bottom track displays the NG Capture-C profile, in which the interactions between the alpha globin promoters and enhancer elements are defined by tens of thousands of unique read pairs. The profiles are an average of 4 replicates in erythroid cells and 3 ES cell replicates, which allows for comparative analysis between the active and inactive cell types and identification of statistically significant interactions when the gene is in an active state (**Supplementary Fig. 1**).

Figure 5. Case study 2: regulation of the SOX2 locus in murine ES (mES) cells.

The *Sox2* gene is situated in a less gene dense locus, with DNase I hypersensitive sites extending over 900kb away from the gene promoter. A cluster of regulatory elements is found 85-111kb from the gene promoter, which reach the criteria for a super-enhancer⁵³, however, there are several other potential regulatory elements distal to this cluster (see **Supplementary Fig, 2** for further chromatin data).

The top panel shows 5C data, which is able to delineate the significant interactions with the super-enhancer, but does not fully cover the entire locus despite the use of several hundred probes⁵⁴(GSE36203). The next panel shows Hi-C data from mES cells (40kb resolution)⁸. At this resolution it is difficult to make out the interaction with the close enhancers but the more distal interactions are more prominent (the interactions defined by 4C and Capture-C are highlighted by black boxes). The 4C (GSM1868926)⁵⁵ and NG Capture-C (GSE67959)¹⁵ profiles at the bottom, from the *Sox2* viewpoint highlighted in red, both show interactions across the gene desert.

568 **Figure 6. Flow diagram of steps required for analysis of high-throughput**
569 **3C-based technologies.**

570 3C-based experiments are usually sequenced using paired-end sequencing.
571 Standard quality control measures and trimming of adapter sequences need to
572 be performed prior to further analysis. When the DNA fragment is smaller than
573 the size of the combined paired-end reads, the central area of overlap can be
574 used to combine the reads into one single read.

575 Reads then need to be separated into the component fragments that ligated to
576 one another at the restriction enzyme cut site. The ligation junction can be
577 determined directly when the restriction cut site has been sequenced. However,
578 it is common to infer that there is a restriction enzyme cut site in the central
579 unsequenced portion of the read. This has the disadvantage that there may be
580 other fragments interposed in the central un-sequenced part of the read.

581 After the reads have been split into the component fragments they are aligned to
582 the genome separately. This can be challenging if the reads are not of sufficient
583 length to allow both parts of the read to be aligned properly (this can be
584 problematic with short single-end reads). The reads build up at the ends of the
585 restriction fragments and so the data need further processing to generate a
586 meaningful signal. The reads can either be mapped to the restriction fragments
587 or data from multiple fragments can be combined either into bins or with a
588 moving window.

589 When possible the data should be filtered to remove PCR duplicates. It is also
590 important to remove signal caused by poor restriction enzyme digestion, mis-
591 mapping to repetitive DNA sequences and off-target capture and mis-priming
592 (depending on the method used).

593 Finally, statistical analysis can be performed. When possible, areas of increased
594 interaction on gene activation can be identified by comparing the interaction
595 profiles in active and inactive states. It is also possible to subtract the inactive
596 background distribution using mathematical modeling, though this is less
597 accurate and not as straightforward to interpret.

598
599

600 References

- 601
- 602 1 Mikkelsen, T. S. *et al.* Genome-wide maps of chromatin state in
603 pluripotent and lineage-committed cells. *Nature* **448**, 553-560,
604 doi:10.1038/nature06008 (2007).
- 605 2 Robertson, G. *et al.* Genome-wide profiles of STAT1 DNA association using
606 chromatin immunoprecipitation and massively parallel sequencing.
607 *Nature methods* **4**, 651-657, doi:10.1038/nmeth1068 (2007).
- 608 3 Hesselberth, J. R. *et al.* Global mapping of protein-DNA interactions in vivo
609 by digital genomic footprinting. *Nature methods* **6**, 283-289,
610 doi:10.1038/nmeth.1313 (2009).
- 611 4 Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y. & Greenleaf, W. J.
612 Transposition of native chromatin for fast and sensitive epigenomic
613 profiling of open chromatin, DNA-binding proteins and nucleosome
614 position. *Nature methods* **10**, 1213-1218, doi:10.1038/nmeth.2688
615 (2013).
- 616 5 Stamatoyannopoulos, J. Connecting the regulatory genome. *Nature*
617 *genetics* **48**, 479-480, doi:10.1038/ng.3553 (2016).
- 618 6 Sanyal, A., Lajoie, B. R., Jain, G. & Dekker, J. The long-range interaction
619 landscape of gene promoters. *Nature* **489**, 109-113,
620 doi:10.1038/nature11279 (2012).
- 621 7 Bulger, M. & Groudine, M. Functional and mechanistic diversity of distal
622 transcription enhancers. *Cell* **144**, 327-339,
623 doi:10.1016/j.cell.2011.01.024 (2011).
- 624 8 Dixon, J. R. *et al.* Topological domains in mammalian genomes identified
625 by analysis of chromatin interactions. *Nature* **485**, 376-380,
626 doi:10.1038/nature11082 (2012).
- 627 9 Nora, E. P. *et al.* Spatial partitioning of the regulatory landscape of the X-
628 inactivation centre. *Nature* **485**, 381-385, doi:10.1038/nature11049
629 (2012).
- 630 10 Tolhuis, B., Palstra, R. J., Splinter, E., Grosveld, F. & de Laat, W. Looping
631 and interaction between hypersensitive sites in the active beta-globin
632 locus. *Molecular cell* **10**, 1453-1465 (2002).
- 633 11 Tan-Wong, S. M. *et al.* Gene loops enhance transcriptional directionality.
634 *Science* **338**, 671-675, doi:10.1126/science.1224350 (2012).
- 635 12 Hsieh, T. H. *et al.* Mapping Nucleosome Resolution Chromosome Folding
636 in Yeast by Micro-C. *Cell* **162**, 108-119, doi:10.1016/j.cell.2015.05.048
637 (2015).
- 638 13 Mukherjee, S., Erickson, H. & Bastia, D. Detection of DNA looping due to
639 simultaneous interaction of a DNA-binding protein with two spatially
640 separated binding sites on DNA. *Proceedings of the National Academy of*
641 *Sciences of the United States of America* **85**, 6287-6291 (1988).
- 642 14 Dekker, J., Rippe, K., Dekker, M. & Kleckner, N. Capturing chromosome
643 conformation. *Science* **295**, 1306-1311, doi:10.1126/science.1067799
644 (2002).
- 645 15 Davies, J. O. *et al.* Multiplexed analysis of chromosome conformation at
646 vastly improved sensitivity. *Nature methods* **13**, 74-80,
647 doi:10.1038/nmeth.3664 (2016).

648 16 Lucas, J. S., Zhang, Y., Dudko, O. K. & Murre, C. 3D trajectories adopted by
649 coding and regulatory DNA elements: first-passage times for genomic
650 interactions. *Cell* **158**, 339-352, doi:10.1016/j.cell.2014.05.036 (2014).

651 17 Nagano, T. *et al.* Single-cell Hi-C reveals cell-to-cell variability in
652 chromosome structure. *Nature* **502**, 59-64, doi:10.1038/nature12593
653 (2013).

654 18 Rao, S. S. *et al.* A 3D map of the human genome at kilobase resolution
655 reveals principles of chromatin looping. *Cell* **159**, 1665-1680,
656 doi:10.1016/j.cell.2014.11.021 (2014).

657 19 Nagano, T. *et al.* Comparison of Hi-C results using in-solution versus in-
658 nucleus ligation. *Genome biology* **16**, 175, doi:10.1186/s13059-015-0753-
659 7 (2015).

660 20 Lieberman-Aiden, E. *et al.* Comprehensive mapping of long-range
661 interactions reveals folding principles of the human genome. *Science* **326**,
662 289-293, doi:10.1126/science.1181369 (2009).

663 21 Spilianakis, C. G. & Flavell, R. A. Long-range intrachromosomal
664 interactions in the T helper type 2 cytokine locus. *Nature immunology* **5**,
665 1017-1027, doi:10.1038/ni1115 (2004).

666 22 Liu, Z. & Garrard, W. T. Long-range interactions between three
667 transcriptional enhancers, active κ gene promoters, and a 3'
668 boundary sequence spanning 46 kilobases. *Molecular and cellular biology*
669 **25**, 3220-3231, doi:10.1128/MCB.25.8.3220-3231.2005 (2005).

670 23 Vernimmen, D. *et al.* Chromosome looping at the human alpha-globin
671 locus is mediated via the major upstream regulatory element (HS -40).
672 *Blood* **114**, 4253-4260, doi:10.1182/blood-2009-03-213439 (2009).

673 24 Vernimmen, D., De Gobbi, M., Sloane-Stanley, J. A., Wood, W. G. & Higgs, D.
674 R. Long-range chromosomal interactions regulate the timing of the
675 transition between poised and active gene expression. *The EMBO journal*
676 **26**, 2041-2051, doi:10.1038/sj.emboj.7601654 (2007).

677 25 Zhou, G. L. *et al.* Active chromatin hub of the mouse alpha-globin locus
678 forms in a transcription factory of clustered housekeeping genes.
679 *Molecular and cellular biology* **26**, 5096-5105, doi:10.1128/MCB.02454-
680 05 (2006).

681 26 Zhao, Z. *et al.* Circular chromosome conformation capture (4C) uncovers
682 extensive networks of epigenetically regulated intra- and
683 interchromosomal interactions. *Nature genetics* **38**, 1341-1347,
684 doi:10.1038/ng1891 (2006).

685 27 Simonis, M. *et al.* Nuclear organization of active and inactive chromatin
686 domains uncovered by chromosome conformation capture-on-chip (4C).
687 *Nature genetics* **38**, 1348-1354, doi:10.1038/ng1896 (2006).

688 28 Stadhouders, R. *et al.* Multiplexed chromosome conformation capture
689 sequencing for rapid genome-scale high-resolution detection of long-
690 range chromatin interactions. *Nature protocols* **8**, 509-524,
691 doi:10.1038/nprot.2013.018 (2013).

692 29 Simonis, M., Kooren, J. & de Laat, W. An evaluation of 3C-based methods to
693 capture DNA interactions. *Nature methods* **4**, 895-901,
694 doi:10.1038/nmeth1114 (2007).

695 30 Schwartzman, O. *et al.* UMI-4C for quantitative and targeted chromosomal
696 contact profiling. *Nature methods* **13**, 685-691, doi:10.1038/nmeth.3922
697 (2016).

698 31 Andrey, G. *et al.* A switch between topological domains underlies HoxD
699 genes collinearity in mouse limbs. *Science* **340**, 1234167,
700 doi:10.1126/science.1234167 (2013).

701 32 de Wit, E. *et al.* The pluripotent genome in three dimensions is shaped
702 around pluripotency factors. *Nature* **501**, 227-231,
703 doi:10.1038/nature12420 (2013).

704 33 Pasquali, L. *et al.* Pancreatic islet enhancer clusters enriched in type 2
705 diabetes risk-associated variants. *Nature genetics* **46**, 136-143,
706 doi:10.1038/ng.2870 (2014).

707 34 Simonis, M. *et al.* High-resolution identification of balanced and complex
708 chromosomal rearrangements by 4C technology. *Nature methods* **6**, 837-
709 842, doi:10.1038/nmeth.1391 (2009).

710 35 Groschel, S. *et al.* A single oncogenic enhancer rearrangement causes
711 concomitant EVI1 and GATA2 deregulation in leukemia. *Cell* **157**, 369-
712 381, doi:10.1016/j.cell.2014.02.019 (2014).

713 36 Dostie, J. *et al.* Chromosome Conformation Capture Carbon Copy (5C): a
714 massively parallel solution for mapping interactions between genomic
715 elements. *Genome research* **16**, 1299-1309, doi:10.1101/gr.5571506
716 (2006).

717 37 Hughes, J. R. *et al.* Analysis of hundreds of cis-regulatory landscapes at
718 high resolution in a single, high-throughput experiment. *Nature genetics*,
719 doi:10.1038/ng.2871 (2014).

720 38 Yaffe, E. & Tanay, A. Probabilistic modeling of Hi-C contact maps
721 eliminates systematic biases to characterize global chromosomal
722 architecture. *Nature genetics* **43**, 1059-1065, doi:10.1038/ng.947 (2011).

723 39 Imakaev, M. *et al.* Iterative correction of Hi-C data reveals hallmarks of
724 chromosome organization. *Nature methods* **9**, 999-1003,
725 doi:10.1038/nmeth.2148 (2012).

726 40 Hu, M. *et al.* HiCNorm: removing biases in Hi-C data via Poisson
727 regression. *Bioinformatics* **28**, 3131-3133,
728 doi:10.1093/bioinformatics/bts570 (2012).

729 41 Jin, F. *et al.* A high-resolution map of the three-dimensional chromatin
730 interactome in human cells. *Nature* **503**, 290-294,
731 doi:10.1038/nature12644 (2013).

732 42 Naumova, N. *et al.* Organization of the mitotic chromosome. *Science* **342**,
733 948-953, doi:10.1126/science.1236083 (2013).

734 43 Zhang, Y. *et al.* Spatial organization of the mouse genome and its role in
735 recurrent chromosomal translocations. *Cell* **148**, 908-921,
736 doi:10.1016/j.cell.2012.02.002 (2012).

737 44 Ma, W. *et al.* Fine-scale chromatin interaction maps reveal the cis-
738 regulatory landscape of human lincRNA genes. *Nature methods* **12**, 71-78,
739 doi:10.1038/nmeth.3205 (2015).

740 45 Schoenfelder, S. *et al.* The pluripotent regulatory circuitry connecting
741 promoters to their long-range interacting elements. *Genome research* **25**,
742 582-597, doi:10.1101/gr.185272.114 (2015).

- 46 Kolovos, P. *et al.* Targeted Chromatin Capture (T2C): a novel high resolution high throughput method to detect genomic interactions and regulatory elements. *Epigenetics & chromatin* **7**, 10, doi:10.1186/1756-8935-7-10 (2014).
- 47 Fullwood, M. J. *et al.* An oestrogen-receptor-alpha-bound human chromatin interactome. *Nature* **462**, 58-64, doi:10.1038/nature08497 (2009).
- 48 Kieffer-Kwon, K. R. *et al.* Interactome maps of mouse gene regulatory domains reveal basic principles of transcriptional regulation. *Cell* **155**, 1507-1520, doi:10.1016/j.cell.2013.11.039 (2013).
- 49 Brackley, C. A. *et al.* Predicting the three-dimensional folding of cis-regulatory regions in mammalian genomes using bioinformatic data and polymer models. *Genome biology* **17**, 59, doi:10.1186/s13059-016-0909-0 (2016).
- 50 Fraser, J. *et al.* Hierarchical folding and reorganization of chromosomes are linked to transcriptional changes in cellular differentiation. *Molecular systems biology* **11**, 852, doi:10.15252/msb.20156492 (2015).
- 51 Li, G. *et al.* Extensive promoter-centered chromatin interactions provide a topological basis for transcription regulation. *Cell* **148**, 84-98, doi:10.1016/j.cell.2011.12.014 (2012).
- 52 van de Werken, H. J. *et al.* Robust 4C-seq data analysis to screen for regulatory DNA interactions. *Nature methods* **9**, 969-972, doi:10.1038/nmeth.2173 (2012).
- 53 Whyte, W. A. *et al.* Master transcription factors and mediator establish super-enhancers at key cell identity genes. *Cell* **153**, 307-319, doi:10.1016/j.cell.2013.03.035 (2013).
- 54 Phillips-Cremins, J. E. *et al.* Architectural protein subclasses shape 3D organization of genomes during lineage commitment. *Cell* **153**, 1281-1295, doi:10.1016/j.cell.2013.04.053 (2013).
- 55 de Wit, E. *et al.* CTCF Binding Polarity Determines Chromatin Looping. *Molecular cell* **60**, 676-684, doi:10.1016/j.molcel.2015.09.023 (2015).

Annotation of references

This is the seminal paper first describing chromosome conformation capture in yeast¹⁴.

This is the first use of 3C to define interactions between regulatory elements in mammalian cells¹⁰.

This paper first described the Hi-C method and it describes large scale organization of chromatin as a fractal globule²⁰.

790

791 This paper uses 4C to delineate the changes in gene regulation and interaction
792 profiles at the HoxD genes during limb development³¹.

793

794 This paper describes the highest possible resolution currently achievable with
795 genome-wide all vs all approaches using Hi-C¹⁸.

796

797 This paper describes the highest resolution and sensitivity available with a one
798 vs all approach using NG Capture-C and is additionally capable of high levels of
799 multiplexing of viewpoints¹⁵.

800

801

Figure 1

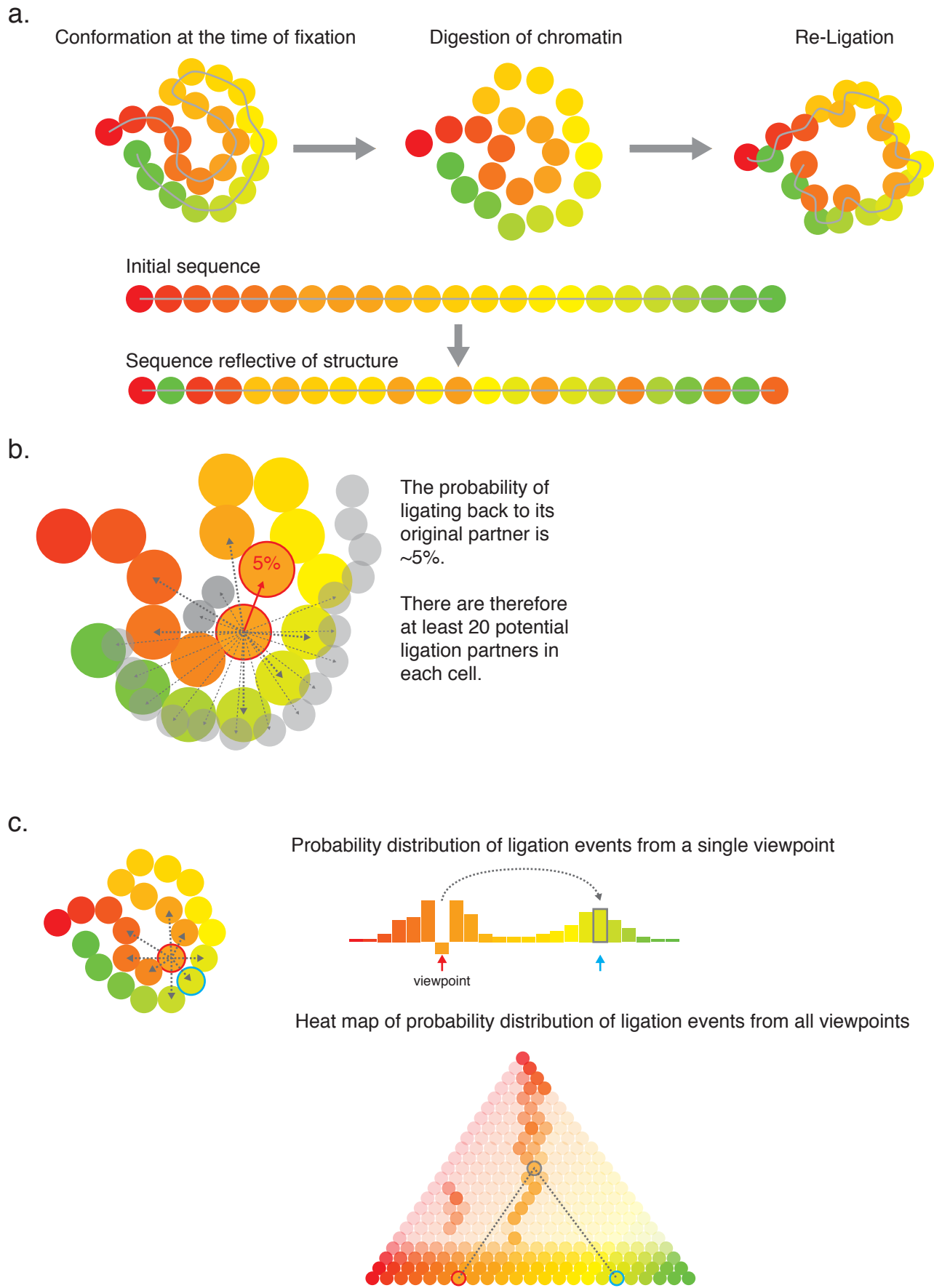


Figure 2

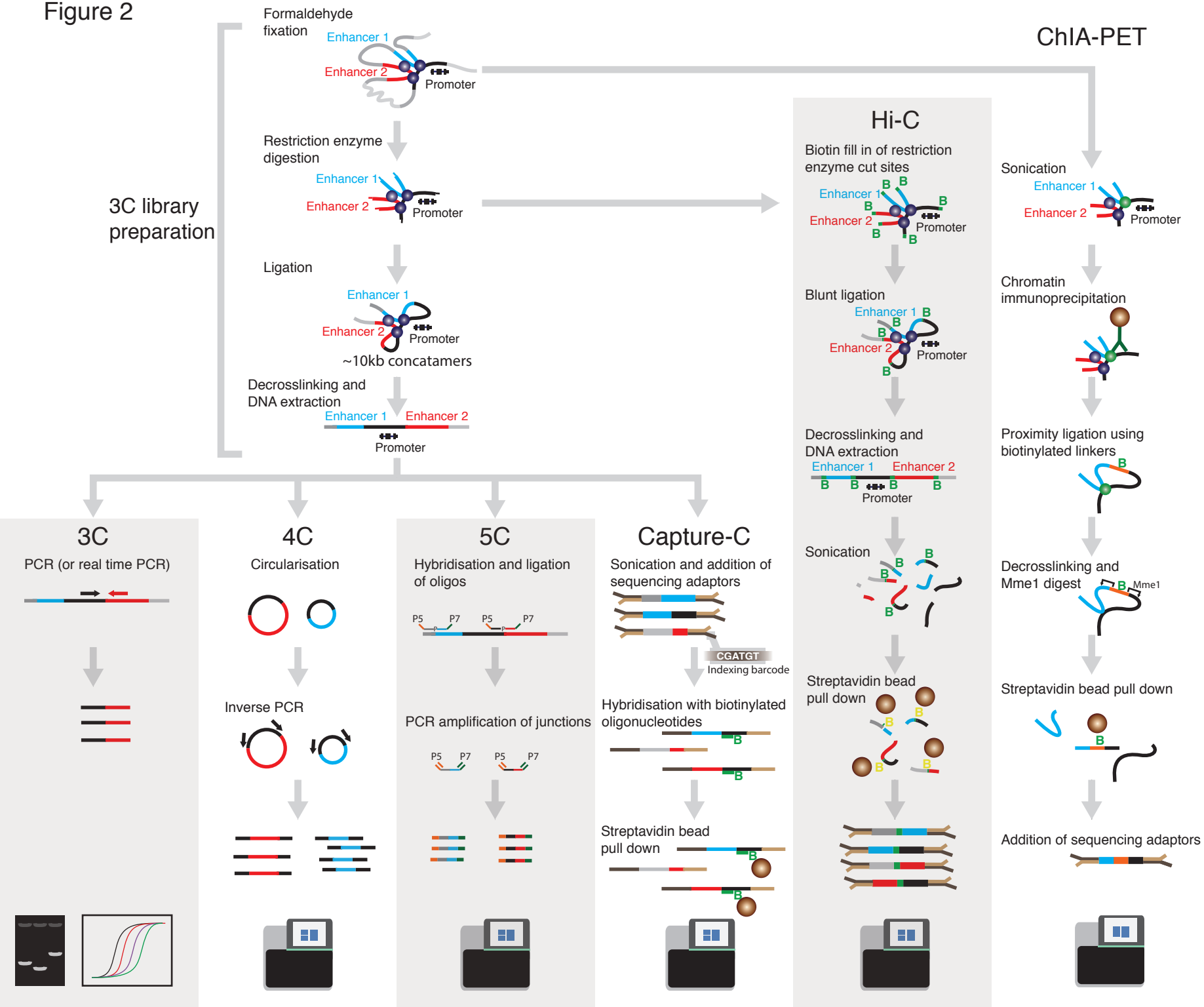
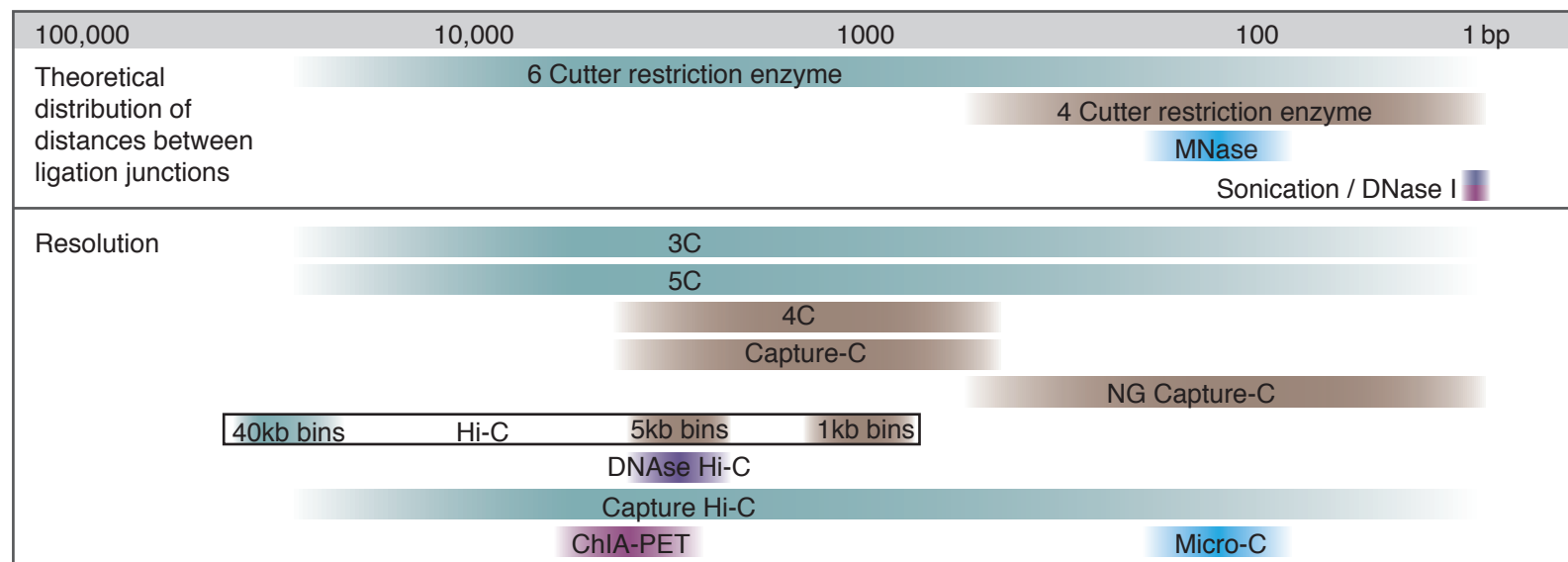
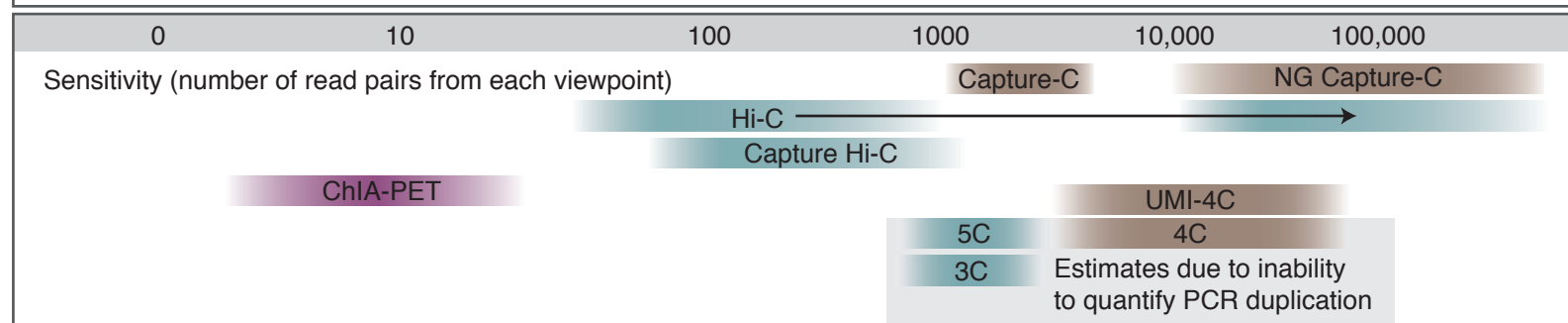


Figure 3

a.



b.



c.

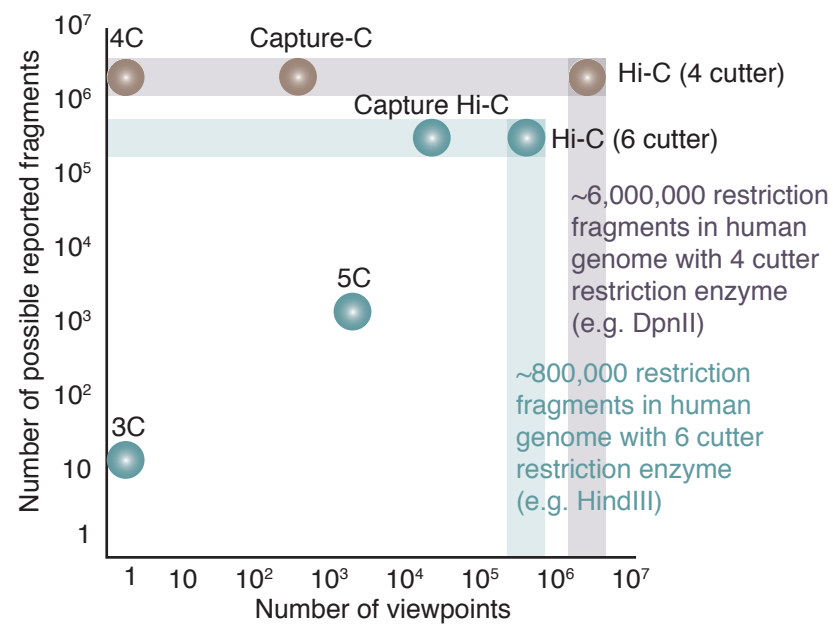


Figure 4

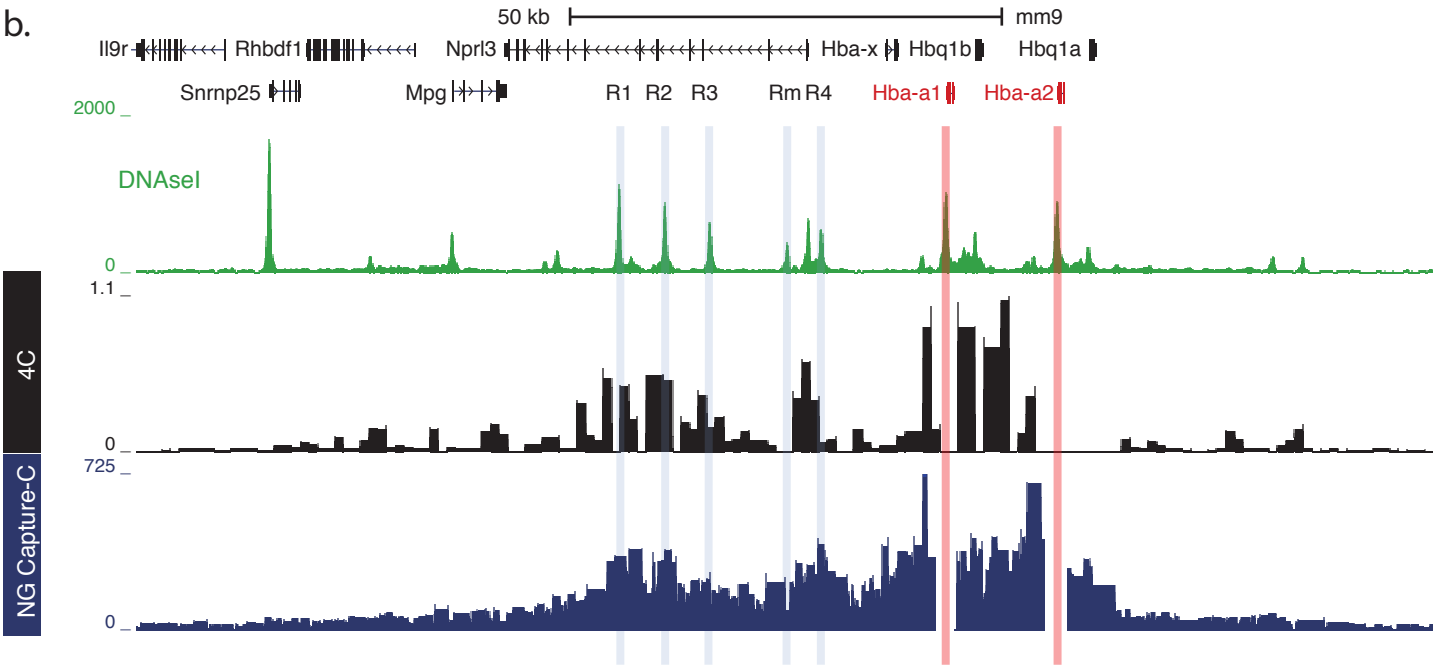
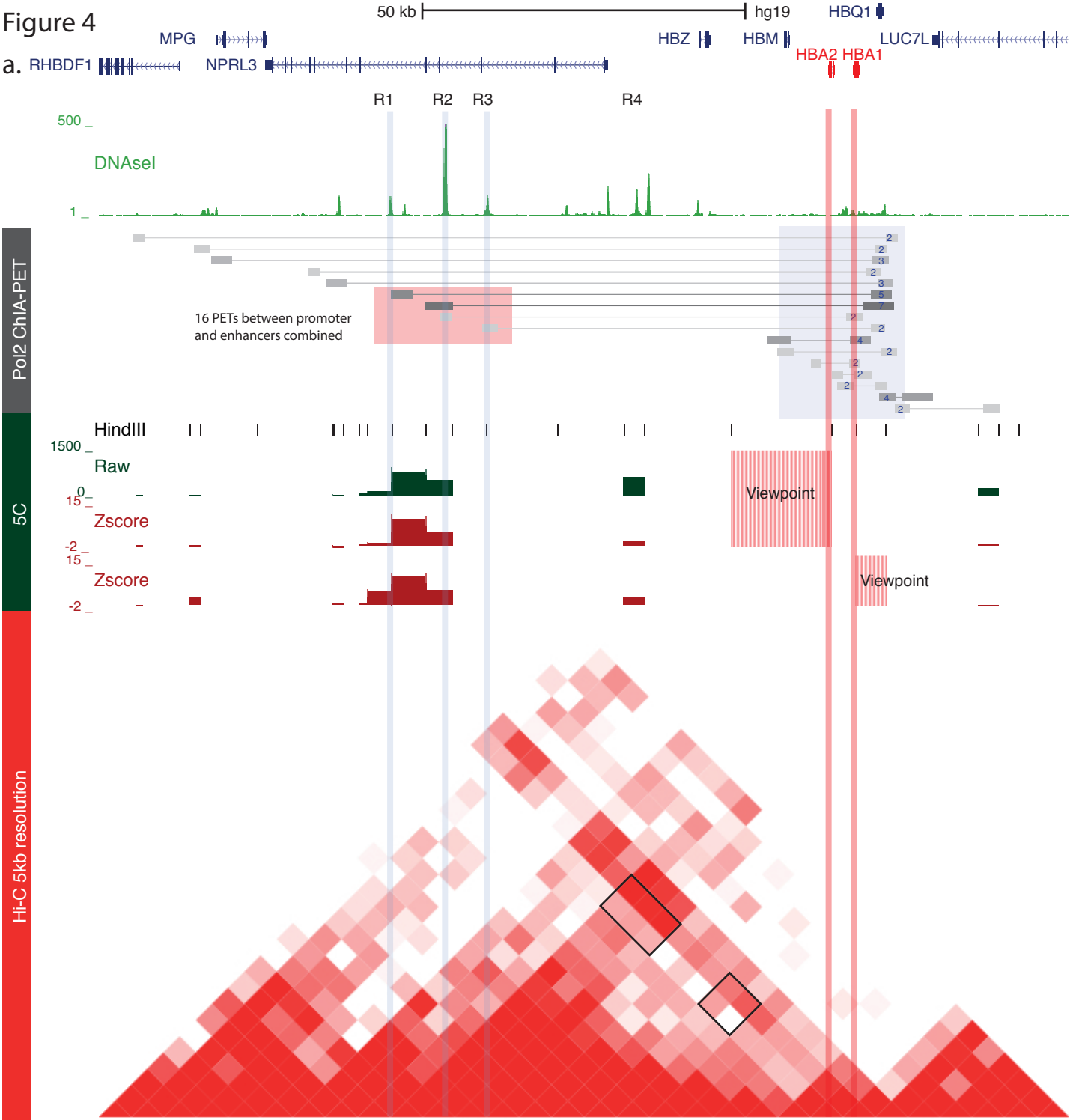


Figure 5

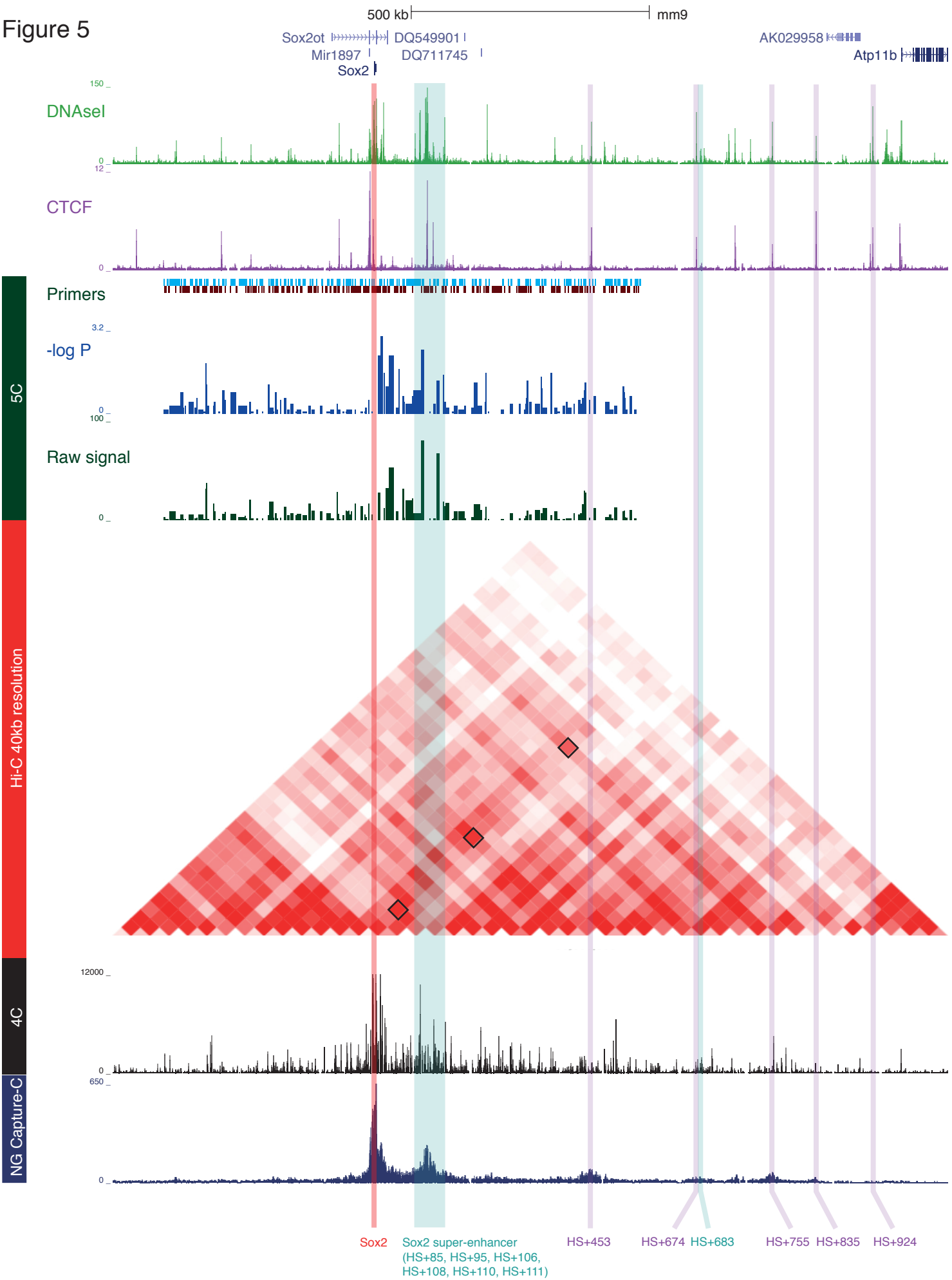
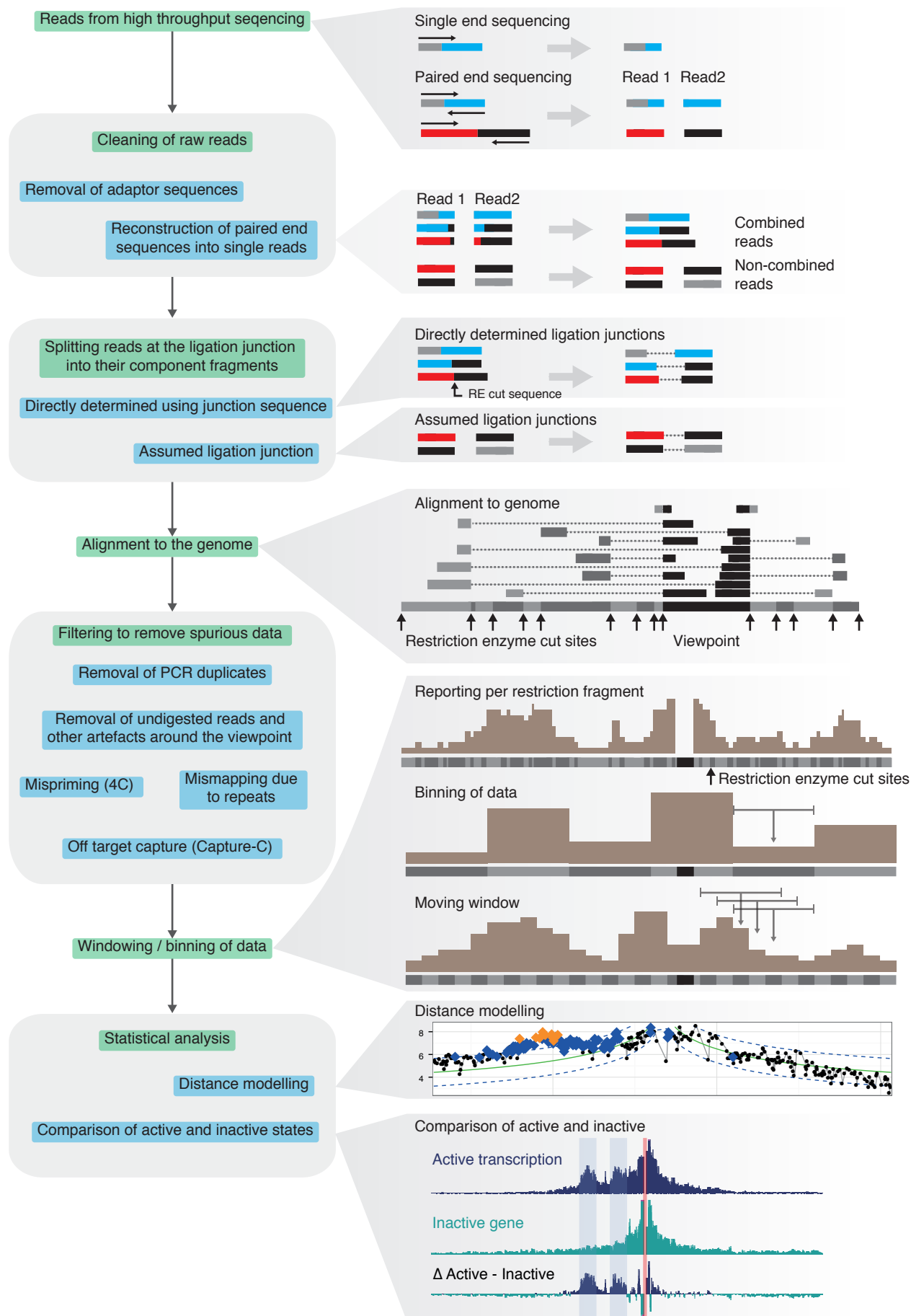
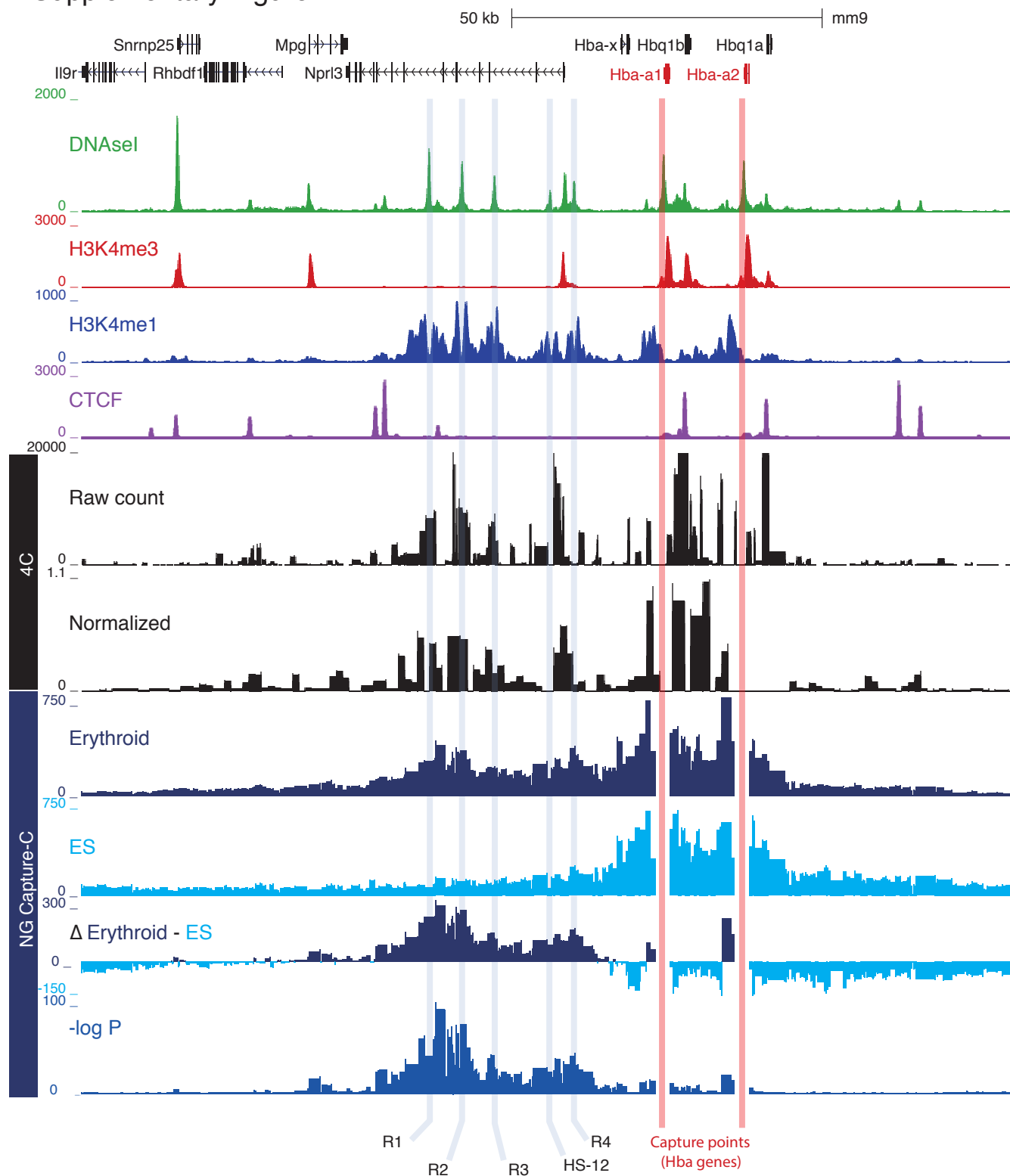


Figure 6



Supplementary Figure 1

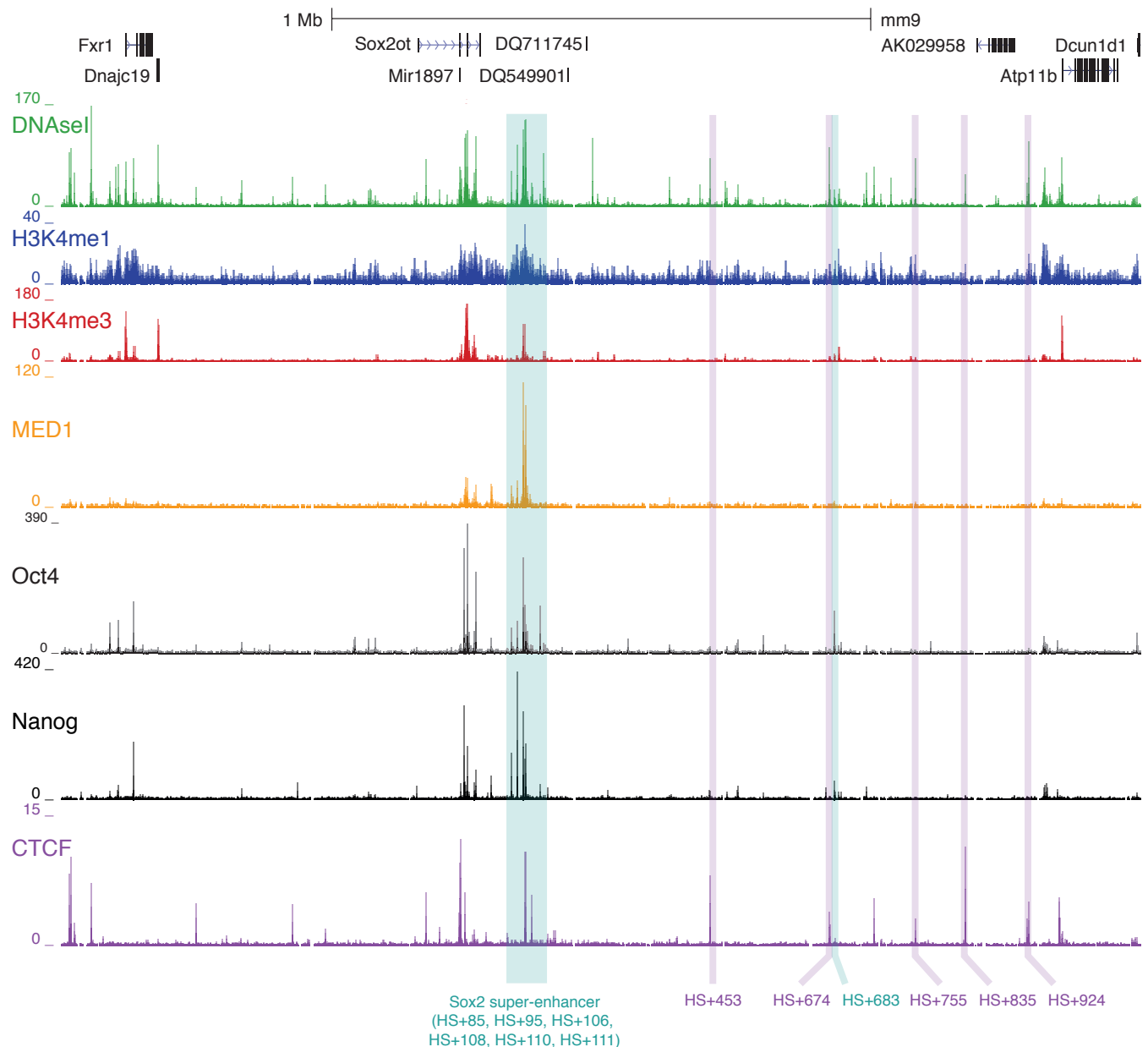


Supplementary Figure 1: Chromatin, 4C and NG Capture-C data for the alpha globin locus in mouse erythroid cells.

The DNaseI hypersensitivity track (green), marking open chromatin, is shown at the top. Below are ChIP-seq profiles for H3K4me3 (red), H3K4me1 (blue), and CTCF (purple), that highlight promoters, enhancers and CTCF binding sites, respectively (Hughes et al., 2014).

The raw and normalized 4C profiles are shown below (Van der Werken et al., 2012). The NG Capture-C profiles are an average of 4 replicates in erythroid cells and 3 ES cell replicates, which allows for comparative analysis between cell types when the gene is active (erythroid) and inactive (ES) and identification of statistically significant interactions when the gene is in an active state (Davies et al, 2016).

Supplementary Figure 2



Supplementary Figure 2: Chromatin data for the Sox2 locus in mES cells.

The DNaseI hypersensitivity track (green), marking open chromatin, is shown at the top. Below are ChIP-seq profiles for H3K4me1 (blue), H3K4me3 (red), mediator (MED, yellow), octamer-binding transcription factor 4 (Oct4, black), Nanog (black) and CTCF (purple), which allow the DNaseI hypersensitive sites to be defined further.

A cluster of hypersensitive sites 85-111 kb (highlighted in green) from the promoter of the gene are defined as a super-enhancer (Whyte et al., 2013). However, the gene has several other cell-type specific regulatory elements, which extend nearly 1 Mb from the promoter. These tend to be CTCF binding sites (highlighted in purple), although there is an additional potential regulatory element bound by Nanog and Oct4 (HS+683, highlighted in green).

ES cell data: DNaseI (ENCODE UW); ChIP-seq H3K4me1 and H3K4me3 (ENCODE/LICR); ChIP-seq MED1 (Young lab GSM1038259); ChIP-seq Oct4 (Young lab GSM1082340); ChIP-seq Nanog (Young lab GSM1082342), ChIP-seq CTCF (LICR GSM918748).