

1 **Robust estimation of Hi-C contact matrices by fused lasso reveals preferential**
2 **insulation of super-enhancers by strong TAD boundaries and a synergistic role in**
3 **cancer**

4 Yixiao Gong^{1,2,+}, Charalampos Lazaris^{1,2,+}, Aurelie Lozano³, Prabhanjan Kambadur⁴,
5 Panagiotis Ntziachristos⁵, Iannis Aifantis^{1,2,*}, Aristotelis Tsirigos^{1,2,6,*}

6 ¹ Department of Pathology, NYU School of Medicine, New York, NY 10016, USA

7 ² Laura and Isaac Perlmutter Cancer Center and Helen L. and Martin S. Kimmel Center for
8 Stem Cell Biology, NYU School of Medicine, New York, NY 10016, USA

9 ³ Center for Computational and Statistical Learning, IBM T.J. Watson Research Center, NY
10 10598, USA

11 ⁴ Bloomberg LP, 731 Lexington Avenue, New York City, NY, USA

12 ⁵ Department of Biochemistry and Molecular Genetics, Feinberg School of Medicine,
13 Northwestern University, Chicago, IL 60611, USA

14 ⁶ Applied Bioinformatics Laboratories, NYU School of Medicine, NY 10016, USA

15 + Equal contribution

16

17 * To whom correspondence should be addressed. Tel: +16465012693; Email:

18 Aristotelis.Tsirigos@nyumc.org; Correspondence may also be addressed to Iannis Aifantis.

19 Tel: +1 212 263 9898, Fax: +1 212 263 9210, E-mail: Ioannis.Aifantis@nyumc.org

20

21

22 **ABSTRACT**

23 The metazoan genome is compartmentalized in megabase-scale areas of highly interacting
24 chromatin known as topologically associating domains (TADs), typically identified by
25 computational analyses of Hi-C sequencing data. TADs are demarcated by boundaries that
26 are largely conserved across cell types and even across species, although, increasing
27 evidence suggests that the seemingly invariant TAD boundaries may exhibit plasticity and
28 their insulating strength can vary. However, a genome-wide characterization of TAD boundary
29 strength in mammals is still lacking. A systematic classification and characterization of TAD
30 boundaries may generate new insights into their function. In this study, we use fused two-
31 dimensional lasso as a machine-learning method to first improve Hi-C contact matrix
32 reproducibility, and, subsequently, categorize TAD boundaries based on their strength. We
33 demonstrate that increased boundary strength is associated with elevated CTCF levels and
34 that TAD boundary insulation scores may differ across cell types. Intriguingly, we observed
35 that super-enhancer elements are preferentially insulated by strong boundaries. Furthermore,
36 a pan-cancer analysis revealed that strong TAD boundaries and super-enhancer elements are
37 frequently co-duplicated. Taken together, our findings suggest that super-enhancers insulated
38 by strong TAD boundaries may be exploited, as a functional unit, by cancer cells to promote
39 oncogenesis.

40

41 **INTRODUCTION**

42 The advent of proximity-based ligation assays has allowed us to probe the three-dimensional
43 chromatin organization at an unprecedented resolution [1, 2]. Hi-C, a high-throughput
44 chromosome conformation variant, has enabled genome-wide identification of chromatin-
45 chromatin interactions [3]. Hi-C has revealed that the metazoan genome is organized in areas
46 of active and inactive chromatin known as A and B compartments respectively [3]. These are
47 further compartmentalized in super-TADs [4], topologically associating domains (TADs) [5–7]
48 and sub-TADs [8], as well as gene neighbourhoods [9]. Several algorithms have been already
49 developed to reveal this hierarchical chromatin organization, including Directionality Index (DI)

50 [5], Armatus [10], TADtree [11], Insulation Index (Crane) [12], IC-Finder [13] and others.
51 However, none of these studies has systematically explored the properties of the hierarchical
52 organization of TADs. Additionally, although TADs are seemingly invariant, mounting evidence
53 suggests that TAD boundaries can vary in strength, ranging from permissive TAD boundaries
54 that allow more inter-TAD interactions to more rigid (strong) boundaries that clearly demarcate
55 adjacent TADs [14]. Recent studies have shown that in *Drosophila*, exposure to heat-shock
56 caused local changes in certain TAD boundaries resulting in TAD merging [15]. A recent study
57 showed that during motor neuron (MN) differentiation in mammals, TAD and sub-TAD
58 boundaries in the *Hox* cluster are not rigid and their plasticity is linked to changes in gene
59 expression during differentiation [16]. It has also been demonstrated that boundary strength
60 is positively associated with the occupancy of structural proteins including CCCTC-binding
61 factor (CTCF) [5]. Despite the fact that there is a handful of studies demonstrating that TAD
62 boundaries can vary in strength in organisms like *Drosophila*, no study has yet addressed the
63 issue of boundary strength in mammals and how it may be related to potential boundary
64 disruptions and aberrant gene activation in cancer. Here we introduce a new method based
65 on fused two-dimensional lasso [17] in order to: (a) robustly estimate Hi-C contact matrices,
66 (b) categorize TAD boundaries based on their insulating strength, (c) characterize TAD
67 boundaries in terms of CTCF binding and other functional elements, and (d) investigate
68 potential genetic alterations of TAD boundaries in cancer. We anticipate that our study will
69 help generate new insights into the significance of TAD boundaries.

70

71 **MATERIALS AND METHODS**

72 **Comprehensive re-analysis of published high-resolution Hi-C datasets**

73 In order to develop a method that successfully handles variation in Hi-C data and improves
74 reproducibility, we carefully selected our Hi-C datasets to represent technical variation due to
75 the execution of the experiments by different laboratories and/or the usage of different
76 restriction enzymes. We identified publicly available human Hi-C datasets that fulfilled the
77 following criteria: (i) availability of two biological replicates and (ii) sufficient sequencing depth

78 to robustly identify topologically-associating domains (TADs) as described in our TAD calling
79 benchmark study [18]. Specifically, we ensured that our datasets included samples with at
80 least ~40 million intra-chromosomal read pairs and that the Hi-C experiment was performed
81 in biological replicates, either by using one restriction enzyme (HindIII or Mbol) (H1 cells and
82 their derivatives [19], K562, KBM7 and NHEK cells [20] and in-house generated CUTLL-1), or
83 two enzymes (HindIII or Mbol) (GM12878 [20], IMR90 [5, 21]), in order to examine the
84 consistency of predicted Hi-C interactions across different enzymes. All datasets were then
85 comprehensively re-analysed using our HiC-bench platform [18]. Quality assessment analysis
86 revealed that the samples varied considerably in terms of total numbers of reads, ranging from
87 ~150 million reads to more than 1.3 billion (**Supplementary Figure 1a**). Mappable reads were
88 over 96% in all samples. The percentages of total accepted reads corresponding to *cis* (ds-
89 accepted-intra, dark green) and *trans* (ds-accepted-inter, light green) (**Supplementary Figure**
90 **1b**) also varied widely, ranging from ~17% to ~56%. Duplicate read pairs (*ds-duplicate-intra*
91 and *ds-duplicate-inter*, red and pink respectively), non-uniquely mappable (*multihit*; light blue),
92 single-end mappable (*single-sided*; dark blue) and unmapped reads (*unmapped*; dark purple)
93 were discarded. Self-ligation products (*ds-same-fragment*; orange) and reads mapping too far
94 (*ds-too-far*, light purple) from restriction sites or too close to one another (*ds-too-close*; orange)
95 were also discarded. Only double-sided uniquely mappable *cis* (*ds-accepted-intra*; dark green)
96 and *trans* (*ds-accepted-inter*; light green) read pairs were used for downstream analysis.
97 Despite the differences in sequencing depth and in the percentages of useful reads across
98 samples, all samples had enough useful reads for TAD. However, due to the wide differences
99 in sequencing depth, and to ensure fair comparisons of Hi-C matrices in this study, all datasets
100 were down-sampled such that the number of usable intra-chromosomal reads pairs was ~40
101 million for each replicate. Finally, to study the effect of sequencing depth, we also resampled
102 at ~80 and ~120 million read pairs, by limiting our evaluation to those samples that had
103 adequate sequencing depth.

104

105

106 **Scaled Hi-C contact matrices**

107 Hi-C contact matrices were scaled by: (a) the total number of (usable) intra-chromosomal read
108 pairs, and (b) the “effective length” of the corresponding pair of interacting bins [22]. More
109 specifically, the scaled Hi-C count corresponding to interactions between the Hi-C matrix bins
110 i, j (y_{ij}) is defined by the formula:

$$111 \quad y_{ij} = \frac{x_{ij}}{eff_i \cdot eff_j \cdot N}$$

112 where x_{ij} is the original number of interactions between the bins i and j , eff_i the effective length
113 for the bin i , eff_j the effective length for the bin j , and N is the total number of read pairs.

114

115 **Distance-normalized Hi-C contact matrices**

116 Genomic loci that are further apart in terms of linear distance on DNA tend to give fewer
117 interactions in Hi-C maps than loci that are closer. For intra-chromosomal interactions, this
118 effect of genomic distance should be taken into account. Consequently, the interactions were
119 distance-normalized using a z-score that was calculated taking into account the mean Hi-C
120 counts for all interactions at a given distance d and the corresponding standard deviation.
121 Thus, the z-score for the interaction between the Hi-C contact matrix bins i and j (z_{ij}) is given
122 the following equation:

$$123 \quad z_{ij} = \frac{y_{ij} - \mu(d)}{\sigma(d)}$$

124 where y_{ij} corresponds to the number of interactions between the bins i and j , $\mu(d)$ to the mean
125 (expected) number of interactions for distance $d=|j-i|$ and $\sigma(d)$ is the corresponding standard
126 deviation of the mean.

127 **Fused two-dimensional lasso**

128 While our naïve scaling approach successfully increased the cross-enzyme and same-
129 enzyme correlation of Hi-C matrices, we sought to improve the correlation even further. We
130 used two-dimensional lasso, an optimization machine learning technique widely used to
131 analyse noisy datasets, especially images [17]. This technique is very-well suited for

132 identifying topological domains based on contact maps generated by Hi-C sequencing
133 experiments for two reasons: (a) Hi-C datasets are inherently noisy, and (b) topological
134 domains are continuous DNA segments of highly interacting loci that would represent solid
135 squares along the diagonal of Hi-C contact matrices. Topological domains map to squares of
136 different length along the diagonal of the Hi-C contact matrix, but they are not solid as they
137 contain several gaps, i.e. scattered regions on those squares that show little or no interaction.
138 Two-dimensional fused lasso [23] addresses the issue by penalizing differences between
139 neighbouring elements in the contact matrix. This is achieved by the penalty parameter λ
140 (lambda), as described in the equation:

$$141 \quad \hat{\beta} = \underset{\beta \in \mathbb{R}^n}{\operatorname{argmin}} \frac{1}{2} \sum_{i=1}^n (y_i - \beta_i)^2 + \lambda \sum_{(i,j) \in E} |\beta_i - \beta_j| ,$$

142 where y is the original (i.e. observed) contact matrix, and $\hat{\beta}$ is the estimated contact matrix
143 such that the objective function described above is minimized. In the interest of computational
144 efficiency, we also applied one-dimensional lasso on the Hi-C contact matrices in order to
145 estimate the matrices for high values of λ ($\lambda \gg 1$) and obtain the full hierarchy of TAD
146 boundaries. Using one-dimensional lasso instead of the two-dimensional version had no
147 negative impact on the correlations of Hi-C contact matrices between replicates
148 **(Supplementary Figure 2)**.

149

150 **Calculation of same-enzyme and cross-enzyme correlations**

151 We calculated two types of correlation for Hi-C matrices, to evaluate the performance of our
152 method: (a) same-enzyme correlation which corresponds to all the Hi-C replicates prepared
153 with the same restriction enzyme, (b) cross-enzyme correlation which corresponds to all the
154 sample pairs where the same Hi-C sample was prepared with two different enzymes (e.g
155 HindIII/MboI). Pearson correlation coefficients were calculated either on the filtered, ICE-
156 corrected [24] or scaled Hi-C contact matrices (Pearson) or the distance-normalized ones
157 (Pearson/z-score).

158

159 **TAD boundary “ratio” insulation score**

160 Given a potential TAD boundary, we denote the “upstream” TAD to the left of the boundary as
161 L , and the “downstream” TAD to the right as R . The cross-TAD interactions between L and R
162 are denoted as X . The “ratio” insulation score is defined as follows:

$$163 \text{ ratio} = \text{intra}_{\max}/\text{inter}$$

164 where:

$$165 \text{ intra}_{\max} = \max(\text{mean}(L), \text{mean}(R))$$

$$166 \text{ inter} = \text{mean}(X)$$

167 For more details, see [18].

168

169 **Classification of boundaries based on fused two-dimensional lasso**

170 We applied two-dimensional fused lasso to categorize TAD boundaries based on their strength.
171 The rationale behind this categorization is that topological domains separated by more
172 “permissive” (i.e. weaker) boundaries [25] will tend to fuse into larger domains when lasso is
173 applied, compared to TADs separated by well-defined, stronger boundaries. We indeed
174 applied this strategy and categorized boundaries into multiple groups ranging from the most
175 permissive to the strongest boundaries. The boundaries that were lost when λ value was
176 increased from 0 to 0.25, fall in the first category ($\lambda=0$), the ones lost when λ was increased to
177 0.5, in the second ($\lambda=0.2$) etc.

178

179 **Association of CTCF levels with boundary strength**

180 We obtained CTCF ChIP-sequencing data for the cell lines utilized in this study (with the
181 exception of KBM7 for which no publicly available dataset was available) and we uniformly re-
182 processed all data using HiC-bench [18]. Total CTCF levels (i.e. aggregated peak intensities
183 from potentially multiple CTCF peaks) at each TAD boundary were calculated and their
184 normalized distributions for each boundary category (weak to strong) were plotted in boxplots
185 in order to demonstrate the association of increased boundary strength with increased levels
186 of CTCF binding. We performed this analysis separately for TSS-only and non-TSS CTCF

187 binding sites. The rationale behind this separate analysis was based on the observation that
188 several TAD boundaries, especially strong boundaries, contain TSSs.

189

190 **Association of boundary strength with super-enhancers**

191 Super-enhancers were called using H3K27ac ChIP-seq data from GEO, ENCODE and in-
192 house generated data. Reads were first aligned with Bowtie2 v2.3.1 [26] and then HOMER
193 v4.6 [27] was used to call super-enhancers, all with standard parameters. For each super-
194 enhancer in each sample, we identified the corresponding TAD and its TAD boundaries. We
195 then calculated (per sample) the percentage of super-enhancers that are surrounded by
196 boundaries belonging in each boundary category, demonstrating that most super-enhancers
197 are insulated by strong boundaries.

198

199 **RESULTS**

200 **Analysis workflow**

201 The overall workflow, including our benchmark strategy and downstream analysis, is
202 summarized in **Figure 1a**. Our analysis starts with unprocessed Hi-C contact matrices
203 (“filtered” matrices). We then generate processed Hi-C matrices using both ICE “correction”
204 and our naïve “scaling” approach. Then, fused two-dimensional lasso is applied either on the
205 actual matrices or, alternatively, on the distance-normalized matrices. Matrix reproducibility
206 between biological replicates is assessed across samples for a variety of parameters, for
207 example, resolution, distance between interacting loci, sequencing depth, etc. Finally,
208 downstream analysis, involves the characterization of TAD boundaries based on their
209 insulating strength, the enrichment in CTCF binding, proximity to repeat elements and super-
210 enhancers, and, finally, their genetic alterations in cancer.

211

212 **Assessment of same-enzyme and cross-enzyme reproducibility of Hi-C contact** 213 **matrices**

214 Hi-C is prone to biases and multiple algorithms have been developed for Hi-C bias correction,
215 including probabilistic modelling methods [22], Poisson or negative binomial normalization [28]
216 and the widely popular Iterative Correction and Eigenvalue decomposition method (ICE) [24],
217 which assumes “equal visibility” of genomic loci. A similar iterative method named Sequential
218 Component Normalization was introduced by Cournac *et al.* [29]. Additional efficient correction
219 methods have been developed to handle high-resolution Hi-C datasets [30]. However,
220 estimating highly reproducible Hi-C contact maps remains a challenging task [31], especially
221 at high resolutions, as we also demonstrate below. Specifically, we focused on multiple factors
222 that may play an important role on reproducibility: first, we separately considered biological
223 replicates of Hi-C libraries generated with the same or different restriction enzymes; second,
224 we studied the impact of Hi-C matrix resolution (i.e. bin size); third, we assessed reproducibility
225 as a function of the distance of interacting loci pairs. Pearson correlation coefficients were
226 calculated for each pair of replicates (same- or cross-enzyme) on Hi-C contact matrices
227 estimated by three methods: (i) naïve filtering (i.e. matrix generation by simply using double-
228 sided accepted intra-chromosomal read pairs from **Supplementary Figure 1a**), (ii) iterative
229 correction (ICE) which has already been demonstrated to improve cross-enzyme correlation,
230 and (iii) our own “naïve” scaling method that only corrects for effective length bias (see
231 Methods for details). Importantly, correlations were computed both on the actual matrices, but
232 also on the distance-normalized matrices (see Methods for details), as Hi-C interactions are
233 typically concentrated around the diagonal of the Hi-C contact matrix, and values are dropping
234 exponentially as the distance between the interacting pairs is increasing (**Supplementary**
235 **Figure 1c**). Distance-normalized matrices account for the expected Hi-C read count as a
236 function of distance and may therefore reveal real distal interactions. The results of our
237 benchmark analysis are summarized in **Figure 1b**: the left panel summarizes the correlations
238 between replicates generated by the same restriction enzyme, whereas the right panel the
239 correlations between replicates generated by a different restriction enzymes. In both scenarios,
240 as expected, correlations drop quickly as finer resolutions (from 100kb to 20kb) are considered,
241 especially in the distance-normalized matrices. The same conclusion applies for increasing

242 distance (from 2Mb to 10Mb) between interacting loci, demonstrating that long-range
243 interactions require ultra-deep sequencing (beyond what is currently available in most of the
244 datasets in this study) in order to be detected reliably. To elaborate on this point, we repeated
245 the analysis after retaining only those samples with two replicates of at least 70 million or 110
246 million usable intra-chromosomal reads and resampling them down to 80 million or 120 million
247 per replicate (**Supplementary Figure 3** and **Supplementary Figure 4** respectively). Both
248 conclusions hold true with the new sequencing depth and are independent of the Hi-C contact
249 matrix estimation method. Finally, bias-correction methods (ICE and our scaling approach)
250 indeed improved cross-enzyme correlation over the naïve filtering method (**Figure 1b**).
251 Interestingly, this improvement came at the expense of lower correlations in the same-enzyme
252 case. More specifically, we observed that the larger the gain in cross-enzyme correlations,
253 the greater the loss in same-enzyme correlations (ICE method) (**Figure 1b**).

254

255 **Fused lasso improves same-enzyme and cross-enzyme correlations of Hi-C contact** 256 **matrices**

257 Motivated by the poor performance of all methods at fine resolutions and by the observation
258 of a surprising trade-off between cross-enzyme and same-enzyme correlations when
259 correcting for enzyme-related biases, we applied fused two-dimensional lasso [23], to obtain
260 improved estimates of Hi-C contact matrices. Briefly, two-dimensional fused lasso introduces
261 a parameter λ which penalizes differences between neighboring values in the Hi-C contact
262 matrix (see Methods for details). The effect of parameter λ is demonstrated in **Figure 2a** where
263 we show an example of the application of fused two-dimensional lasso on a Hi-C contact
264 matrix focused on an 8Mb locus on chromosome 8 for different values of parameter λ . To
265 evaluate the performance of fused lasso, we calculated same-enzyme and cross-enzyme
266 Pearson correlations between Hi-C contact matrices generated from different replicates.
267 Pearson correlation coefficients were calculated either for iteratively-corrected (ICE) or scaled
268 Hi-C contact matrices (at different λ values) and compared to the naïve filtering approach. The
269 results are summarized in **Figure 2b**. Increasing λ improves correlation independent of

270 resolution, restriction enzyme and bias-correction method, demonstrating the robustness of
271 our approach. Similarly, fused two-dimensional lasso improves the reproducibility of distance-
272 normalized matrices as demonstrated in **Figure 2c**. In all cases, as the value of λ increases,
273 the relative improvements in correlation are diminished. This observation can guide the
274 selection of λ , however a minimum of two replicates per sample are necessary to compute the
275 correlation and implement this strategy. Instead, we propose the use of degrees of freedom
276 as described in [23]. As demonstrated in **Supplementary Figure 2b**, the degrees of freedom
277 are rapidly decreasing for small values of λ , and quickly reaching a plateau with a moderate
278 increase in λ .

279

280 **Fused lasso reveals a TAD hierarchy linked to TAD boundary strength**

281 After demonstrating that parameter λ improves reproducibility of Hi-C contact matrices
282 independent of the bias-correction method, we hypothesized that increased values of λ may
283 also define distinct classes of TADs with different properties. For this reason, we now allowed
284 λ to range from 0 to 5 (after a finite value of λ , $\lambda \gg 5$, the entire Hi-C matrix attains a constant
285 value independent of the value of λ). For efficient computation, we used a one-dimensional
286 approximation of the two-dimensional lasso solution (see Methods for details and
287 **Supplementary Figure 2**). We then identified TADs at multiple λ values using HiC-bench,
288 and we observed that the number of TADs is monotonically decreasing with the value of λ
289 (**Figure 3a**), suggesting that by increasing λ , we are effectively identifying larger TADs
290 encompassing smaller TADs detected at lower λ values. Equivalently, certain TAD boundaries
291 “disappear” as λ is increased. Therefore, we hypothesized that TAD boundaries that disappear
292 at lower values of λ are weaker (i.e. lower insulation score), whereas boundaries that
293 disappear at higher values of λ are stronger (i.e. higher insulation score). To test this
294 hypothesis, we identified the TAD boundaries that are “lost” at each value of λ , and generated
295 the distributions of the insulation scores for each λ . As insulation score, we used the Hi-C
296 “ratio” score (see Methods), which was shown to outperform other TAD calling methods [18].
297 Indeed, as hypothesized, TAD boundaries lost at higher values of parameter λ are associated

298 with higher TAD insulation scores (**Figure 3b**). We then stratified TAD boundaries into five
299 classes (numbers 1 through 5 in **Figure 3b**; zero corresponds to lack of boundary) according
300 to their strength, independently in each Hi-C dataset used in this study. A heatmap
301 representation including all TAD boundaries and their associated class across all samples is
302 depicted in **Figure 3c**. Unbiased hierarchical clustering correctly grouped replicates and
303 related cell types independent of enzyme biases or batch effects related to the lab that
304 generated the Hi-C libraries, suggesting that TAD boundary strength can be used to
305 distinguish cell types. Equivalently, this finding suggests that, although TAD boundaries have
306 been shown to be largely invariant across cell types, a certain subset of TAD boundaries may
307 exhibit varying degrees of strength in different cell types. As expected, TAD boundary strength
308 was found to be positively associated with CTCF levels, suggesting that stronger CTCF
309 binding confers stronger insulation. Since we noticed that several TAD boundaries contain
310 TSSs, this analysis was done separately for all CTCF peaks (data not shown) and TSS-only
311 CTCF peaks (**Figure 3d**). Both approaches revealed the same trend, with the exception of the
312 class of strongest boundaries, where CTCF levels in TSS regions were significantly higher
313 compared to non-TSS regions, suggesting that the strongest boundaries are formed by CTCF-
314 mediated loops at gene promoters. SINE elements have also been shown to be enriched at
315 TAD boundaries [5], and besides confirming this finding, we now demonstrate that Alu
316 elements (the most abundant type of SINE elements) are enriched at stronger TAD boundaries
317 (**Supplementary Figure 5**, top-left panel). A comprehensive analysis of all major repetitive
318 element subtypes can be found in **Supplementary Figure 5**.

319

320 **Super-enhancers are preferentially localized within TADs demarcated by at least one** 321 **strong boundary**

322 We then explored what type of functional elements are localized within TADs demarcated by
323 strong TAD boundaries. Specifically, we tested super-enhancers identified in matched
324 samples (see Methods for details). Super-enhancers are key regulatory elements thought to
325 be defining cell identity [9, 32], and are usually found near the center of TADs [33]. Our

326 analysis determined that they are significantly more frequently localized within TADs insulated
327 by at least one strong TAD boundary (**Figure 3e**). Further analysis revealed that, in many
328 cases, super-enhancers are insulated by strong boundaries both in the upstream and
329 downstream directions (~3 times more likely compared to a strong/weak boundary
330 combination). We then mined the tissue-based map of the human proteome [34], a collection
331 of ubiquitously expressed genes as well as genes of variable tissue-specificity. Remarkably,
332 our analysis revealed that the genes closest to strong boundaries are significantly enriched in
333 the class of ubiquitously expressed genes (**Supplementary Figure 6a**). However, and
334 consistent with previous findings, tissue-specific genes are more enriched further away from
335 the TAD boundaries, in the vicinity of super-enhancers. Taken together, our findings suggest
336 that, because of their significance in gene regulation, super-enhancers should only target
337 genes confined in the “correct” TAD or neighborhood, while remaining strongly insulated from
338 genes in adjacent TADs. This is conceivably achieved by the strong TAD boundaries we have
339 identified in this study. At the same time, ubiquitously expressed genes are insulated from
340 enhancer elements in adjacent TADs by the same strong TAD boundaries in order to maintain
341 proper expression levels, unaffected by regulation from enhancer elements in adjacent TADs.

342

343 **Strong TAD boundaries are co-duplicated with super-enhancers and oncogenes in** 344 **cancer**

345 To further investigate the importance of variable boundary strength, we asked whether TAD
346 boundaries are prone to genetic alterations in cancer. To this end, we mined structural variants
347 released by the International Cancer Genome Consortium (ICGC) [35]. A summary of the
348 reported variant types across all cancer types available on ICGC, is presented in
349 **Supplementary Figure 6b**. First, for each focal (up to 1Mb) deletion event, we identified the
350 TAD boundaries closest to the breakpoints, and calculated the frequency of deletions by
351 boundary strength. We observed that the frequency of deletions monotonically decreased with
352 increasing boundary strength (**Figure 4a**). This suggests that strong TAD boundaries are less
353 frequently lost in cancer, as they may “safeguard” functional elements that are necessary for

354 proliferation. By contrast, the frequency of tandem duplications (up to 1Mb) increased with
355 increasing boundary strength (**Figure 4b**). Both results were robust to various cutoffs on the
356 sizes of the structural variants, within the usual range of TAD sizes (from 250kb to 2.5Mb).
357 Then, to further clarify the connection between super-enhancers, strong TAD boundaries and
358 cancer, we studied tandem duplication events where super-enhancers (obtained from the
359 largest available collection of super-enhancers [36]) are co-duplicated with adjacent strong
360 boundaries. As demonstrated in **Figure 4c**, super-enhancers are indeed co-duplicated with
361 strong TAD boundaries. This suggests that, in cancer, not only are strong boundaries
362 protected from deletions, but they are also co-duplicated with super-enhancer elements.
363 Intriguingly, this observation raises the possibility that pairs of super-enhancers and
364 oncogenes represent functional entities encapsulated by strong boundaries, that are
365 frequently duplicated or perhaps amplified in malignancies. Such an example of an oncogene
366 is shown in **Figure 4d**: *MYC*, a well-known oncogene that is typically overexpressed in cancer,
367 is localized next to a strong TAD boundary and is co-duplicated with the boundary as well as
368 with several proximal super-enhancers. Taken together, these observations highlight the
369 importance of strong TAD boundaries in the context of cancer.

370

371 **DISCUSSION**

372 Multiple recent studies have revealed that the metazoan genome is compartmentalized in
373 boundary-demarcated functional units known as topologically associating domains (TADs).
374 TADs are highly conserved across species and cell types. A few studies, however, provide
375 compelling evidence that specific TADs, despite the fact that they are largely invariant, exhibit
376 some plasticity. Given that TAD boundary disruption has been recently linked to aberrant gene
377 activation and multiple disorders including developmental defects and cancer, categorization
378 of boundaries based on their strength and identification of their unique features becomes of
379 particular importance. In this study, we developed a method based on fused two-dimensional
380 lasso in order to categorize TAD boundaries based on their strength. We demonstrated that
381 our method: (a) improves the correlation of Hi-C contact matrices irrespective of the Hi-C bias

382 correction method used, (b) reveals multiple levels of chromatin organization and (c)
383 successfully identifies boundaries of variable strength and that strong predicted boundaries
384 exhibit certain expected features, such as elevated CTCF levels and increased insulating
385 capacity. By performing an integrative analysis of estimated boundary strength with super-
386 enhancers in matched samples, we observed that super-enhancers are preferentially
387 insulated by strong boundaries, suggesting that super-enhancers and strong boundaries may
388 represent a biologically relevant entity. Motivated by this observation, we examined the
389 frequency of structural alterations involving strong boundaries and super-enhancers. We
390 found that not only strong boundaries are “protected” from deletions, but, more importantly,
391 they are co-duplicated together with super-enhancers. Recently, it has been shown that
392 genetic or epigenetic alterations near enhancers may lead to aberrant activation of oncogenes
393 [37–40]. Our results, expand on these studies and highlight a synergistic role of super-
394 enhancers and TAD boundaries in cancer.

395

396 **AUTHOR CONTRIBUTIONS**

397 YG and CL performed computational analyses and generated figures. AT, AL and PK
398 conceived this study. PN performed the CUTLL-1 Hi-C experiments. PN and IA offered
399 biological insights and helped with the interpretation of Hi-C data. AT designed and
400 implemented the method. CL and AT wrote the manuscript. All authors read and approved the
401 final manuscript.

402

403

404 **ACKNOWLEDGEMENTS**

405 We would like to thank all members of the Tsirigos and Aifantis Laboratories for critical
406 evaluation of the manuscript. We would like to thank the Applied Bioinformatics Laboratories
407 (ABL) at the NYU School of Medicine for providing bioinformatics support and helping with the
408 analysis and interpretation of the data. This work has used computing resources at the NYU

409 High Performance Computing Facility (HPCF). We also thank the Genome Technology Center
410 (GTC) for expert library preparation and sequencing. This shared resource is partially
411 supported by the Cancer Center Support Grant, P30CA016087, at the Laura and Isaac
412 Perlmutter Cancer Center.

413

414 **FUNDING**

415 The study was supported by the American Cancer Society [RSG-15-189-01-RMC to AT] and
416 a Leukemia & Lymphoma Society New Idea Award [8007-17 to AT]. NYU Genome Technology
417 Center (GTC) is a shared resource, partially supported by the Cancer Center Support Grant
418 [P30CA016087] at the Laura and Isaac Perlmutter Cancer Center.

419

420 **TABLE AND FIGURES LEGENDS**

421 **Figure 1. (a)** Overall workflow and benchmark strategy, **(b)** Comparison of Hi-C contact
422 matrices between biological replicates generated from Hi-C library using the same or different
423 restriction enzyme; Hi-C matrices were estimated using three methods (naïve filtering, iterative
424 correction and simple scaling); assessment was performed using Pearson correlation on the
425 actual or distance-normalized Hi-C matrices at resolutions ranging from 100kb to 20kb and
426 maximum distances of 2Mb, 6Mb and 10Mb between interacting pairs

427

428 **Figure 2. Fused two-dimensional lasso improves reproducibility of Hi-C contact**
429 **matrices. (a)** Example of application of fused two-dimensional lasso on a Hi-C contact matrix
430 focused on a 8Mb locus on chromosome 8 for different values of parameter λ (top
431 panel=original matrix; bottom panel=distance-normalized matrix), **(b)** Hi-C contact matrix
432 correlations are improved by increasing the value of fused lasso parameter λ both for matrices
433 estimated by ICE as well as by our simple scaling method. **(c)** Hi-C contact matrix correlations
434 of distance-normalized matrices. Correlations of Hi-C contact matrices generated by the naïve
435 filtering method are marked by the red line in each panel.

436

437 **Figure 3. Classification and characterization of TAD boundaries according to insulation**
438 **score. (a)** Number of TADs for λ values ranging from 0 to 5, **(b)** TAD boundaries lost at higher
439 values of parameter λ are associated with higher TAD insulation scores, **(c)** Heatmap
440 representation of TAD boundary insulation strength across samples; hierarchical clustering
441 correctly groups replicates and related cell types independent of enzyme biases or batch
442 effects related to the lab that generated the Hi-C libraries, **(d)** TAD boundary strength is
443 associated with CTCF levels, **(e)** Fraction of super-enhancer elements in the vicinity of
444 boundaries of variable strength. The gradient of blue corresponds to λ values with darker blue
445 denoting higher λ value.

446

447 **Figure 4. Pan-cancer analysis of strong vs weak TAD boundaries. (a)** Schematic of pan-
448 cancer analysis (left panel) and classification of focally deleted boundaries in cancer according
449 to their strength (right panel), **(b)** Schematic of pan-cancer analysis (left panel) and
450 classification of focally duplicated boundaries in cancer according to their strength (right panel),
451 **(c)** Schematic of pan-cancer analysis (left panel) and co-duplications of TAD boundaries with
452 super-enhancers in cancer (right panel), **(d)** Snapshot of the *MYC* locus: a strong boundary
453 (black bar) is frequently co-duplicated with *MYC* and potential super-enhancers in cancer
454 patients (highlighted area). IGV tracks from top to bottom: boundary score (gray), strong
455 boundaries (black bars), super-enhancer track from SEA (blue bars), RefSeq genes,
456 duplication frequency (red graph) and ICGC patient tandem duplications (red bars).

457

458 **Supplementary Figure 1. Quality assessment of Hi-C datasets. (a)** Counts of Hi-C read
459 pairs in various read categories: dark and light green indicate read pairs that were not
460 designated as artifacts and can be used in downstream analyses, **(b)** Percentages of Hi-C
461 reads in each category, **(c)** Average scaled Hi-C read pair count as a function of distance
462 between interacting loci.

463

464 **Supplementary Figure 2. Fused one-dimensional lasso improves reproducibility of**
465 **distance-normalized Hi-C contact matrices. (a)** distance-normalized Hi-C contact matrix
466 correlations are improved by increasing the value of fused lasso parameter λ both for matrices
467 estimated by ICE as well as by our simple scaling method; correlations of distance-normalized
468 Hi-C contact matrices generated by the naïve filtering method are marked by the red line in
469 each panel. The gradient of blue corresponds to λ values with darker blue denoting higher λ
470 value, **(b)** degrees of freedom as a function of λ .

471

472 **Supplementary Figure 3. Comparison of Hi-C contact matrices between biological**
473 **replicates generated from Hi-C library using the same restriction enzyme.** Three
474 methods (naïve filtering, iterative correction and simple scaling) were used for estimation.
475 Assessment was performed using Pearson correlation on the actual or distance-normalized
476 Hi-C matrices at resolutions ranging from 100kb to 20kb and maximum distances of 2Mb, 6Mb
477 and 10Mb between interacting pairs. Only samples with approximately 80 million usable intra-
478 chromosomal reads were considered.

479

480 **Supplementary Figure 4. Comparison of Hi-C contact matrices between biological**
481 **replicates generated from Hi-C library using the same restriction enzyme.** Three
482 methods (naïve filtering, iterative correction and simple scaling) were used for estimation.
483 Assessment was performed using Pearson correlation on the actual or distance-normalized
484 Hi-C matrices at resolutions ranging from 100kb to 20kb and maximum distances of 2Mb, 6Mb
485 and 10Mb between interacting pairs. Only samples with approximately 120 million usable intra-
486 chromosomal reads were considered.

487

488 **Supplementary Figure 5.** Normalized numbers of repeat elements in proximity to boundaries
489 of certain boundary strength. Darker blue in the blue colour gradient denotes higher boundary
490 strength.

491

492 **Supplementary Figure 6. (a)** Fraction of ubiquitous genes as well as genes of increasing
493 tissue-specificity in the vicinity of boundaries of variable strength, **(b)** Distribution of somatic
494 structural alterations in the ICGC database.

495

496

497 REFERENCES

498 1. Dekker J, Marti-Renom MA, Mirny LA: **Exploring the three-dimensional organization of**
499 **genomes: interpreting chromatin interaction data.** *Nat Rev Genet* 2013, **14**:390–403.

500 2. Schmitt AD, Hu M, Ren B: **Genome-wide mapping and analysis of chromosome**
501 **architecture.** *Nat Rev Mol Cell Biol* 2016, **17**:743–755.

502 3. Lieberman-Aiden E, van Berkum NL, Williams L, Imakaev M, Ragooczy T, Telling A, Amit I,
503 Lajoie BR, Sabo PJ, Dorschner MO, Sandstrom R, Bernstein B, Bender MA, Groudine M,
504 Gnirke A, Stamatoyannopoulos J, Mirny LA, Lander ES, Dekker J: **Comprehensive**
505 **mapping of long-range interactions reveals folding principles of the human genome.**
506 *Science* 2009, **326**:289–293.

507 4. Fraser J, Ferrai C, Chiariello AM, Schueler M, Rito T, Laudanno G, Barbieri M, Moore BL,
508 Kraemer DCA, Aitken S, Xie SQ, Morris KJ, Itoh M, Kawaji H, Jaeger I, Hayashizaki Y,
509 Carninci P, Forrest ARR, FANTOM Consortium, Semple CA, Dostie J, Pombo A, Nicodemi
510 M: **Hierarchical folding and reorganization of chromosomes are linked to**
511 **transcriptional changes in cellular differentiation.** *Mol Syst Biol* 2015, **11**:852.

512 5. Dixon JR, Selvaraj S, Yue F, Kim A, Li Y, Shen Y, Hu M, Liu JS, Ren B: **Topological**
513 **domains in mammalian genomes identified by analysis of chromatin interactions.**
514 *Nature* 2012, **485**:376–380.

515 6. Nora EP, Lajoie BR, Schulz EG, Giorgetti L, Okamoto I, Servant N, Piolot T, van Berkum
516 NL, Meisig J, Sedat J, Gribnau J, Barillot E, Blüthgen N, Dekker J, Heard E: **Spatial**
517 **partitioning of the regulatory landscape of the X-inactivation centre.** *Nature* 2012,
518 **485**:381–385.

519 7. Sexton T, Yaffe E, Kenigsberg E, Bantignies F, Leblanc B, Hoichman M, Parrinello H,
520 Tanay A, Cavalli G: **Three-dimensional folding and functional organization principles of**
521 **the Drosophila genome.** *Cell* 2012, **148**:458–472.

522 8. Phillips-Cremins JE, Sauria MEG, Sanyal A, Gerasimova TI, Lajoie BR, Bell JSK, Ong C-
523 T, Hookway TA, Guo C, Sun Y, Bland MJ, Wagstaff W, Dalton S, McDevitt TC, Sen R,
524 Dekker J, Taylor J, Corces VG: **Architectural protein subclasses shape 3D organization**
525 **of genomes during lineage commitment.** *Cell* 2013, **153**:1281–1295.

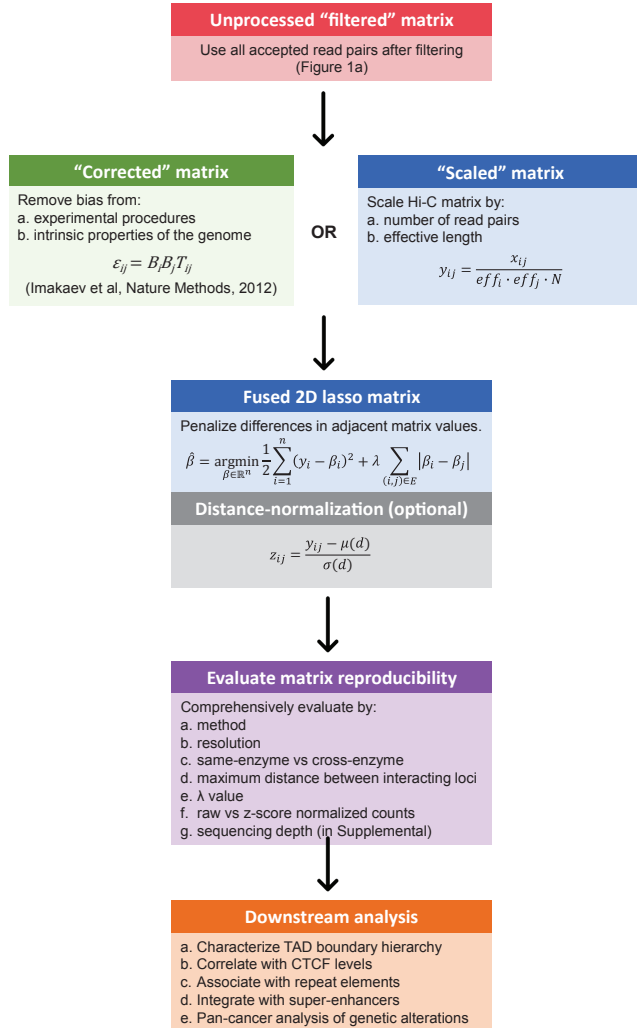
- 526 9. Downen JM, Fan ZP, Hnisz D, Ren G, Abraham BJ, Zhang LN, Weintraub AS, Schuijers J,
527 Lee TI, Zhao K, Young RA: **Control of cell identity genes occurs in insulated**
528 **neighborhoods in mammalian chromosomes.** *Cell* 2014, **159**:374–387.
- 529 10. Filippova D, Patro R, Duggal G, Kingsford C: **Identification of alternative topological**
530 **domains in chromatin.** *Algorithms Mol Biol* 2014, **9**:14.
- 531 11. Weinreb C, Raphael BJ: **Identification of hierarchical chromatin domains.**
532 *Bioinformatics* 2016, **32**:1601–1609.
- 533 12. Crane E, Bian Q, McCord RP, Lajoie BR, Wheeler BS, Ralston EJ, Uzawa S, Dekker J,
534 Meyer BJ: **Condensin-driven remodelling of X chromosome topology during dosage**
535 **compensation.** *Nature* 2015, **523**:240–244.
- 536 13. Haddad N, Vaillant C, Jost D: **IC-Finder: inferring robustly the hierarchical**
537 **organization of chromatin folding.** *Nucleic Acids Res* 2017.
- 538 14. Cubeñas-Potts C, Corces VG: **Topologically Associating Domains: An invariant**
539 **framework or a dynamic scaffold?** *Nucleus* 2015, **6**:430–434.
- 540 15. Li L, Lyu X, Hou C, Takenaka N, Nguyen HQ, Ong C-T, Cubeñas-Potts C, Hu M, Lei EP,
541 Bosco G, Qin ZS, Corces VG: **Widespread rearrangement of 3D chromatin organization**
542 **underlies polycomb-mediated stress-induced silencing.** *Mol Cell* 2015, **58**:216–231.
- 543 16. Narendra V, Bulajić M, Dekker J, Mazzoni EO, Reinberg D: **CTCF-mediated**
544 **topological boundaries during development foster appropriate gene regulation.** *Genes*
545 *Dev* 2016, **30**:2657–2662.
- 546 17. Friedman J, Hastie T, Höfling H, Tibshirani R: **Pathwise coordinate optimization.** *Ann*
547 *Appl Stat* 2007, **1**:302–332.
- 548 18. Lazaris C, Kelly S, Ntziachristos P, Aifantis I, Tsirigos A: **HiC-bench: comprehensive**
549 **and reproducible Hi-C data analysis designed for parameter exploration and**
550 **benchmarking.** *BMC Genomics* 2017, **18**:22.
- 551 19. Dixon JR, Jung I, Selvaraj S, Shen Y, Antosiewicz-Bourget JE, Lee AY, Ye Z, Kim A,
552 Rajagopal N, Xie W, Diao Y, Liang J, Zhao H, Lobanenkov VV, Ecker JR, Thomson JA, Ren
553 B: **Chromatin architecture reorganization during stem cell differentiation.** *Nature* 2015,
554 **518**:331–336.
- 555 20. Rao SSP, Huntley MH, Durand NC, Stamenova EK, Bochkov ID, Robinson JT, Sanborn
556 AL, Machol I, Omer AD, Lander ES, Aiden EL: **A 3D map of the human genome at**
557 **kilobase resolution reveals principles of chromatin looping.** *Cell* 2014, **159**:1665–1680.
- 558 21. Jin F, Li Y, Dixon JR, Selvaraj S, Ye Z, Lee AY, Yen C-A, Schmitt AD, Espinoza CA, Ren
559 B: **A high-resolution map of the three-dimensional chromatin interactome in human**
560 **cells.** *Nature* 2013, **503**:290–294.
- 561 22. Yaffe E, Tanay A: **Probabilistic modeling of Hi-C contact maps eliminates**
562 **systematic biases to characterize global chromosomal architecture.** *Nat Genet* 2011,
563 **43**:1059–1065.

- 564 23. Tibshirani RJ, Taylor J: **The solution path of the generalized lasso.** *The Annals of*
565 *Statistics* 2011, **39**:1335–1371.
- 566 24. Imakaev M, Fudenberg G, McCord RP, Naumova N, Goloborodko A, Lajoie BR, Dekker
567 J, Mirny LA: **Iterative correction of Hi-C data reveals hallmarks of chromosome**
568 **organization.** *Nat Methods* 2012, **9**:999–1003.
- 569 25. Rocha PP, Raviram R, Bonneau R, Skok JA: **Breaking TADs: insights into**
570 **hierarchical genome organization.** *Epigenomics* 2015, **7**:523–526.
- 571 26. Langmead B, Salzberg SL: **Fast gapped-read alignment with Bowtie 2.** *Nat Methods*
572 2012, **9**:357–359.
- 573 27. Heinz S, Benner C, Spann N, Bertolino E, Lin YC, Laslo P, Cheng JX, Murre C, Singh H,
574 Glass CK: **Simple combinations of lineage-determining transcription factors prime cis-**
575 **regulatory elements required for macrophage and B cell identities.** *Mol Cell* 2010,
576 **38**:576–589.
- 577 28. Hu M, Deng K, Selvaraj S, Qin Z, Ren B, Liu JS: **HiCNorm: removing biases in Hi-C**
578 **data via Poisson regression.** *Bioinformatics* 2012, **28**:3131–3133.
- 579 29. Cournac A, Marie-Nelly H, Marbouty M, Koszul R, Mozziconacci J: **Normalization of a**
580 **chromosomal contact map.** *BMC Genomics* 2012, **13**:436.
- 581 30. Knight PA, Ruiz D: **A fast algorithm for matrix balancing.** *IMA Journal of Numerical*
582 *Analysis* 2013, **33**:1029–1047.
- 583 31. Forcato M, Nicoletti C, Pal K, Livi CM, Ferrari F, Bicciato S: **Comparison of**
584 **computational methods for Hi-C data analysis.** *Nat Methods* 2017.
- 585 32. Whyte WA, Orlando DA, Hnisz D, Abraham BJ, Lin CY, Kagey MH, Rahl PB, Lee TI,
586 Young RA: **Master transcription factors and mediator establish super-enhancers at key**
587 **cell identity genes.** *Cell* 2013, **153**:307–319.
- 588 33. Schmitt AD, Hu M, Jung I, Xu Z, Qiu Y, Tan CL, Li Y, Lin S, Lin Y, Barr CL, Ren B: **A**
589 **compendium of chromatin contact maps reveals spatially active regions in the human**
590 **genome.** *Cell Rep* 2016, **17**:2042–2059.
- 591 34. Uhlén M, Fagerberg L, Hallström BM, Lindskog C, Oksvold P, Mardinoglu A, Sivertsson
592 Å, Kampf C, Sjöstedt E, Asplund A, Olsson I, Edlund K, Lundberg E, Navani S, Szigartyo
593 CA-K, Odeberg J, Djureinovic D, Takanen JO, Hober S, Alm T, Edqvist P-H, Berling H,
594 Tegel H, Mulder J, Rockberg J, Nilsson P, Schwenk JM, Hamsten M, von Feilitzen K,
595 Forsberg M, et al.: **Proteomics. Tissue-based map of the human proteome.** *Science*
596 2015, **347**:1260419.
- 597 35. Zhang J, Baran J, Cros A, Guberman JM, Haider S, Hsu J, Liang Y, Rivkin E, Wang J,
598 Whitty B, Wong-Erasmus M, Yao L, Kasprzyk A: **International Cancer Genome**
599 **Consortium Data Portal--a one-stop shop for cancer genomics data.** *Database (Oxford)*
600 2011, **2011**:bar026.
- 601 36. Wei Y, Zhang S, Shang S, Zhang B, Li S, Wang X, Wang F, Su J, Wu Q, Liu H, Zhang Y:
602 **SEA: a super-enhancer archive.** *Nucleic Acids Res* 2016, **44**:D172–9.

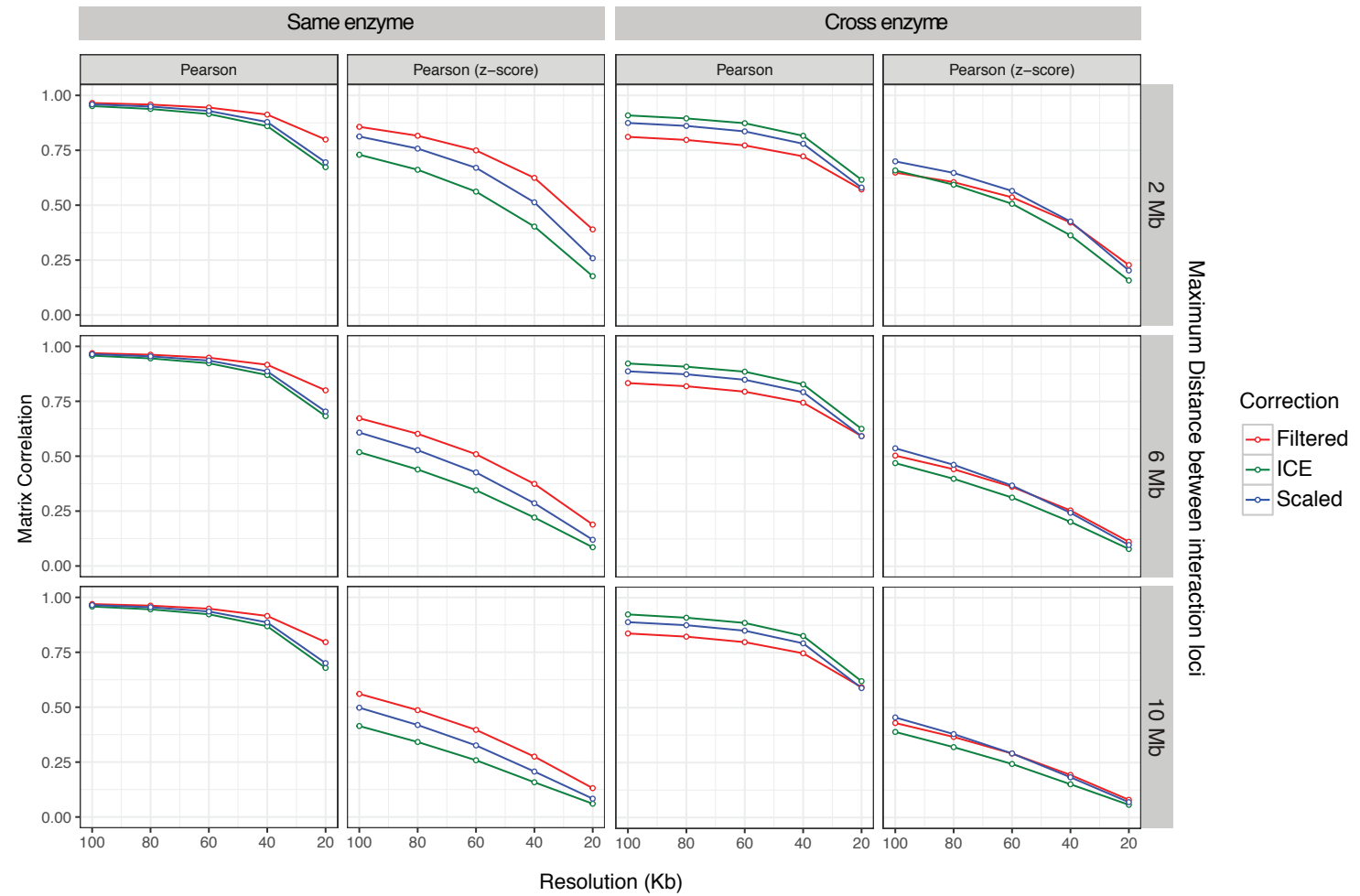
- 603 37. Lupiáñez DG, Kraft K, Heinrich V, Krawitz P, Brancati F, Klopocki E, Horn D, Kayserili H,
604 Opitz JM, Laxova R, Santos-Simarro F, Gilbert-Dussardier B, Wittler L, Borschiwer M, Haas
605 SA, Osterwalder M, Franke M, Timmermann B, Hecht J, Spielmann M, Visel A, Mundlos S:
606 **Disruptions of topological chromatin domains cause pathogenic rewiring of gene-**
607 **enhancer interactions.** *Cell* 2015, **161**:1012–1025.
- 608 38. Flavahan WA, Drier Y, Liao BB, Gillespie SM, Venteicher AS, Stemmer-Rachamimov
609 AO, Suvà ML, Bernstein BE: **Insulator dysfunction and oncogene activation in IDH**
610 **mutant gliomas.** *Nature* 2016, **529**:110–114.
- 611 39. Hnisz D, Weintraub AS, Day DS, Valton A-L, Bak RO, Li CH, Goldmann J, Lajoie BR,
612 Fan ZP, Sigova AA, Reddy J, Borges-Rivera D, Lee TI, Jaenisch R, Porteus MH, Dekker J,
613 Young RA: **Activation of proto-oncogenes by disruption of chromosome**
614 **neighborhoods.** *Science* 2016, **351**:1454–1458.
- 615 40. Weischenfeldt J, Dubash T, Drainas AP, Mardin BR, Chen Y, Stütz AM, Waszak SM,
616 Bosco G, Halvorsen AR, Raeder B, Efthymiopoulos T, Erkek S, Siegl C, Brenner H,
617 Brustugun OT, Dieter SM, Northcott PA, Petersen I, Pfister SM, Schneider M, Solberg SK,
618 Thunissen E, Weichert W, Zichner T, Thomas R, Peifer M, Helland A, Ball CR, Jechlinger M,
619 Sotillo R, et al.: **Pan-cancer analysis of somatic copy-number alterations implicates**
620 **IRS4 and IGF2 in enhancer hijacking.** *Nat Genet* 2017, **49**:65–74.
- 621

Figure 1

a



b



a bioRxiv preprint doi: <https://doi.org/10.1101/141481>; this version posted June 30, 2017. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

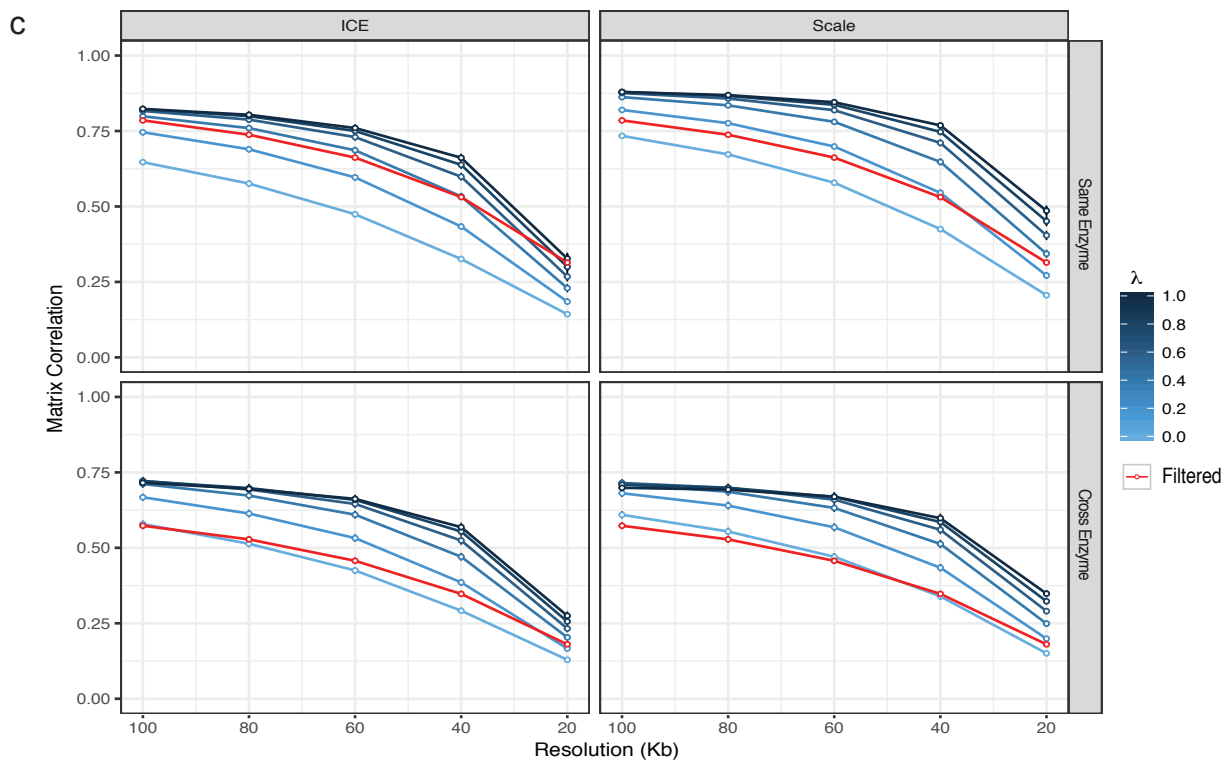
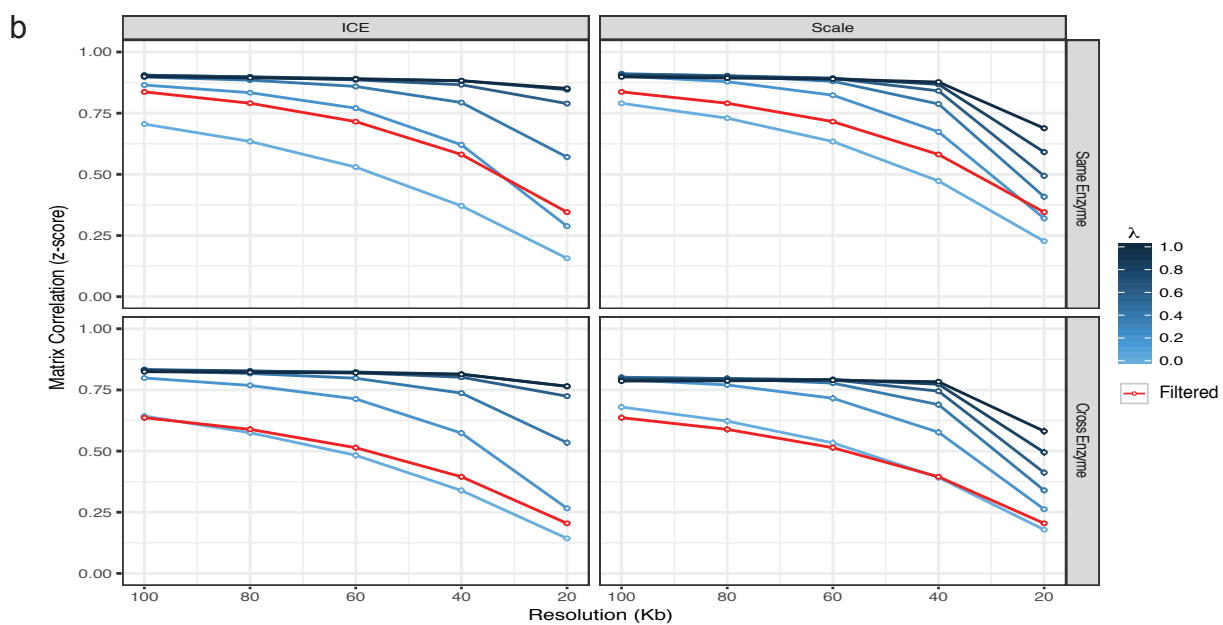
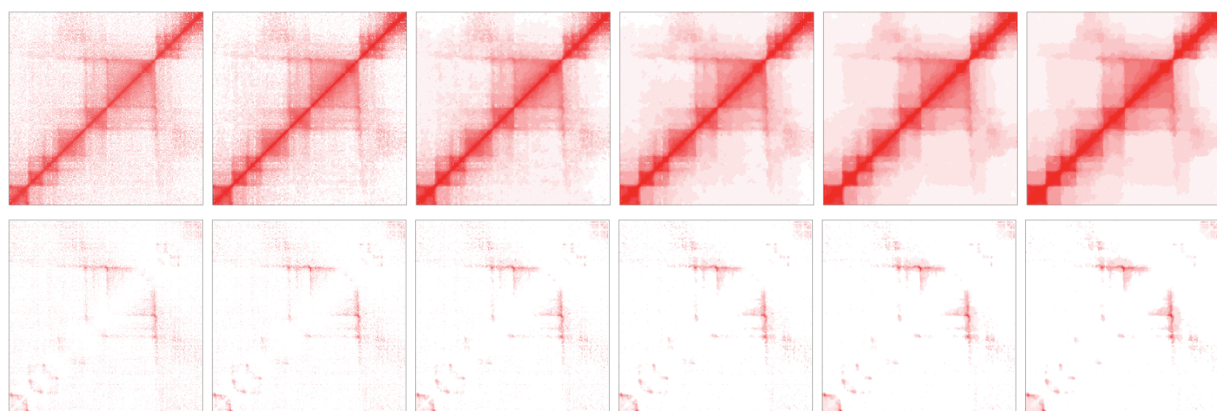
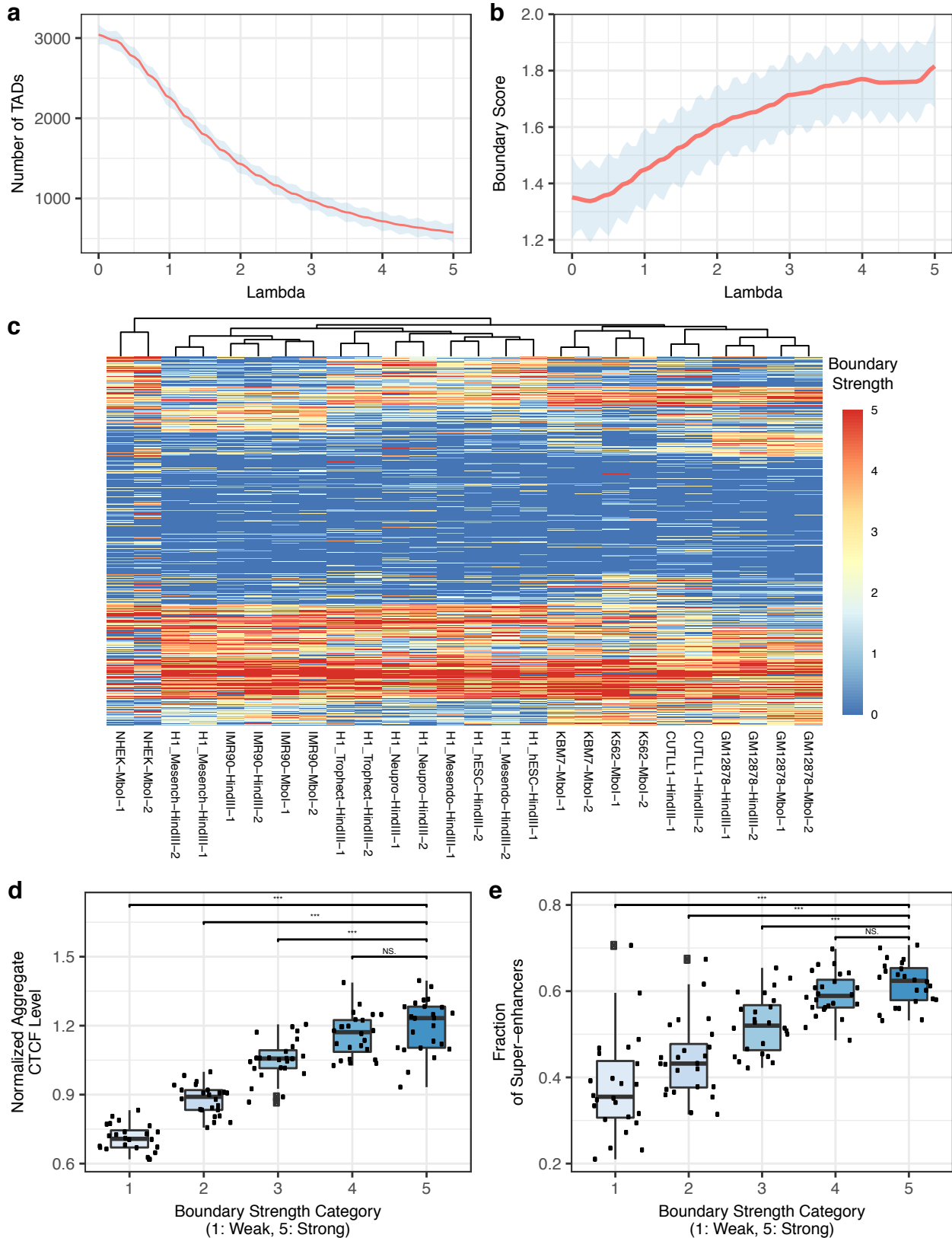
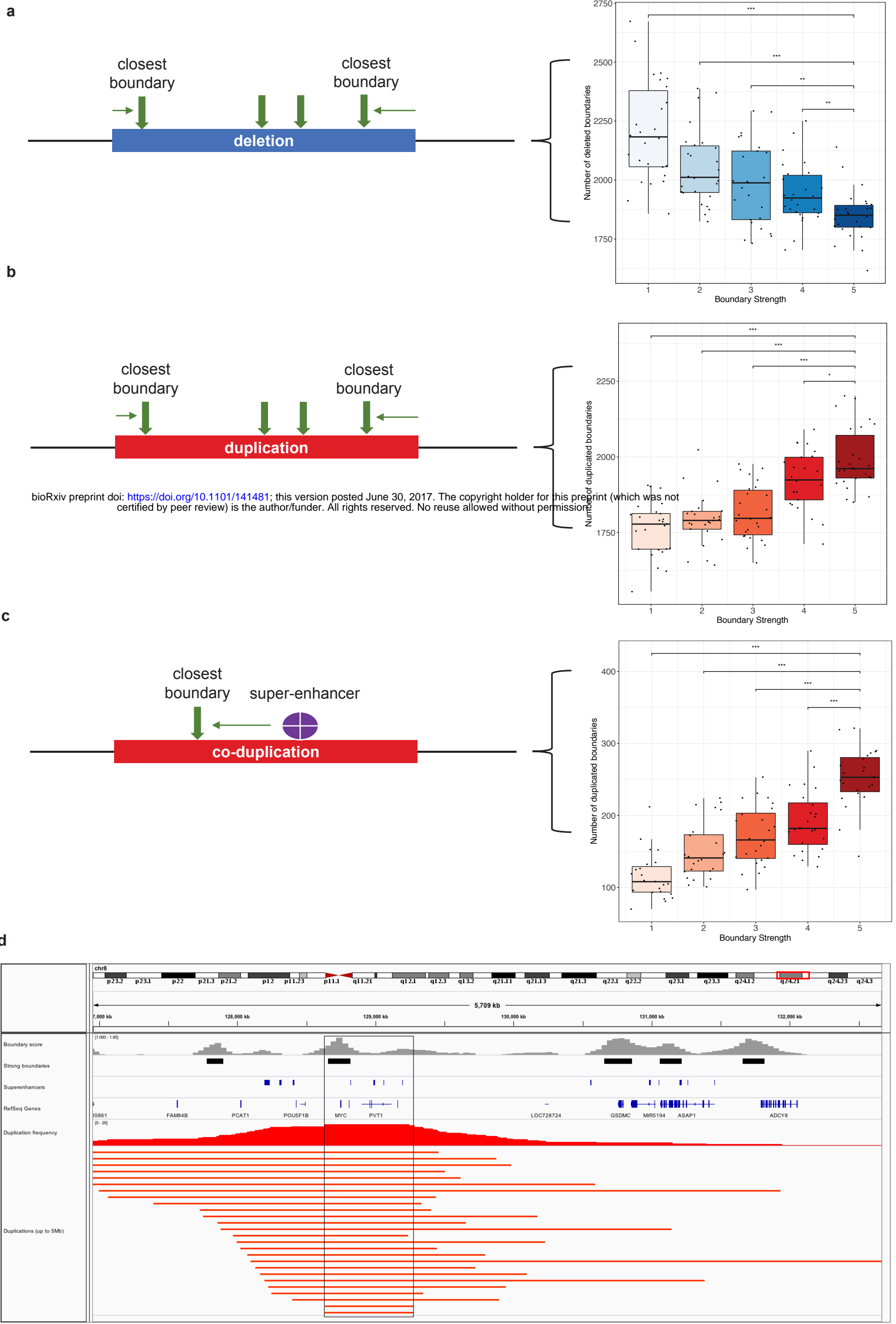
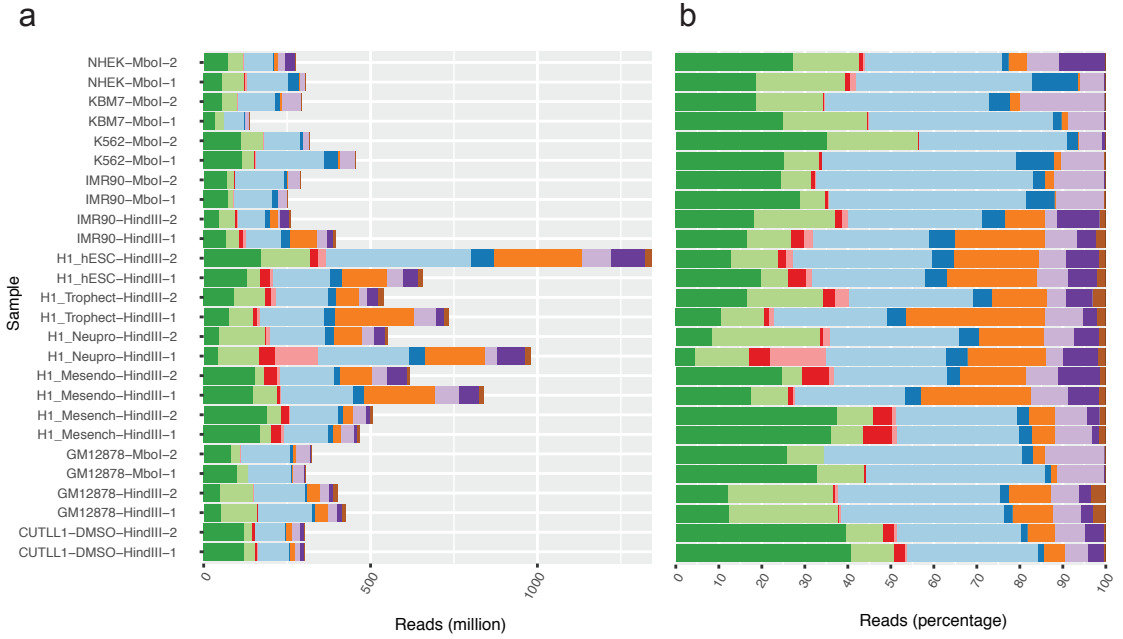
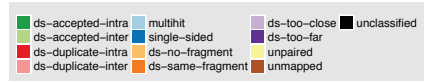


Figure 3

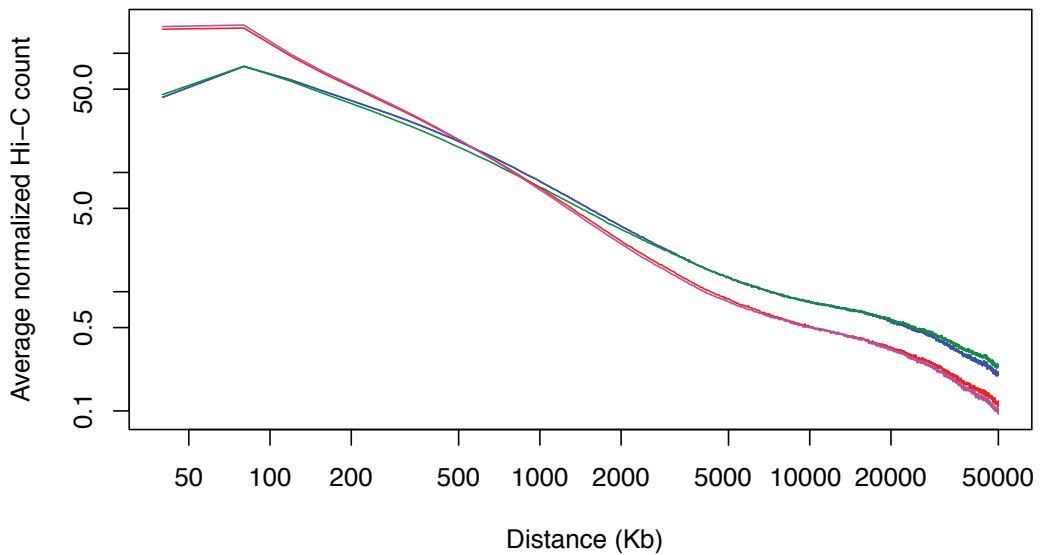


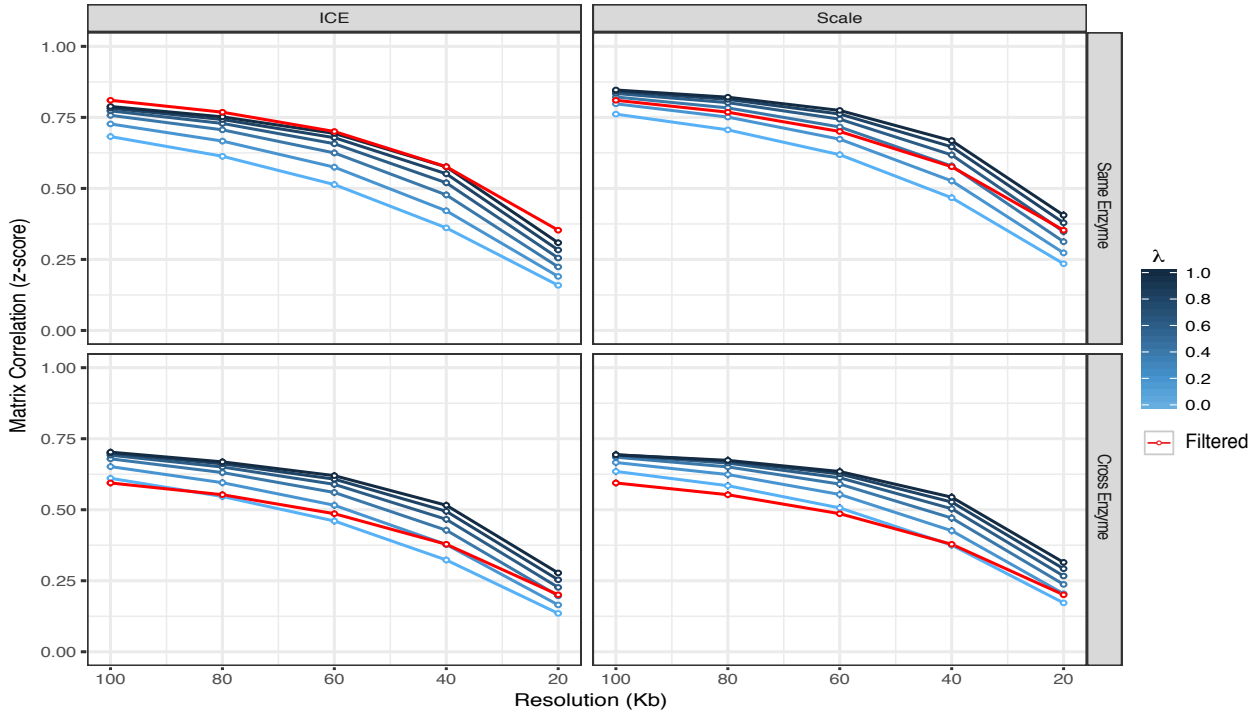
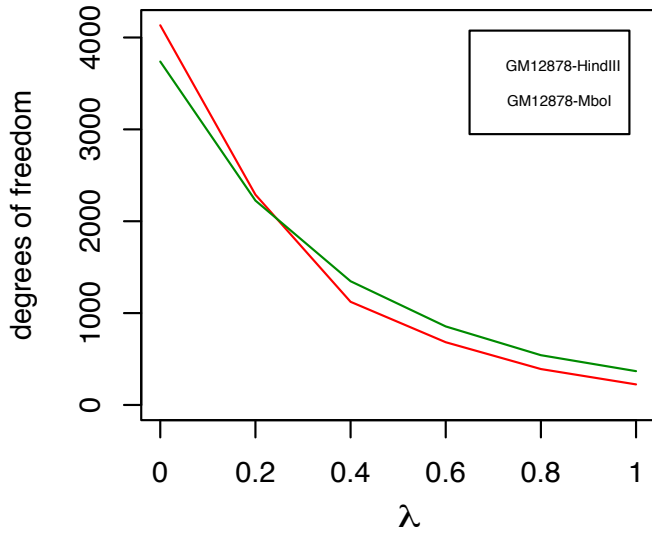
bioRxiv preprint doi: <https://doi.org/10.1101/141481>; this version posted June 30, 2017. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.



c

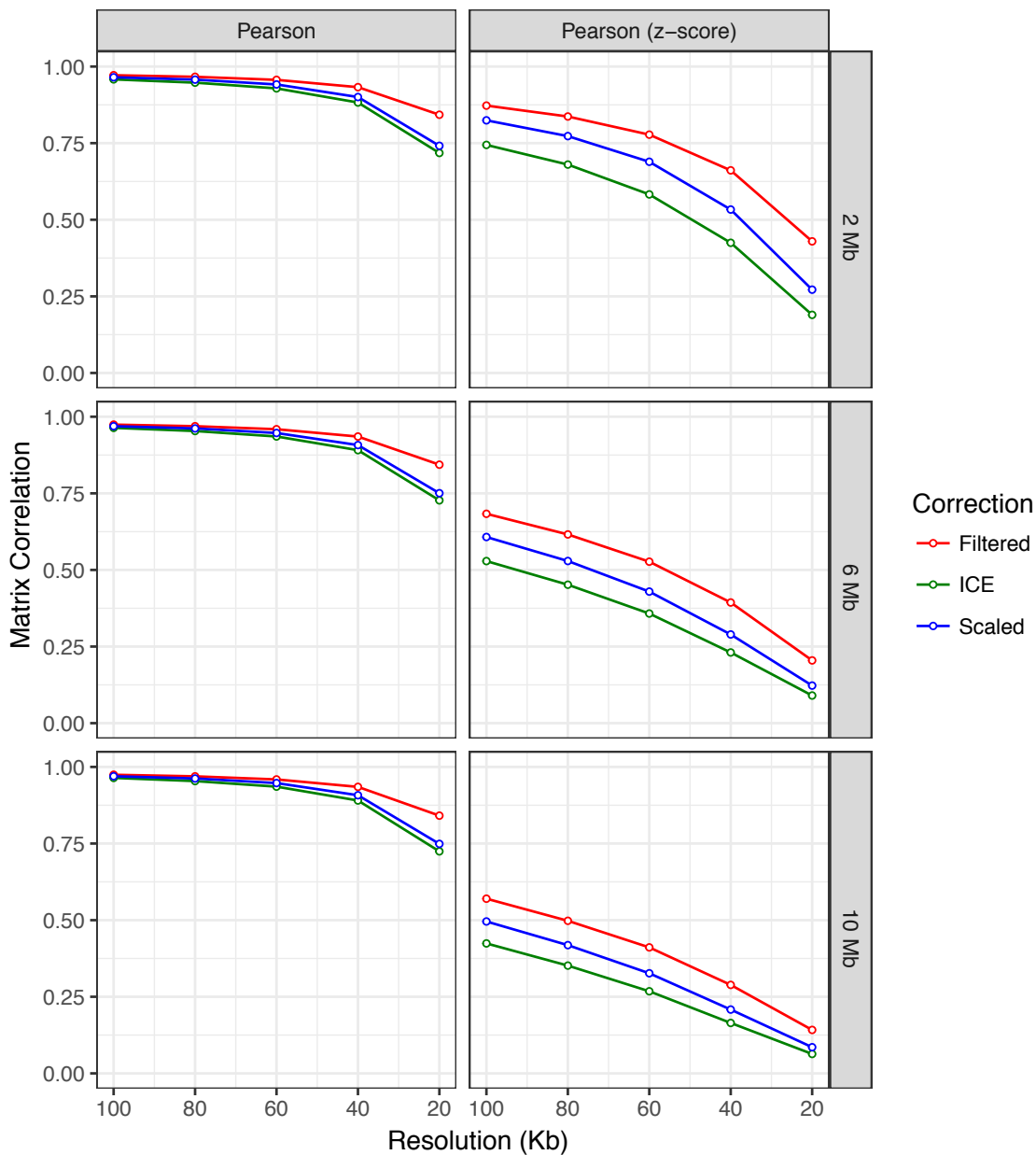
Hi-C count as a function of distance



a**b**

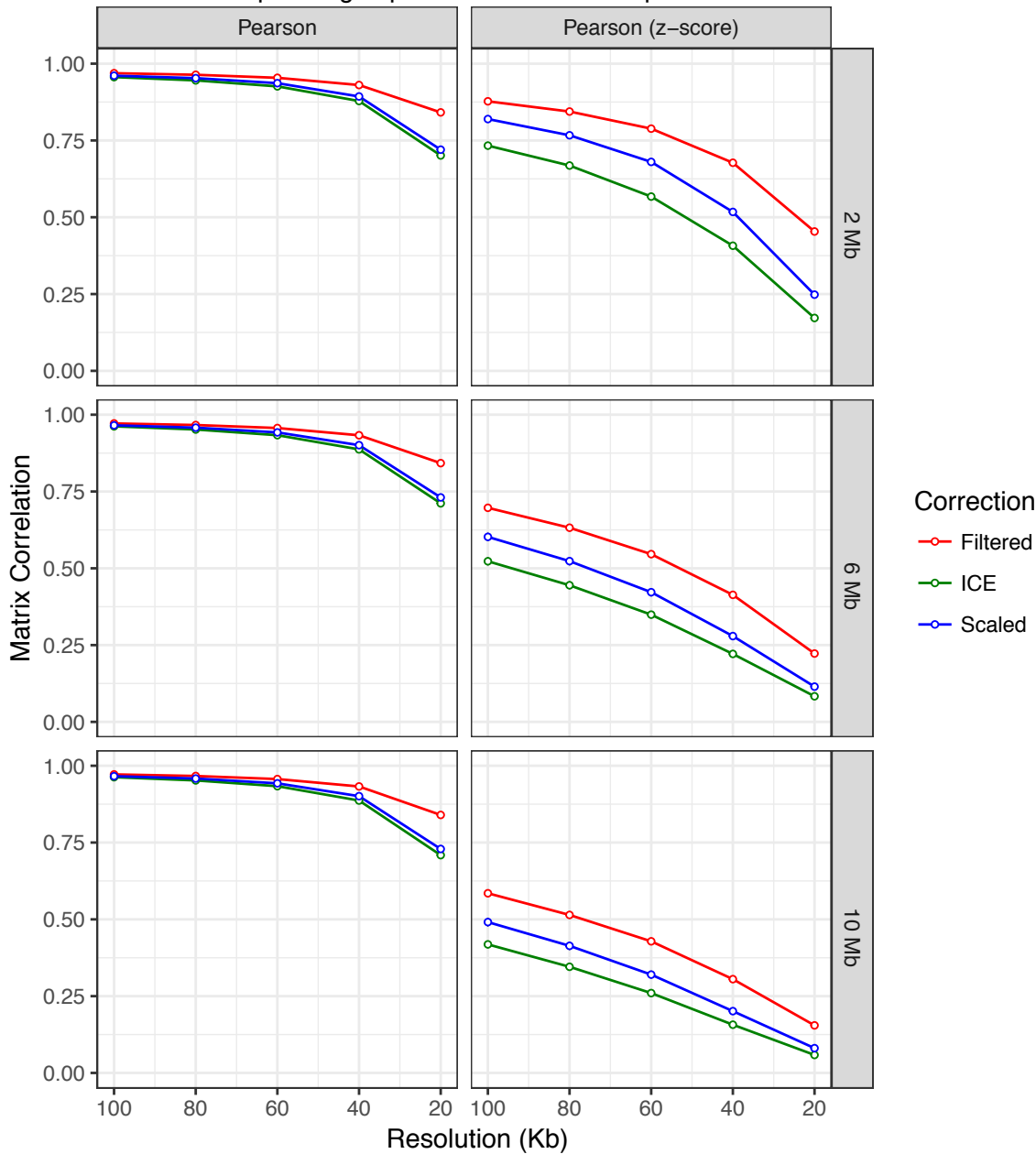
Supplementary Figure 3

Sequencing depth = 80 million read pairs

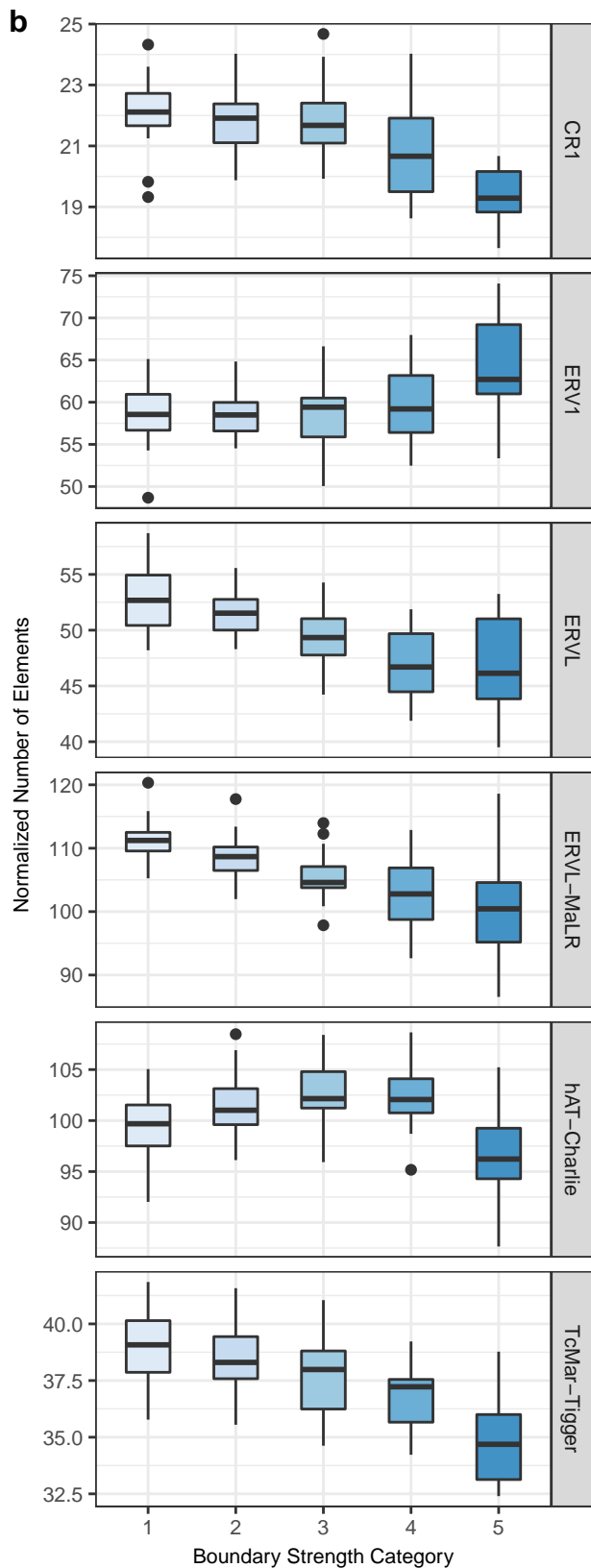
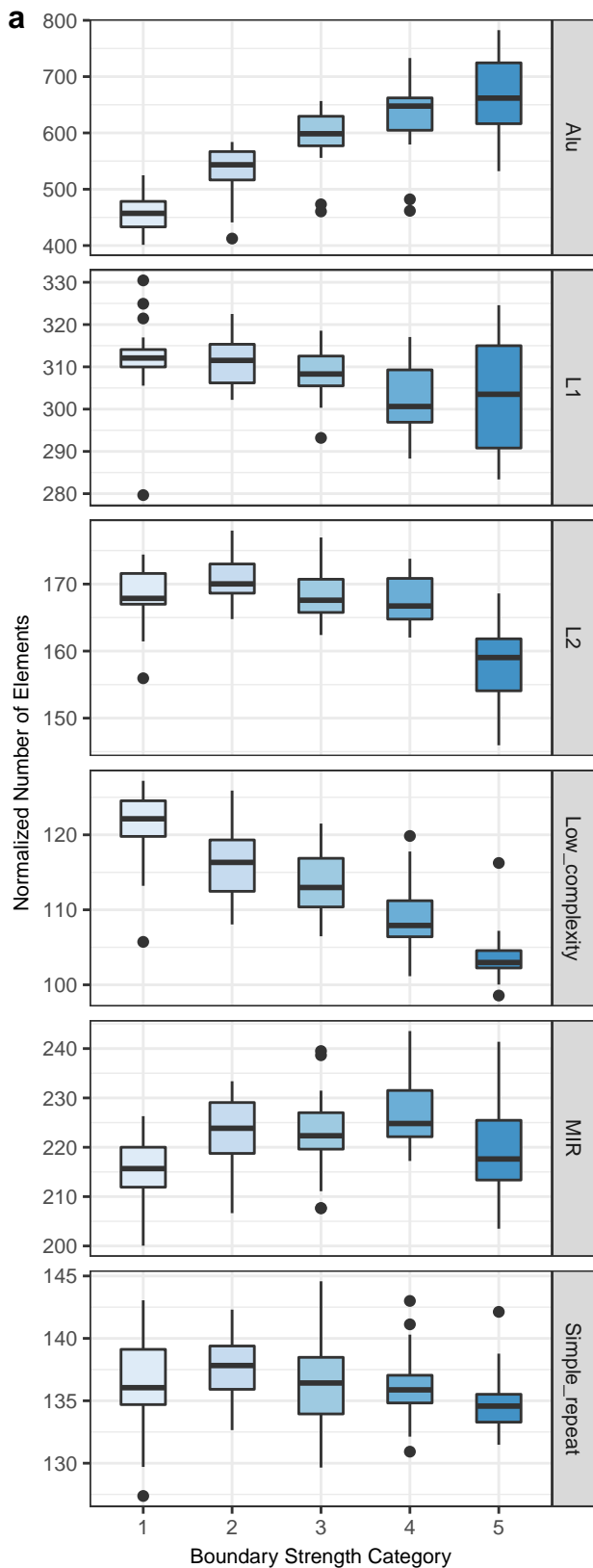


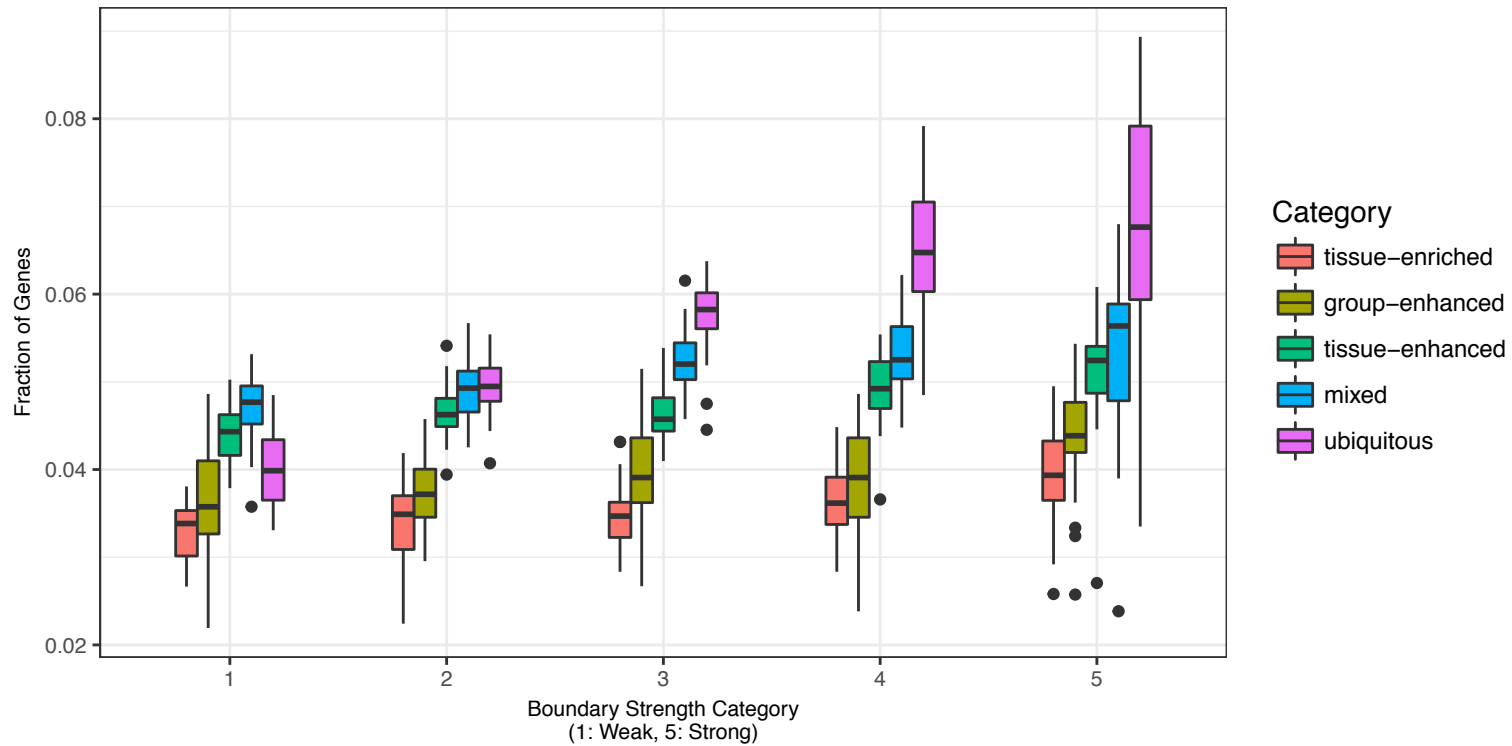
Supplementary Figure 4

Sequencing depth = 120 million read pairs



Supplementary Figure 5



a**b**