

1 **A pathway-centric view of spatial proximity in the 3D nucleome across cell lines**

2

3 Hiren Karathia¹, Carl Kingsford², Michelle Girvan³, Sridhar Hannenhalli^{1*}

4 ¹Center for Bioinformatics and Computational Biology, University of Maryland, College Park,
5 MD

6 ²Computational Biology Department, Carnegie Mellon University, Pittsburgh, PA

7 ³Department of Physics, University of Maryland, College Park, MD

8

9 Authors' emails:

10 Hiren Karathia – hiren@umiacs.umd.edu

11 Carl Kingsford – carlk@cs.cmu.edu

12 Michelle Girvan - mgirvan@gmail.com

13

14 *Corresponding author

15

16 Sridhar Hannenhalli

17 3104G Biomolecular Sciences Building (#296)

18 University of Maryland, College Park, MD 20742, USA

19 301 405 8219 (v) 301 314 1341 (f)

20 sridhar@umiacs.umd.edu

21

22 **Key words:** Hi-C, Chromatin structure, Pathways, Transcriptional regulation

23 **Abstract**

24 Spatial organization of the genome is critical for condition-specific gene expression. Previous
25 studies have shown that functionally related genes tend to be spatially proximal. However,
26 these studies have not been extended to multiple human cell types, and the extent to which
27 context-specific spatial proximity of a pathway is related to its context-specific activity is not
28 known. We report the first pathway-centric analyses of spatial proximity in six human cell lines. We
29 find that spatial proximity of genes in a pathway tends to be context-specific, in a manner consistent
30 with the pathway's context-specific expression and function; housekeeping genes are ubiquitously
31 proximal to each other, and cancer-related pathways such as p53 signaling are uniquely proximal in
32 hESC. Intriguingly, we find a correlation between the spatial proximity of genes and interactions of
33 their protein products, even after accounting for the propensity of co-pathway proteins to interact.
34 Related pathways are also often spatially proximal to one another, and housekeeping genes tend to be
35 proximal to several other pathways suggesting their coordinating role. Further, the spatially proximal
36 genes in a pathway tend to be the drivers of the pathway activity and are enriched for transcription,
37 splicing and transport functions. Overall, our analyses reveal a pathway-centric organization of the 3D
38 nucleome whereby functionally related and interacting genes, particularly the initial drivers of pathway
39 activity, but also genes across multiple related pathways, are in spatial proximity in a context-specific
40 way. Our results provide further insights into the role of differential spatial organization in cell type-
41 specific pathway activity.

42 **Introduction**

43 Recent advances in Chromosome Confirmation Capture (3C) and its high throughput
44 derivative, Hi-C, have enabled genome-wide identification of spatially proximal genomic
45 regions [1-3]. Comparative analysis of Hi-C data across cell lines and species reveals a
46 conserved framework of 3D architecture, represented by topologically associating domains
47 (TADs) and further context-specific variation in distal interactions [4].

48 Among other things, these 3D maps of chromosomes help explain, in part, spatio-temporal
49 regulation of gene expression by distal enhancers, aided by long-range DNA looping [5-7].
50 Similarly, previous studies have shown that groups of spatially clustered enhancers exhibit co-
51 activity across cell types and this co-activity is reflected in co-expression of proximal genes,
52 which are often functionally related [8, 9]. More specifically, genes involved in the same
53 pathway have been shown to be spatially proximal in *Saccharomyces cerevisiae* [10, 11],
54 *Plasmodium falciparum* [12] and *Homo sapiens* lymphoblastoid cell lines [13]. However,
55 these previous studies have not been extended to multiple human cell lines, and it is not clear
56 to what extent spatio-temporal activity of pathways is related to the spatial proximity of the
57 constituent genes. More generally, the broader characterization of physical proximity of
58 genes in the context of functional pathways is missing and could reveal organizing principles
59 underlying spatial proximity of pathway genes as they relate to pathway activity.

60 In this work, we perform a comparative pathway-centric analysis of Hi-C-derived spatial
61 proximity data in 6 ENCODE [14] cell types - *HEK293* [15], *hESC* [4], *IMR90*, *BT483*[16],
62 *GMO6990*[16], and *RWPE1*[17], each with replicate data. Our analysis of two large sets of
63 pathways – KEGG [18], and NetPath [19] reveals several properties of spatial proximity of
64 pathway genes: We find that in general, genes in a pathway tend to be spatially proximal and
65 this tendency is even greater for gene pairs that belong to multiple pathways. Our expression
66 analysis shows that genes that are co-localized in nuclear space with other genes have higher
67 expression, and this effect is especially prominent when they are proximal to a gene in the
68 same pathway. We also found that spatial proximity of pathway genes is strongly correlated
69 with cell type-specific pathway activity. As an expected corollary, housekeeping genes, by
70 virtue of being ubiquitously active, exhibit ubiquitous spatial proximity. Surprisingly though,
71 we found that the protein products of spatially proximal genes in a pathway have a
72 significantly greater tendency to physically interact than various controls. Functional
73 enrichment analysis suggests that spatially proximal pathway genes are enriched for specific
74 functional classes such as transcription factor and transmembrane genes, and they occupy

75 higher levels in the regulatory hierarchy. Finally, we look at higher-level spatial organization
76 of functional pathways by quantifying spatial proximity for all pairs of pathways. Using this,
77 we identify a network of spatially proximal pathways that is consistent with their functional
78 roles.

79 Overall, this first comprehensive pathway-centric analysis of spatial proximity in multiple
80 human cell lines shows a strong link between spatial proximity and context-specific gene
81 expression and pathway activity. Our analysis also reveals surprising links between spatial
82 proximity and interaction between the corresponding protein products. Functional analysis of
83 proximal genes within pathways strongly suggests a regulatory hierarchical bias in physical
84 proximity of pathway genes. Taken together, these results are consistent with a mechanism in
85 which early regulatory components of a pathway are brought into spatial proximity in a
86 condition specific manner.

87 **Results**

88 **1. Software pipeline for Hi-C processing - overview**

89 Fig. 1 shows the overall pipeline that, starting from the raw reads obtained from a Hi-C
90 experiment, produces significant pair-wise gene interactions. Details of the pipeline are
91 provided in the Methods section, and we highlight a few pertinent features here. The pipeline
92 allows the user to select a resolution at which significant interactions are identified. We
93 performed our analysis at 100 kb resolution because we are interested in gene-centric
94 interactions and 100 kb is expected to cover ~1 gene; smaller resolution yields fewer
95 significant interactions due to loss in power and larger resolution results in ambiguous gene-
96 segment mapping. We have further discussed the choice of the resolution later in the
97 discussion section. The pipeline uses the normalization step of the Homer tool [20] to control
98 for the genomic distance-dependent features of Hi-C counts (that proximal genomic regions
99 are more likely to interact). Using this pipeline, we processed 6 sets of pooled replicates for 6
100 ENCODE cell lines - HEK293, hESC, IMR90, BT483, GM06990, and RWPE1. Table 1 shows
101 the data obtained for the 6 cell types at the default interaction p-value threshold of 0.001 and
102 $FDR \leq 0.1$.

103 **Assessment of spatial proximity of pathway genes**

104 We quantify cell-line-specific spatial proximity of genes in a pathway using *edge fraction (EF)*,
105 which is essentially the fraction of all possible gene pairs in the pathway that are spatially
106 proximal. This measure was previously shown to be effective [10]. We then quantify

107 significance of *EF* based on a sampling procedure (see Materials & methods), obtaining a Z-
108 score and the corresponding p-value and multiple testing corrected q-value; a higher Z-score
109 is indicative of spatial proximity of the pathway genes above expectation. We studied 164
110 KEGG pathways [21] with at least 10 genes and estimated their spatial proximity Z-scores in 6
111 cell lines. Fig. 2 shows that overall biological pathways tend to be spatially proximal,
112 consistent with previous reports [10]. Interestingly, as shown in Fig. 2, we found that spatial
113 proximity for subsets of genes shared between two pathways is even greater, suggesting that
114 such genes, coordinating multiple pathways, may be under a greater constraint to be spatial
115 proximal.

116 Fig. 3 shows the cell-type-specific Z-scores for a representative set of pathways and
117 Supplementary Fig. 1 shows the same for all pathways. Consistent with the fact that the KEGG
118 database is dominated by essential and broad cellular processes, we found that spatially
119 proximity of KEGG pathways are not only generally high (Fig. 2), but a large fraction of
120 pathways exhibit a significant level of spatial proximity in many cell types (Fig. 4). In
121 particular, given the ubiquitous expression and function of housekeeping genes, we tested
122 whether these genes tend to be ubiquitously spatially proximal or whether their ubiquitous
123 expression is decoupled from their spatial proximity to one other. Based on 3800
124 housekeeping genes [22], we found that housekeeping genes exhibit significant spatial
125 proximity to other housekeeping genes in 5 out of 6 cell lines tested (Fig. 3). Other
126 ubiquitously proximal pathways include *Metabolism of xenobiotics by cytochrome P450, 1-*
127 *and 2-Methylnaphthalene degradation* and *Gamma-hexachlorocyclohexane degradation*, all
128 involved in drug metabolisms in animals.

129 Several cases of cell-type-restricted spatial proximity are worth noting. For instance, ‘Cell
130 cycle’ genes are expected to be active in pluripotent stem cells. This pathway is significantly
131 proximal in only two cell types, one of which is the human embryonic stem cell (hESC).
132 Cytokine-mediated signaling is critical in immune response and consistently, Cytokine-
133 cytokine receptor interactions are uniquely proximal in immune B-cell (GMO6990); B-cell-
134 specific proximity of cytokines CCL23 and CCL4 is consistent with their known role in
135 increased monocyte recruitment during inflammation [23]. Likewise we found the *Androgen*
136 *and estrogen metabolism* pathways to be proximal in breast cancer cell lines, where the role
137 of this pathway is well known [24]. Interestingly, we found the *Androgen-Estrogen receptor*
138 pathway to be proximal in Kidney cell line as well (Z-score = 6.5, FDR = 0.12), consistent with
139 the role of this pathway in glucuronidation activity that involves communication between
140 thyroid and kidney [25]. Unexpectedly, we see proximity of ‘*Type-II diabetic mellitus*’ in lung

141 fibroblast-derived IMR90. This is however consistent with recently observed connections
142 between diabetes and lung functions [26]. Finally, one the only two cell lines where
143 Melanogenesis genes is found to be proximal are the prostate epithelial RWEP1 and mammary
144 epithelial BT483; melanogenesis in prostate epithelial cells has been previously reported [27].

145 As an additional set of pathways, we examined a set of annotated cancer-related pathways
146 from NetPath [19]. As shown in Supplementary Fig. 2, genes in cancer-related pathways
147 exhibit much more subdued spatial proximity patterns in cell lines not derived from primary
148 tumors. Another noticeable trend is that many pathways exhibit high intra-pathway gene
149 proximity in hESC, consistent with previously noted similarities between cancer and stem cell
150 processes [28]. We found that *IL3-signaling* pathway is uniquely proximal in hESC
151 consistent with its known role in cell cycle control [29].

152 As a useful resource, in Supplementary Table 1 we have provided, for all pathways considered,
153 the specific gene pairs that were proximal in different cell lines.

154

155 **2. Spatial proximity and gene expression**

156 Next, we assessed whether spatial proximity of genes is correlated with their expression levels.
157 Furthermore, we also assessed among the spatially proximal genes whether belonging to same
158 pathway has any association with expression level. This analysis was done in all 6 cell lines
159 using RNA-seq data available in GEO (see Materials & Methods). We compared cell type
160 specific expression levels for three disjoint groups of genes (Fig. 5). The first group consisted
161 of genes that are proximal to another gene in the same pathway (*proximal-intra-pathway*).
162 The second group consisted of genes proximal to another gene but excluding *proximal-intra-*
163 *pathway* genes; this group was designed to assess whether shared pathway membership
164 impinges on gene expression. The last group consisted of all other genes not proximal to any
165 gene (*non-proximal-generic*). Fig. 6 shows that in pooled result from all 6 tissues, while genes
166 that are spatially proximal to other genes have a greater expression than non-proximal genes,
167 the expression is greater for genes that are proximal to a gene in the same pathway. The
168 results for individual cell lines are qualitatively similar and are shown in Supplementary Fig.
169 3.

170 We directly assessed the correlation between cell type specific pathway proximity and pathway
171 activity. We used two measures to approximate pathway activity (see Materials & Methods).

172 The first measure captured the mean expression of all genes in the pathway and the second
173 measure captured the ratio of mean expressions of proximal and non-proximal genes in the
174 pathway. Each measure was converted into a based on random sampling (see Materials &
175 Methods). As shown in Fig. 7, spatial proximity is highly correlated with mean pathway
176 expression (Spearman rho = 0.77, p = 5.9e-62). Note that this analysis does not rely on any Z-
177 score cutoff.

178 Our observed highly positive correlation between spatial proximity of a pathway genes and the
179 pathway activity would imply that the pathway genes are localized in active compartments and
180 not in repressive or inactive compartments (in which case spatial proximity would result in
181 pathway suppression – an opposite effect). A previous paper [48] has shown that in
182 GM06990, relative to the inactive (B) compartment, the active (A) compartment (1) is highly
183 enriched for genes, (2) has higher expression, (3) has more accessible chromatin, and (4) is
184 loosely packed, i.e., has fewer interactions. To assess our hypothesis that interacting pathway
185 genes are preferentially in A-compartment, we estimated the compartmentalization in
186 GM06990 cell line using HOMER tool and found 18256 (75%) genes in A-compartment and
187 6371 (25%) genes in B-compartment, consistent with previous report. Next, we compared for
188 4 groups of genes (defined in Fig. 5) the tendency to belong to A-compartment. As shown in
189 Supplementary Fig. 4 there is a robust monotonic trend whereby the genes proximal to
190 another gene in the same pathway have the greatest tendency to belong to A-compartment
191 and both spatial proximity and co-pathway membership contribute to this tendency.
192 Specifically there a large difference between the first class (proximal-intra-pathway) and the
193 non-proximal genes (p-value = 1.6E-53, Odds-ratio = 6.3). Thus, it seems that our observed
194 positive correlation between pathway proximity and activity is due to the fact that most
195 within-pathway proximal genes are in A-compartment.

196

197 **3. Spatial proximity and protein-protein interaction**

198 We have found that genes in a pathway tend to be spatially proximal. Previous studies have
199 shown that the proteins in a pathway have a greater tendency to physically interact with each
200 other [30]. We therefore directly assessed the correlation between spatial proximity of a gene
201 pair and the physical interaction of their products. We obtained the protein-protein
202 interactions (PPI) from HPRD database [31] and STRING [32]. Fig. 8 shows, for pooled data
203 from all the 6 cell lines, for each of the 5 groupings of gene pairs (defined in Fig. 5), the
204 fraction of all gene pairs that have evidence for physical interaction (fractional PPI). Results

205 for other individual cell lines are provided in Supplementary Fig. 5. Taken together, these data
206 suggest that, both pathway membership and spatial proximity of a gene pair is equally
207 associated with PPI between their products, and the effect is partly independent, that is, PPI
208 tendency is much greater for gene pairs that are both in the same pathway and physically
209 proximal.

210 **4. Functional enrichment for spatially proximal genes in pathways**

211 Only a small fraction of genes in a pathway are spatially proximal to other genes in the same
212 pathway. This could partly be due to low coverage and false negatives in the Hi-C derived
213 interactions. However, given the differences between proximal and non-proximal genes above,
214 we investigated whether specific functional terms are enriched among *proximal-intra-*
215 *pathway* genes relative to non-proximal-intra-pathway genes (see Fig. 5 for definition). We
216 pooled the data for all pathways into two groups and performed the functional enrichment
217 analysis using GOrilla software [33]. Fig. 9 suggests that regulatory (both RNA processing,
218 and splicing) and protein binding functions are highly enriched among the spatially proximal
219 pathway genes. We emphasize that our background set of genes in this analysis included those
220 that are in the pathways and are spatially proximal to some other gene not belonging to the
221 same pathway. Thus the observed functional enrichment is not due to spatial proximity alone
222 and is a specific property of spatially proximal genes within pathways.

223 **5. Spatial proximity and regulatory hierarchy**

224 The comparative analyses of biological properties of spatially proximal pathway genes relative
225 to other pathway genes thus far suggests that the upstream genes in a pathway may be more
226 likely to be spatially proximal, ensuring their robust expression and consequently robust
227 pathway activity. We derived the hierarchical level of all pathway genes based on directed
228 pathway edges (see Materials & Methods), and compared the hierarchical levels of proximal
229 and non-proximal genes, pooled overall pathways. Our approach for assigning hierarchical
230 level does not partition the pathway into a strict hierarchy, and can accommodate cycles. We
231 found that, in 3 of the 6 cell lines, the hierarchical levels of proximal genes were higher than
232 the rest; Wilcoxon test p-values: IMR90 ($p = 4.4E-04$), GM06990 ($p = 9.7E-04$) and hESC (p
233 $= 0.05$). Lack of significance in other cell lines (p -values = 0.07, 0.16, and 0.2) may be
234 attributed to due to insufficient data; when we pool the data from the other 3 cell lines, the
235 result is significant ($p = 1.41E-05$). If we apply Chi-square test based on hierarchy level of 2 as
236 the partitioning criterion, an additional cell line yields significance and importantly the odds
237 ratio in all six cell lines range favorably from 1.7 to 3.2.

238 **6. Spatial proximity between pathways**

239 Next, we investigated, whether certain pairs of pathways might occupy neighboring spaces in
240 the nucleus, suggesting their functional relatedness. Analogous to intra-pathway spatial
241 proximity estimation, for each pair of pathways, after excluding the common genes, we
242 obtained the number of interactions between genes across pathways and estimated its Z-score
243 based on 1000 controlled random samplings (for computational tractability) of 2 sets of genes
244 representing the 2 pathways. Using Z-score > 2 as the threshold, we identified a total of 3109
245 pathway pairs across all cell lines, 73 of which were proximal in at least 4 tissues and 20 were
246 proximal in at least 5 tissues (Supplementary Table 2).

247 Next, analogous to intra-pathway proximity analysis, we estimated inter-pathway spatial
248 proximity between housekeeping genes and all other KEGG pathways. We found that
249 housekeeping genes are significantly proximal to 95 out of 164 KEGG pathways, at Z-score > 2
250 threshold in at least 1 cell line; the full distribution of number of housekeeping-proximal
251 pathways shared in 1 or more cell lines is provided in Supplementary Fig. 6, and pathways
252 proximal to housekeeping genes in 4 or more cell lines are provided in Supplementary Table
253 2. These results suggest a central role housekeeping genes may play in the genome
254 organization as well as in coordinating the activities of other pathways.

255 Several of the pathway pairs deemed to be spatially proximal in our analysis (Supplementary
256 Table 2) are very likely to be functionally related, for instance, pathways for metabolism of
257 various amino acids and derivatives thereof. Importantly, however, these data also reveal non-
258 trivial relationships, which have some support in literature. We discuss a few next.

- 259 • Consistent with the proximity of steroid hormone metabolism and proteolysis, direct
260 functional links between these two processes have been noted [34].
- 261 • Links between amino acid metabolism and cell cycle and cancer (Chronic myloid leukemia
262 in our case) is not surprising given the metabolic requirements during cell division and
263 growth.
- 264 • JAK-STAT pathway and SNARE complexes are found to be proximal. While we did not
265 find a direct link between the two, they are known to be co-targeted by TGF-beta signaling
266 [35].
- 267 • Abnormal *steroid hormone metabolism* has been reported in *Myeloid leukemia* patients
268 [36], which is consistent with the detected proximity of these two pathways.

- 269 • *Butanoate metabolism* and *Glioma* are deemed to be spatially proximal in our analysis.
270 Treatment with butanoate (sodium butyrate) is known to induce differentiation of c6
271 Glioma cells [37].
- 272 • While a direct link between *Cholera* and *DRPLA*, deemed to be proximal, is not clear, we
273 note that *DRPLA* is a spinocerebral degeneration disease, and *cholera toxin subunit b* is
274 known to be dispersed to brain and spinal cord neurons [38].
- 275 • *Synthesis and degradation of ketone bodies* was found to be proximal to genes involved in
276 several different cancers, including pancreatic cancer, and cell cycle. Previous papers have
277 shown growth inhibitory effects of *ketone* bodies on pancreatic cancer could be mediated
278 by reduced c-Myc expression [39]. Moreover c-Myc is overexpressed in pancreatic cancer
279 [40] and its inhibition has been shown to result in regression of *lung cancer* [41] and
280 pancreatic cancer [42].
- 281 • *Glycosylphosphatidylinositol* (GPI) and *Huntington* (a neurological disease) were deemed
282 proximal. While a direct link between these two pathways is not clear, GPI-anchor cleavage
283 is known to modulate notch signaling and promoter neurogenesis [43].

284

285 Overall, these results suggest that spatial proximity between pathway genes is somewhat
286 associated with functional interactions between the pathways.

287 **Discussion**

288 In this work, we have presented the first comprehensive analysis of intra- and inter-pathway
289 spatial proximity in multiple *Homo sapiens* cell lines. Previous studies have shown that in
290 *Saccharomyces cerevisiae* [10], *Plasmodium falciparum* [12] and *H.sapiens* lymphoblastoid
291 cell line [13], broadly, functionally related genes tend to be spatially proximal. Our goal here
292 was to not only extend these previous observations to multiple human cell lines and assess the
293 relationships between spatial proximity and pathway activity based on gene expression, but
294 equally importantly, to further functionally characterize proximal genes within pathways and
295 examine higher-order physical and functional interactions between pathways.

296 Previous similar analysis in *S. cerevisiae* [10] are based on only inter-chromosomal segment
297 interactions, and *H. sapiens* lymphoblastoid cell line results [13] are based on low-resolution
298 (1 Mb) segments, which can result in spurious interactions at the level of individual gene loci.
299 Importantly, however, a greater tendency for genomically proximal regions to be spatially
300 proximal, i.e., autocorrelation, unless appropriately controlled for, can result in false positives
301 in inferring significant spatial proximity from Hi-C data. The absence of effective tools to

302 control for autocorrelation has forced previous studies to exclude intra-chromosomal
303 interactions from consideration, significantly impacting their statistical power [10]. In
304 contrast the Homer tool [20] satisfactorily controls of autocorrelation in estimating
305 significance, consequently enabling us to include all significant interactions, both inter- and
306 intra-chromosomal, in our analyses, while obviating an explicit control for inter-gene
307 distances in random sampling procedures. We note that despite an explicit control for
308 autocorrelation, genomically proximal regions (within 500 kb) have a slightly higher tendency
309 to be spatially proximal (data not shown). However, we explicitly tested if this biases our
310 pathway proximity assessment as follows. We compared intra-pathway gene distances with
311 those for randomly selected genes from the same chromosome and found that for none of the
312 pathways there was a significant difference between the two sets of distances.

313 We have performed our analyses based on 100 kb resolution to detect interactions. However,
314 we have also assessed the impact of using a higher resolution of 10 kb. We found that the
315 number of interactions detected was much greater when using 100 kb resolution, especially
316 the inter-chromosomal interactions, due to a greater statistical power for interaction detection
317 at this resolution. For example, for cell type HEK293, only 7889 intra- and only 14 inter-
318 chromosomal gene-gene interactions are detected using 10kb resolution and 43439 intra- and
319 1538 inter-chromosomal interactions at 100 kb. Moreover, at 10 kb resolution, large fractions
320 of detected interactions are within a gene, which does not contribute to our analyses. We note
321 that using a 100 kb resolution does not substantively influence the gene-gene interaction
322 inference, as only 4% of 100 kb segments have multiple genes. Thus, to maximize statistical
323 power with relative small fraction of ambiguous (but not necessarily biologically wrong) gene-
324 gene interaction calls, we chose to perform all downstream analyses based on 100 kb
325 resolution.

326 Hi-C data, like most genome-wide datasets, comes with a level of false positives, as does the
327 pathway data. This is an important issue, and one that cannot be addressed by computational
328 means alone. We have relied on the published interaction data that have undergone quality
329 control measures, and used robust tool with recommended controls to perform the analyses.
330 For pathway proximity, we have relied on a well-controlled randomized gene set to assess the
331 significance of proximity. Despite the controls, a certain fraction of data is likely to be false
332 positive. However, the noise in the data, as long as the tests are properly controlled, is not
333 likely to generate strong consistent signals across multiple cell types simply by chance.

334 We performed a number of checks to ensure the robustness of our conclusions against several
335 potential biases. First, note that in quantifying the significance of the pathway spatial
336 proximity, we control for the lengths of the pathway genes. Second, we ensured that our ability
337 to detect the spatial proximity of a pathway is independent of the number of genes in the
338 pathway (Spearman correlation between pathway size and spatial proximity Z-score was
339 statistically insignificant). A recent paper [9] has suggested a link between detection of Hi-C
340 interaction of a gene and a gene's codon usage (which is related to its expression) and the GC
341 content of the gene's genomic locus. Third, we ascertained that the codon-usage does not bias
342 the detection of Hi-C interactions (Spearman correlation = -0.09).

343 Finally, with regards to the potential GC bias, indeed highly interacting loci tend to have lower
344 GC composition, as noted previously [11]. Although the GC composition near the restriction
345 sites can present a technical bias [44], several papers also suggest that GC composition may be
346 an inherent property of the physical proximity [45]. Therefore, it seems that an explicit control
347 for GC content may not be ideal. However to explore the extent to which GC-controlled
348 analyses would affect our results, for HEK293 cell line we reprocessed the data with specific
349 GC control option provided in HOMER tool and compared the downstream results of our
350 analyses with and without the GC control. We found that 80% of interactions detected without
351 GC control are also detected with GC control and overall correlation between gene-wise
352 interaction degrees between the two is 0.92. Next, we re-estimated the pathway proximity Z-
353 scores based on GC-controlled interactions and found those to be highly correlated with the Z-
354 score based on uncontrolled interaction detection (Spearman rho = 0.85). Lastly, we re-
355 calculated the correlation between pathway proximity and pathway activity and found that too
356 is as significant as the estimated without GC control (Spearman rho = 0.80, relative to 0.83
357 without GC control). Thus, our conclusions are not substantively biased by these various
358 potentially confounding factors.

359 Our primary resource of biological pathways – KEGG, is dominated by essential and broadly
360 utilized cellular pathways, and therefore it is encouraging to see that by and large KEGG
361 pathways are not only highly significantly proximal (Fig. 2), but a large fraction of these
362 pathways are proximal in multiple cell lines (Fig. 4). However, as we show, this is not true for
363 a different set of pathways relevant to cancer, where the overall z-scores are much more
364 subdued, and spatial proximity is less ubiquitous (Supplementary Fig. 2). Despite, general
365 ubiquity of spatial proximity of KEGG pathways, we still see a strong correlation between cell
366 type-specific spatial proximity and pathway activity as approximated by gene expression (Fig.
367 7).

368 Previous studies have noted greater transcription in spatially clustered regions [8].
369 Independently, earlier studies have shown the existence of so called transcription factories
370 [46] - nuclear locales with enriched core transcriptional machinery components where
371 transcripts are synthesized and processed. Moreover, links between transcription factories
372 and chromatin organization have been noted [47]. Taken together, these previous results are
373 consistent with our observation that genes in spatial proximity with other genes have much
374 higher expression. However, interestingly, in addition to spatial proximity alone, functional
375 relationship between the spatially proximal genes, i.e., membership in the same pathway,
376 makes a small but significant additional contribution to the gene expression level (Fig. 6).
377

378 Our analysis reveals an unexpected association between spatial proximity of a gene pair and
379 the interaction between their protein products. Among the physically proximal gene pairs, the
380 well-annotated genes, i.e., annotated in some KEGG pathway, have much greater PPI
381 propensity than genes that do not belong to an annotated pathway ((III) vs. (IV) in Fig. 8);
382 this may be explained by a greater representation of well-studied genes in PPI databases.
383 Functionally related genes have been shown to have a greater propensity to physically interact
384 [49], consistent with our findings (compare (I) and (III) in Fig. 8). However, we found that
385 spatial proximity is independently associated with protein interaction in both pathway ((I)
386 versus (II) in Fig. 8) and non-pathway ((IV) versus (V) in Fig. 8) contexts. The gene-pairs that
387 are both spatially proximal and belong to same pathway have the highest PPI propensity ((I)
388 in Fig. 8). These trends are identical in all 6 cell lines, suggesting that both spatial proximity
389 and pathway membership contribute independently to PPI propensity.

390 Scrutinizing each of the pathways, we found the spatially proximal genes in a pathway to have
391 distinguishing functional characteristics relative to other pathway genes that are not spatially
392 proximal to any other gene in the same pathway. In terms of biological processes (Fig. 9), such
393 genes are overwhelmingly involved in transcription, splicing and intracellular transport and
394 localization. Interestingly, the genes in a pathway that are spatially proximal to other genes in
395 the same pathway tend to occupy a higher level in the regulatory hierarchy, related to other
396 genes in the pathway, that also are spatially proximal to other genes but none in the same
397 pathway. Overall, these results ascribe, for the first time, a special functional status to spatially
398 proximal genes in pathways – such genes tend to perform higher-level regulatory functions.

399 We found that housekeeping genes, consistent with their ubiquitous expression and activity,
400 tend to be broadly and highly spatially proximal. Previous studies have observed clustering of
401 housekeeping genes into so call transcription factories [50], and have suggested that

402 interactions between housekeeping genes may play a role in the spatial organization of the
403 chromatin [51]. Our results confirm these previous observations through the first genome-
404 wide assessment of spatial proximity of housekeeping genes in multiple cell lines. In addition
405 to spatial proximity of housekeeping genes, we also found that as a group housekeeping genes
406 are ‘centrally’ located in the nucleus and act as a link between numerous other pathways. A
407 mechanistic interpretation of this intriguing observation, as well as of the causal link between
408 the expression of housekeeping genes and their spatial proximity, will require further analysis.

409 **Conclusion**

410 Overall, based on the first comprehensive pathway-centric analysis of spatial proximity in
411 multiple cell lines, our results suggest that (i) context-specific regulation of pathways is
412 associated with their context-specific spatial proximity; in doing so, our analysis provides
413 mechanistic insights into cell type-specific activity of certain pathways, (i) spatial-proximity of
414 pathway genes is associated with physical interaction among their gene products, and (iii)
415 specific classes of genes within pathways, likely occupying higher regulatory levels, have a
416 greater tendency to be spatially proximal. Our results also provide insights into correlated
417 activity of multiple pathways by showing that the genes in these pathways are spatially
418 proximal.

419

420 **Methods**

421 **Hi-C processing pipeline.** We downloaded paired-end Hi-C raw reads FASTQ files of
422 sample replicates for the following tissues from GEO database (www.ncbi.nlm.nih.gov/geo):
423 (i) HEK293 (GSM1081530, GSM1081531) [15], (ii) IMR90 (GSM1055800, GSM1055801) [52],
424 (iii) hESC (GSM862723, GSM892306) [4], (iv) GMO6990 (GSM1340639) [16], (v) RWPE1
425 (GSM927076) [17] and (vi) BT483 (GSM1340638, GSM1340637) [16]. We mapped the reads
426 onto hg19 human genome using *BWA* tools [53] with default parameters. The resulting SAM
427 files were converted to BAM files using “*samtools view*” program [54] and processed to
428 removing PCR duplicates using “*samtools sort*” and “*Picard*” tools.

429 We then processed the non-redundant reads for Hi-C analysis using various HOMER tools
430 [20]: we ran “*makeTagDirectory*” program using options “*tbp -1*” (to ensure that any
431 genomic location is mapped by a unique read), “*-restrictionSite*” (only keep reads if both ends
432 of the paired-end read have a restriction site within the fragment length estimate 3' to the
433 read), “*-removePEbg*” (removing read pairs separated by less than 1.5x the sequencing insert

434 fragment length) , “*-removeSelfLigation*” (remove re-ligation events), and “*-removeSpikes*”
435 (remove high tag density regions).

436 **Normalization of Hi-C interactions.** Having output from the previous steps, next we
437 normalized the data to create background of the Hi-C interactions at 100kb resolution using
438 HOMER “*analyzeHiC*” program. The program (i) divides the genome into 100kb regions, (ii)
439 calculates total read coverage in each region, (iii) calculates the fraction of interactions
440 spanning any given distance with respect to read depth, (iv) optimizes a read count model to
441 assign expected interaction counts in regions with uneven sequencing depth and (v) calculates
442 variation in interaction frequencies as a function of distance. For fragments i and j , the
443 procedure the estimated expected number of reads supporting the interaction as:

$$e_{i,j} = \frac{f(i - j) * n_i * n_j}{N}$$

444 where N = total number of reads, n = number of reads in a region, and f represents the
445 expected frequency of Hi-C reads as a function of distance.

446 The background model were created for the entire data for a given sample then applied for
447 selection of Hi-C interacting reads at default p-value cutoff of 0.001 using “*-interaction*”
448 parameters using “*analyzeHiC*” program. We further filtered the Hi-C interactions using FDR
449 cutoff ≤ 0.1 and considered the resulting set as significant Hi-C interactions. The significant
450 Hi-C interacting reads were then passed through “*annotateInteractions*” program for
451 mapping them onto the annotated genomic features (i.e., 5’UTR, CDS, introns, exons,
452 intergenic etc.), from which, we selected only those interactions mapping to gene regions (i.e.,
453 Promoter, 5’UTR, exon, intron and 3’UTR), resulting in a comprehensive set of spatially
454 proximal gene pairs (SGP) for each tissue.

455 **Codon usage.** To assess the effect of codon bias on interaction detection, we used
456 Biopython’s “*SeqeUtils::ModuleCodonUsage*” package and calculated Codon Adaptation
457 Index (CAI) for each gene of the *H. sapiens* genome.

458 **Pathway datasets.** A ‘pathway’ for our purpose is a set of genes. We downloaded two sets of
459 pathways from KEGG [21] and NetPath [19] databases, only retaining those with at least 10
460 genes, resulting in 164 and 32 pathways respectively. We also included the set of 3800
461 housekeeping genes [22], as an additional ‘pathway’.

462 **Defining classes of edges and nodes relative to pathways.** For various analyses we
463 have defined 5 disjoint sets of genes and gene pairs in the context of a pathway and spatial
464 proximity (see Figure 5 for illustration). All spatially proximal gene pairs were partitioned into
465 three groups: *proximal-intra-pathway*: intra-pathway spatially proximal gene pairs,
466 *proximal-inter-pathway*: Spatially proximal pair of genes where each gene is in a different
467 pathway, and *proximal-generic*: Spatially proximal pair of genes such that at least one of the
468 genes is not in any pathway. Similarly all non-proximal gene pairs were partitioned into two
469 categories: *non-proximal-intra-pathway*: a non-proximal gene pair within a pathway, and
470 *non-proximal-generic*: any pair of genes that are not spatially proximal to any other gene and
471 not within any pathway.

472 **Edge fraction and its significance in estimating pathway spatial proximity.** Intra-
473 pathway gene spatial proximity was estimated as *Edge Fraction (EF)* – number of pairwise
474 gene-gene interaction in the pathway normalized by the number of total possible interactions.
475 To quantify significance of the *EF* for pathway with N genes, we randomly sampled N genes,
476 such that number of genes in each chromosome is identical to real pathway, and each sampled
477 gene's length was within 20% of the matched pathway gene's length (this controls for length-
478 based bias in interaction detection). We generated 1000 such samples and calculated 1000
479 corresponding *EFs*. We then obtained the Z-score corresponding to the *EF* for actual pathway
480 relative to 1000 controls. We estimated the Z-score for each pathway in each cell line resulting
481 in 165 x 6 Z-score matrix.

482 **Inter-pathway proximity.** Analogous to the intra-pathway proximity analysis above, for a
483 pair of pathways, after excluding the shared genes, we estimate the edges between genes in
484 two pathways, and estimate its significance based on randomly sampling 2 gene sets (instead
485 of 1 as above), with identical controls as above. We thus estimated a Z-score for inter-pathway
486 proximity for all pairs of pathways.

487 **Processing RNA-Seq data.** We downloaded raw RNA-Seq FASTQ files for three of the
488 tissues in which matching RNA-Seq was available: (i) HEK293 (GSM1081534, GSM1081535)
489 [15], (ii) IMR90 (GSM1154029) [52] and (iii) RWPE1 (GSM927074) [17], (iv) hESC
490 (GSM758566) [55], (v) GM06990 (GSM958747) [56] and (vi) BT483 (GSM1172854) [57]. The
491 raw FASTQ reads were mapped and processed up to de-duplication steps using the same
492 pipeline that we used to apply for Hi-C data analysis, and then used to quantify expression
493 levels using cufflinks [58] tool with default parameters, yielding gene-wise RPKM values.

494 **Significance of pathway activity.** For genes with multiple transcripts, we take the
495 maximum expression over all transcripts for the genes. For a pathway, we estimated its
496 activity as the average expression of the genes in the pathway; we selected only the *proximal-*
497 *intra-pathway* genes to estimate activity. To quantify significance of the pathway activity, we
498 followed a sampling approach similar to that for estimating EF significance above, yielding a
499 z-score for each pathway and cell line.

500 **Estimating regulatory hierarchy.** The regulatory analysis of a pathway requires a
501 directed graph in which all defined interactions suggest direction of the signal flow. KEGG
502 does not provide the directed graph by default. Therefore, we downloaded *KGML* (KEGG
503 Markup Language) and parsed the files using *KEGGgraph* package in Bioconductor and
504 created directed graph for each pathway using *NetworkX* package in Python. In order to
505 investigate hierarchy of genes we first assigned a synthetic root connecting to all pathway
506 nodes with zero in-degree, and then calculated the shortest path length (SPL) from root to
507 every other node. Low SPL indicates higher level of hierarchy. For Chi-square and Fisher tests
508 (see Results), we created two gene sets: (i) genes with SPL = 1 (top level) and (ii) genes with
509 SPL ≥ 3 (lower hierarchy). We pooled these two sets across all pathways.

510 **Author contributions**

511 S.H. conceived the project with input from C.K. and M.G. H.K. designed and implemented the
512 software pipeline and performed all analyses. S.H. wrote the manuscript with help from all
513 authors.

514 **Additional Data Files**

515 “Supplementary File” contains supplementary figures and tables. File ‘Supplementary Fig1’
516 contains the Supplementary Fig. 1.

517 **Acknowledgement**

518 This work was funded by NIH HG007104 to C.K. and NIH GM100335 to S.H.

519 **Competing interests**

520 None of the authors has any competing interests.

521

522

523 **Figure Legends**

524 **Figure 1. Overview of the Hi-C processing pipeline and flow of downstream**
525 **analysis (see Methods for details).**

526 **Figure 2. Z-score distribution of intra-pathway spatial proximity.** The figure shows
527 the distributions of spatial-proximity z-scores for three sets of gene sets, pooled from 6 cell
528 lines. Blue: KEGG pathway. Red: Random gene-sets matching each KEGG pathway controlled
529 for gene lengths and chromosomal distributions. Green: intersection of each pair of KEGG
530 pathways.

531 **Figure 3. Spatial proximity of *intra-pathway* genes in selected KEGG pathways**
532 **across six cell lines.** Note the ubiquitous proximity of Housekeeping genes, hESC-specific
533 proximity of p53 signaling pathway.

534 **Figure 4. *Intra-pathway* pathway proximity is shared across tissues.** The figure
535 shows number of pathways (Y-axis) with high *intra-pathway* proximity (Z-score ≥ 2) in
536 different number of cell lines (X-axis). The table values show number of pathways whose
537 *intra-pathway* genes proximity (Z-score ≥ 2) is unique to a cell line (diagonal) or shared
538 between a pair of cell lines (off-diagonal).

539 **Figure 5. Schematic defining sets of genes and gene-pairs.** Outermost rectangle
540 represents all genes. Inner rectangle represents all genes that are in spatial proximity to at
541 least one other gene. Circles represent annotated pathways. Nodes represent genes and edges
542 represent spatially proximal gene pairs. All spatially proximal gene pairs are partitioned into
543 three groups: ***proximal-intra-pathway***: intra-pathway spatially proximal gene pairs
544 (orange nodes and edges), ***proximal-inter-pathway***: Spatially proximal pair of genes
545 where each gene is in a different pathway (light blue nodes and edges), and ***proximal-***
546 ***generic***: Spatially proximal pair of genes such that at least one of the genes is not in any
547 pathway (red nodes and edges). Similarly all non-proximal gene pairs are partitioned into two
548 categories: ***non-proximal-intra-pathway***: a non-proximal gene pair within a pathway,
549 and ***non-proximal-generic***: any pair of genes that are not spatially proximal to any other
550 gene and not within any pathway (pairs of dark blue nodes).

551 **Figure 6. Spatial proximity and gene expression.** The figure shows box-plots of gene
552 expression (FPKM) values of the genes in three different groups (see Figure-5) pooled from all
553 6 cell lines. A~B Wilcoxon test p-values = $3.4E-35$. A~C Wilcoxon test p-value = $8.43E-72$. See
554 Supplementary Fig. 4 for results of 6 individual cells.

555 **Figure 7. Spatial proximity versus mean pathway expression.** This figure shows
556 scatter plot between proximity z-scores of pathways versus z-scores of expression values of the
557 proximal-intra-pathway genes pooled from all 6 cells. Spearman rho = 0.77, p-value: 5.90e-
558 62.

559 **Figure 8. Pathway membership, spatial proximity and PPI.** The figure shows 6 cells
560 pooled fraction of gene pairs (Y-axis) in different gene groups (X-axis; see Figure 5) whose
561 protein products physically interact. The inset shows Fishers tests p-values for various pair-
562 wise comparisons. See supplementary Figure 6 for other tissues.

563 **Figure 9. GO enrichment analysis.** This figure enriched GO terms ($-\log(q\text{-value}) \geq 2$)
564 ***proximal-intra-pathway*** genes relative to other spatially proximal genes (see Fig. 5).

565

566 **Table 1. Hi-C analysis summary.**

Cell line	Rep-1 Hi-C mapped fragments in million	Rep-2 Hi-C mapped fragments in million	# Unique genes in HiC interactions	# Unique genes in HiC interactions
IMR90	66.69	299.86	14550	101502
HEK293	92.03	210	11978	44977
hESC	42.49	234.33	10713	34531
GM06990	58.62	96.33	9019	11447
BT483	41.8	64.7	5948	5792
RWPE1	73.65	82.33	9416	14124

567

568

569 References

- 570 1. Shaw PJ: **Mapping chromatin conformation.** *F1000 biology reports* 2010, **2**.
- 571 2. Burton JN, Adey A, Patwardhan RP, Qiu R, Kitzman JO, Shendure J: **Chromosome-scale**
- 572 **scaffolding of de novo genome assemblies based on chromatin interactions.**
- 573 *Nature biotechnology* 2013, **31**:1119-1125.
- 574 3. van Berkum NL, Lieberman-Aiden E, Williams L, Imakaev M, Gnirke A, Mirny LA, Dekker J,
- 575 Lander ES: **Hi-C: a method to study the three-dimensional architecture of genomes.**
- 576 *Journal of visualized experiments : JoVE* 2010.
- 577 4. Dixon JR, Selvaraj S, Yue F, Kim A, Li Y, Shen Y, Hu M, Liu JS, Ren B: **Topological domains**
- 578 **in mammalian genomes identified by analysis of chromatin interactions.** *Nature*
- 579 2012, **485**:376-380.
- 580 5. Aguilar-Arnal L, Sassone-Corsi P: **Chromatin landscape and circadian dynamics:**
- 581 **Spatial and temporal organization of clock transcription.** *Proceedings of the National*
- 582 *Academy of Sciences of the United States of America* 2014.
- 583 6. Duggal G, Wang H, Kingsford C: **Higher-order chromatin domains link eQTLs with the**
- 584 **expression of far-away genes.** *Nucleic acids research* 2014, **42**:87-96.
- 585 7. Sanyal A, Lajoie BR, Jain G, Dekker J: **The long-range interaction landscape of gene**
- 586 **promoters.** *Nature* 2012, **489**:109-113.
- 587 8. Malin J, Aniba MR, Hannenhalli S: **Enhancer networks revealed by correlated DNase**
- 588 **hypersensitivity states of enhancers.** *Nucleic acids research* 2013, **41**:6828-6838.
- 589 9. Diamant A, Pinter RY, Tuller T: **Three-dimensional eukaryotic genomic organization**
- 590 **is strongly correlated with codon usage expression and function.** *Nature*
- 591 *communications* 2014, **5**:5876.
- 592 10. Wang H, Duggal G, Patro R, Girvan M, Hannenhalli S, Kingsford C: **Topological properties**
- 593 **of chromosome conformation graphs reflect spatial proximities within**
- 594 **chromatin.** In *Proceedings of the International Conference on Bioinformatics,*
- 595 *Computational Biology and Biomedical Informatics.* pp. 306-315. Washington DC, USA: ACM;
- 596 2013:306-315.
- 597 11. Dekker J: **GC- and AT-rich chromatin domains differ in conformation and histone**
- 598 **modification status and are differentially modulated by Rpd3p.** *Genome biology*
- 599 2007, **8**:R116.
- 600 12. Ay F, Bunnik EM, Varoquaux N, Bol SM, Prudhomme J, Vert JP, Noble WS, Le Roch KG:
- 601 **Three-dimensional modeling of the P. falciparum genome during the erythrocytic**
- 602 **cycle reveals a strong connection between genome architecture and gene**
- 603 **expression.** *Genome research* 2014, **24**:974-988.
- 604 13. Thevenin A, Ein-Dor L, Ozery-Flato M, Shamir R: **Functional gene groups are**
- 605 **concentrated within chromosomes, among chromosomes and in the nuclear**
- 606 **space of the human genome.** *Nucleic acids research* 2014, **42**:9854-9861.
- 607 14. **An integrated encyclopedia of DNA elements in the human genome.** *Nature* 2012,
- 608 **489**:57-74.
- 609 15. Zuin J, Dixon JR, van der Reijden MI, Ye Z, Kolovos P, Brouwer RW, van de Corput MP, van de
- 610 Werken HJ, Knoch TA, van IWF, et al: **Cohesin and CTCF differentially affect**
- 611 **chromatin architecture and gene expression in human cells.** *Proceedings of the*
- 612 *National Academy of Sciences of the United States of America* 2014, **111**:996-1001.
- 613 16. Dryden NH, Broome LR, Dudbridge F, Johnson N, Orr N, Schoenfelder S, Nagano T, Andrews
- 614 S, Wingett S, Kozarewa I, et al: **Unbiased analysis of potential targets of breast cancer**
- 615 **susceptibility loci by Capture Hi-C.** *Genome research* 2014, **24**:1854-1868.
- 616 17. Rickman DS, Soong TD, Moss B, Mosquera JM, Dlabal J, Terry S, MacDonald TY, Tripodi J,
- 617 Bunting K, Najfeld V, et al: **Oncogene-mediated alterations in chromatin**
- 618 **conformation.** *Proceedings of the National Academy of Sciences of the United States of*
- 619 *America* 2012, **109**:9083-9088.
- 620 18. Kanehisa M, Goto S, Sato Y, Kawashima M, Furumichi M, Tanabe M: **Data, information,**
- 621 **knowledge and principle: back to metabolism in KEGG.** *Nucleic acids research* 2014,
- 622 **42**:D199-205.

- 623 19. Kandasamy K, Mohan SS, Raju R, Keerthikumar S, Kumar GS, Venugopal AK, Telikicherla D,
624 Navarro JD, Mathivanan S, Pecquet C, et al: **NetPath: a public resource of curated signal**
625 **transduction pathways**. *Genome biology* 2010, **11**:R3.
- 626 20. Heinz S, Benner C, Spann N, Bertolino E, Lin YC, Laslo P, Cheng JX, Murre C, Singh H, Glass
627 CK: **Simple combinations of lineage-determining transcription factors prime cis-**
628 **regulatory elements required for macrophage and B cell identities**. *Molecular cell*
629 2010, **38**:576-589.
- 630 21. Kanehisa M, Goto S: **KEGG: kyoto encyclopedia of genes and genomes**. *Nucleic acids*
631 *research* 2000, **28**:27-30.
- 632 22. Eisenberg E, Levanon EY: **Human housekeeping genes, revisited**. *Trends in genetics :*
633 *TIG* 2013, **29**:569-574.
- 634 23. Starr AE, Dufour A, Maier J, Overall CM: **Biochemical analysis of matrix**
635 **metalloproteinase activation of chemokines CCL15 and CCL23 and increased**
636 **glycosaminoglycan binding of CCL16**. *The Journal of biological chemistry* 2012,
637 **287**:5848-5860.
- 638 24. Yager JD, Davidson NE: **Estrogen carcinogenesis in breast cancer**. *The New England*
639 *journal of medicine* 2006, **354**:270-282.
- 640 25. Basu G, Mohapatra A: **Interactions between thyroid disorders and kidney disease**.
641 *Indian journal of endocrinology and metabolism* 2012, **16**:204-213.
- 642 26. I AAE-A, Hamdy G, Amin M, Rashad A: **Pulmonary function changes in diabetic lung**.
643 *Egyptian Journal of Chest Diseases and Tuberculosis* 2013, **62**:513-517.
- 644 27. Kudva R, Hegde, P.: **Blue nevus of the prostate**. *Indian J Urol* 2010, **26**:301-302.
- 645 28. Dreesen O, Brivanlou AH: **Signaling pathways in cancer and embryonic stem cells**.
646 *Stem cell reviews* 2007, **3**:7-17.
- 647 29. Seita J, Weissman IL: **Hematopoietic stem cell: self-renewal versus differentiation**.
648 *Wiley interdisciplinary reviews Systems biology and medicine* 2010, **2**:640-653.
- 649 30. Kelley R, Ideker T: **Systematic interpretation of genetic interactions using protein**
650 **networks**. *Nature biotechnology* 2005, **23**:561-566.
- 651 31. Mishra GR, Suresh M, Kumaran K, Kannabiran N, Suresh S, Bala P, Shivakumar K, Anuradha
652 N, Reddy R, Raghavan TM, et al: **Human protein reference database--2006 update**.
653 *Nucleic acids research* 2006, **34**:D411-414.
- 654 32. Franceschini A, Szklarczyk D, Frankild S, Kuhn M, Simonovic M, Roth A, Lin J, Minguez P,
655 Bork P, von Mering C, Jensen LJ: **STRING v9.1: protein-protein interaction networks,**
656 **with increased coverage and integration**. *Nucleic acids research* 2013, **41**:D808-815.
- 657 33. Eden E, Navon R, Steinfeld I, Lipson D, Yakhini Z: **GORilla: a tool for discovery and**
658 **visualization of enriched GO terms in ranked gene lists**. *BMC bioinformatics* 2009,
659 **10**:48.
- 660 34. Umpleby AM, Russell-Jones DL: **The hormonal control of protein metabolism**.
661 *Bailliere's clinical endocrinology and metabolism* 1996, **10**:551-570.
- 662 35. Asirvatham AJ, Gregorie CJ, Hu Z, Magner WJ, Tomasi TB: **MicroRNA targets in immune**
663 **genes and the Dicer/Argonaute and ARE machinery components**. *Molecular*
664 *immunology* 2008, **45**:1995-2006.
- 665 36. Gallagher TF, Hellman L, Zumoff B, Miller DG: **Steroid Hormone Metabolism in Chronic**
666 **Myelogenous Leukemia**. *Blood* 1965, **25**:743-748.
- 667 37. Sun SH, Ou HC, Jang TH, Lin LB, Huang HM: **Altered phospholipid metabolism in**
668 **sodium butyrate-induced differentiation of C6 glioma cells**. *Lipids* 1997, **32**:273-282.
- 669 38. Alisky JM, van de Wetering CI, Davidson BL: **Widespread dispersal of cholera toxin**
670 **subunit b to brain and spinal cord neurons following systemic delivery**.
671 *Experimental neurology* 2002, **178**:139-146.
- 672 39. Shukla SK, Gebregiworgis T, Purohit V, Chaika NV, Gunda V, Radhakrishnan P, Mehla K,
673 Pipinos, II, Powers R, Yu F, Singh PK: **Metabolic reprogramming induced by ketone**
674 **bodies diminishes pancreatic cancer cachexia**. *Cancer & metabolism* 2014, **2**:18.
- 675 40. Buchholz M, Schatz A, Wagner M, Michl P, Linhart T, Adler G, Gress TM, Ellenrieder V:
676 **Overexpression of c-myc in pancreatic cancer caused by ectopic activation of**
677 **NFATc1 and the Ca²⁺/calcineurin signaling pathway**. *The EMBO journal* 2006,
678 **25**:3714-3724.

- 679 41. Soucek L, Whitfield J, Martins CP, Finch AJ, Murphy DJ, Sodir NM, Karnezis AN, Swigart LB,
680 Nasi S, Evan GI: **Modelling Myc inhibition as a cancer therapy.** *Nature* 2008, **455**:679-
681 683.
- 682 42. Kunnumakkara AB, Guha S, Krishnan S, Diagaradjane P, Gelovani J, Aggarwal BB: **Curcumin
683 potentiates antitumor activity of gemcitabine in an orthotopic model of
684 pancreatic cancer through suppression of proliferation, angiogenesis, and
685 inhibition of nuclear factor-kappaB-regulated gene products.** *Cancer research* 2007,
686 **67**:3853-3861.
- 687 43. Park S, Lee C, Sabharwal P, Zhang M, Meyers CL, Sockanathan S: **GDE2 promotes
688 neurogenesis by glycosylphosphatidylinositol-anchor cleavage of RECK.** *Science*
689 2013, **339**:324-328.
- 690 44. Yaffe E, Tanay A: **Probabilistic modeling of Hi-C contact maps eliminates systematic
691 biases to characterize global chromosomal architecture.** *Nature genetics* 2011,
692 **43**:1059-1065.
- 693 45. Vinogradov AE: **DNA helix: the importance of being GC-rich.** *Nucleic acids research*
694 2003, **31**:1838-1844.
- 695 46. Jackson DA, Hassan AB, Errington RJ, Cook PR: **Visualization of focal sites of
696 transcription within human nuclei.** *The EMBO journal* 1993, **12**:1059-1065.
- 697 47. Rieder D, Trajanoski Z, McNally JG: **Transcription factories.** *Frontiers in genetics* 2012,
698 **3**:221.
- 699 48. Lieberman-Aiden E, van Berkum NL, Williams L, Imakaev M, Ragozy T, Telling A, Amit I,
700 Lajoie BR, Sabo PJ, Dorschner MO, et al: **Comprehensive mapping of long-range
701 interactions reveals folding principles of the human genome.** *Science* 2009,
702 **326**:289-293.
- 703 49. Maier T, Guell M, Serrano L: **Correlation of mRNA and protein in complex biological
704 samples.** *FEBS letters* 2009, **583**:3966-3973.
- 705 50. Razin SV, Gavrilov AA, Yarovaya OV: **Transcription factories and spatial organization
706 of eukaryotic genomes.** *Biochemistry Biokhimiia* 2010, **75**:1307-1315.
- 707 51. Gushchanskaya ES, Artemov, A.V., Ulyanov, S.V., Penin, A.A., Logacheva, M.D., Razin, S.V.,
708 Gavrilov, A.A.: **Spatial organization of housekeeping genes in interphase nuclei.**
709 *Molecular Biology* 2014, **48**:886-895.
- 710 52. Jin F, Li Y, Dixon JR, Selvaraj S, Ye Z, Lee AY, Yen CA, Schmitt AD, Espinoza CA, Ren B: **A
711 high-resolution map of the three-dimensional chromatin interactome in human
712 cells.** *Nature* 2013, **503**:290-294.
- 713 53. Li H, Durbin R: **Fast and accurate short read alignment with Burrows-Wheeler
714 transform.** *Bioinformatics* 2009, **25**:1754-1760.
- 715 54. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R:
716 **The Sequence Alignment/Map format and SAMtools.** *Bioinformatics* 2009, **25**:2078-
717 2079.
- 718 55. Pervouchine DD, Djebali S, Breschi A, Davis CA, Barja PP, Dobin A, Tanzer A, Lagarde J,
719 Zaleski C, See LH, et al: **Enhanced transcriptome maps from multiple mouse tissues
720 reveal evolutionary constraint in gene expression.** *Nature communications* 2015,
721 **6**:5903.
- 722 56. Djebali S, Davis CA, Merkel A, Dobin A, Lassmann T, Mortazavi A, Tanzer A, Lagarde J, Lin W,
723 Schlesinger F, et al: **Landscape of transcription in human cells.** *Nature* 2012, **489**:101-
724 108.
- 725 57. Daemen A, Griffith OL, Heiser LM, Wang NJ, Enache OM, Sanborn Z, Pepin F, Durinck S,
726 Korkola JE, Griffith M, et al: **Modeling precision treatment of breast cancer.** *Genome
727 biology* 2013, **14**:R110.
- 728 58. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL, Wold BJ,
729 Pachter L: **Transcript assembly and quantification by RNA-Seq reveals
730 unannotated transcripts and isoform switching during cell differentiation.** *Nature
731 biotechnology* 2010, **28**:511-515.

732

733

Figure 1.

Tissue ID	Tissue Source		DNA	RNA
HEK293	Kidney Cell Line	(Replicate 1 & 2)	■	■
hESC	Embryonic Stem Cell Line	(Replicate 1 & 2)	■	
IMR90	Lung Fibroblast Cell Line	(Replicate 1 & 2)	■	■
BT483	Mammary Gland Cell Line	(Replicate 1 & 2)	■	
GM06990	B-Lymphocyte Cell Line	(Replicate 1 & 2)	■	
RWPE1	Prostate Epithelial Cell Line	(Replicate 1 & 2)	■	■

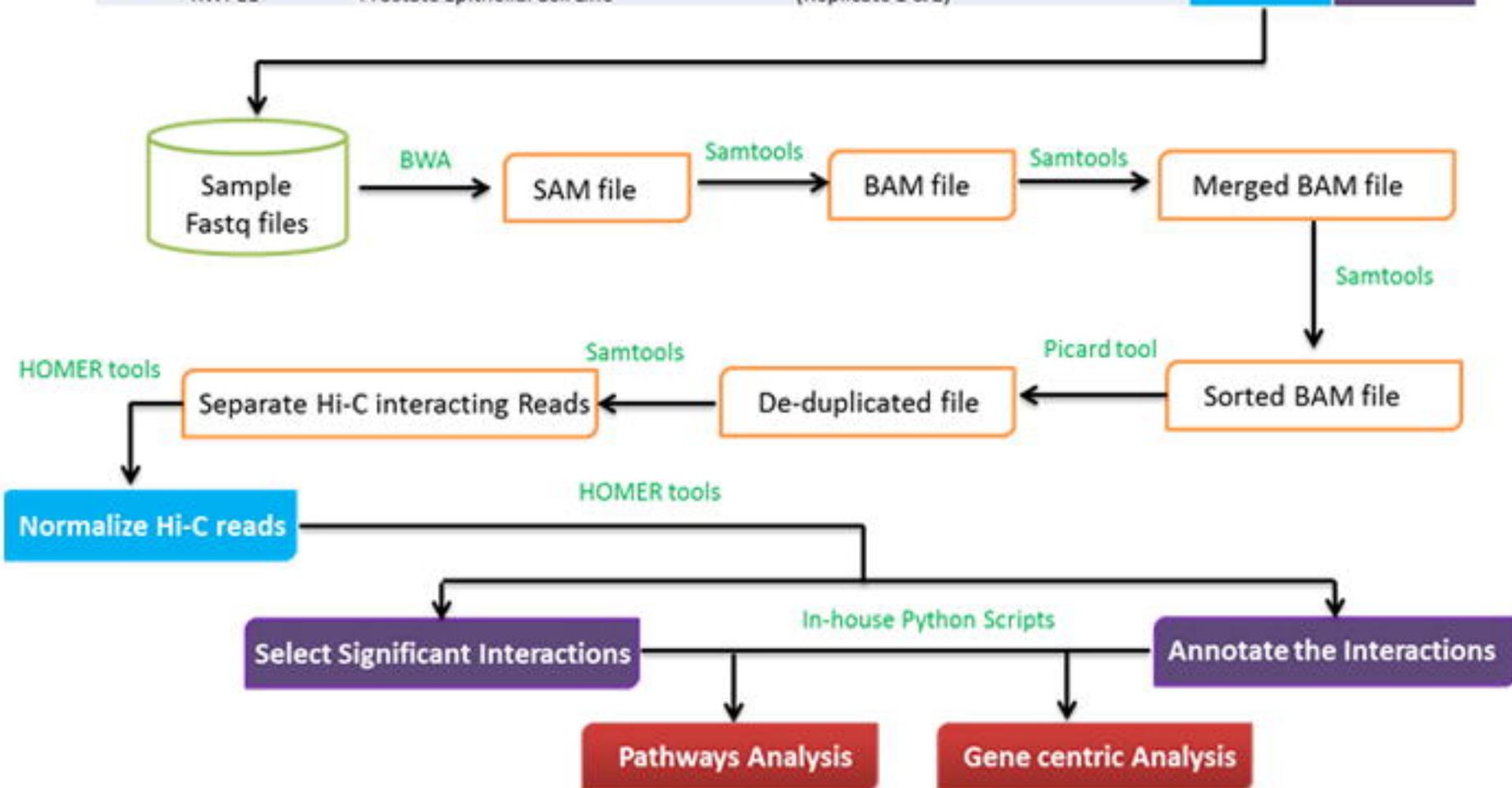
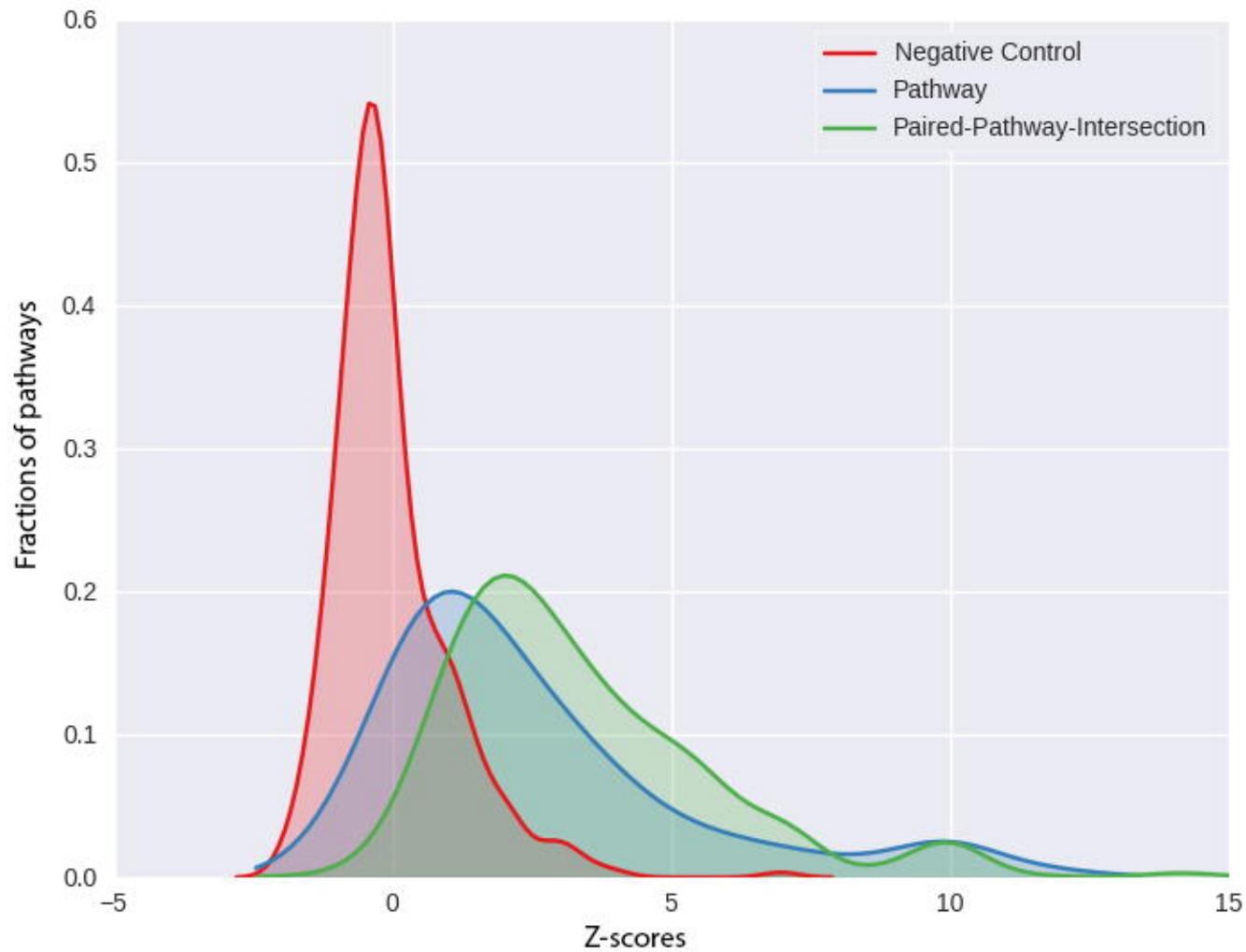
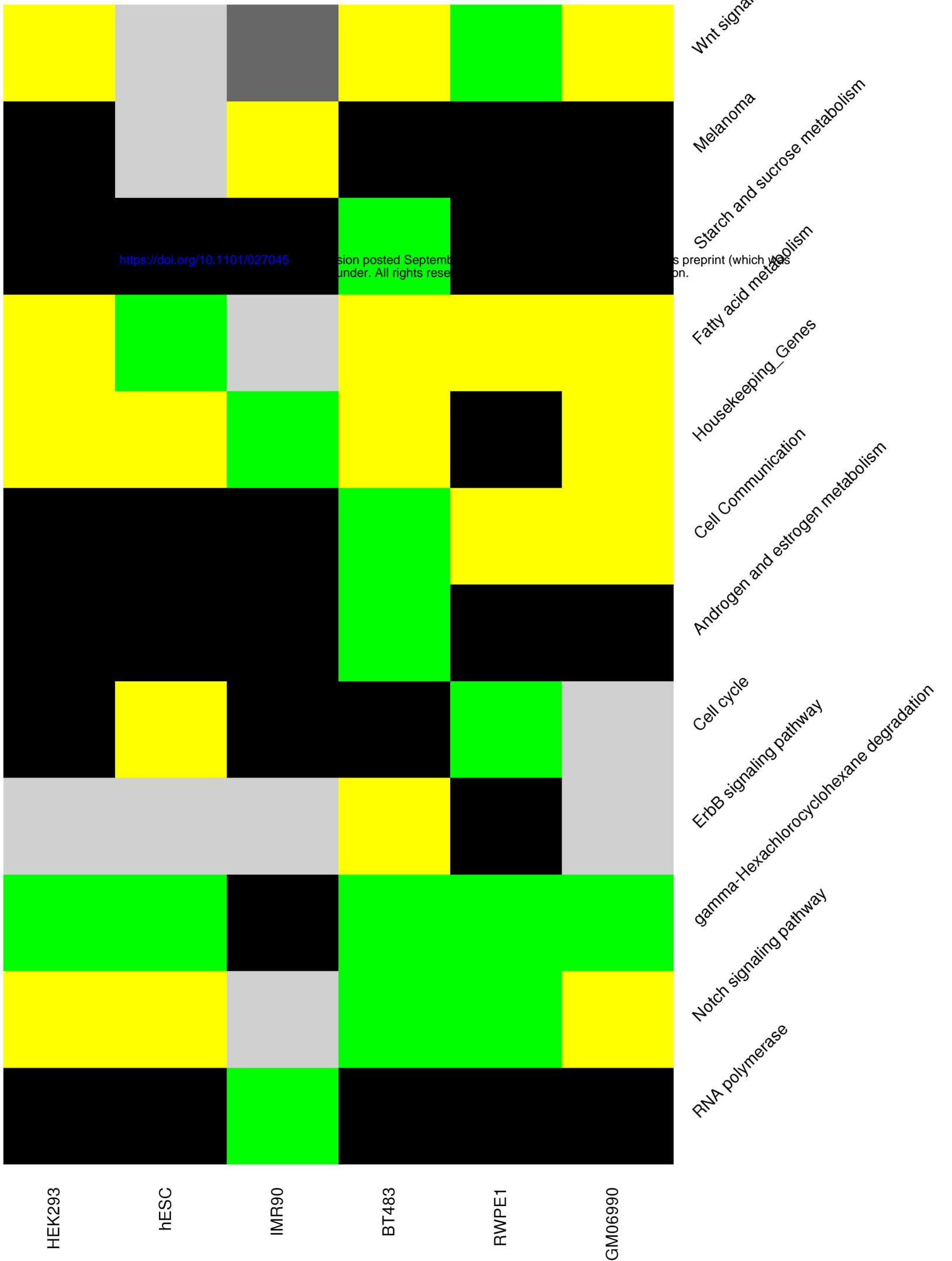


Figure 2.



KEGG PATHWAYS Z-SCORE HEATMAP



<https://doi.org/10.1101/027045> bioRxiv preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

Color Key

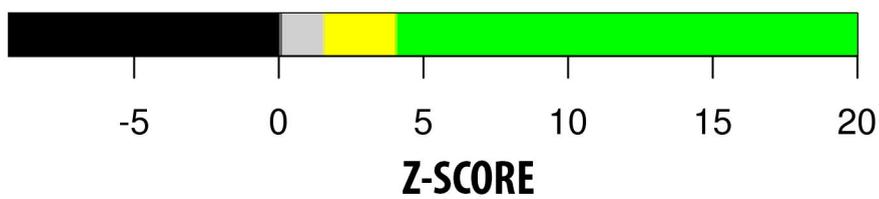
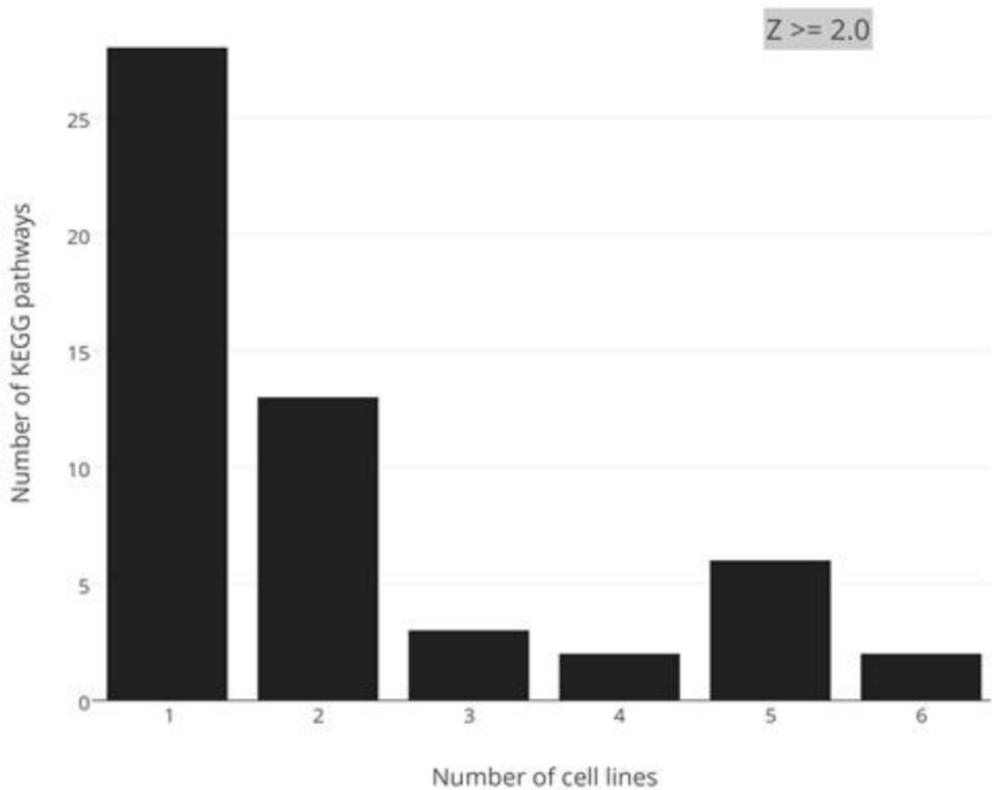


Figure 4.

	HEK293	IMR90	hESC	RWPE1	GM06990	BT483
HEK293	2	8	4	11	10	9
IMR90		2	4	13	10	8
hES			6	5	3	3
RWPE1				15	12	13
GM06990					2	11
BT483						3
Z-score ≥ 2.0						

Figure 5.

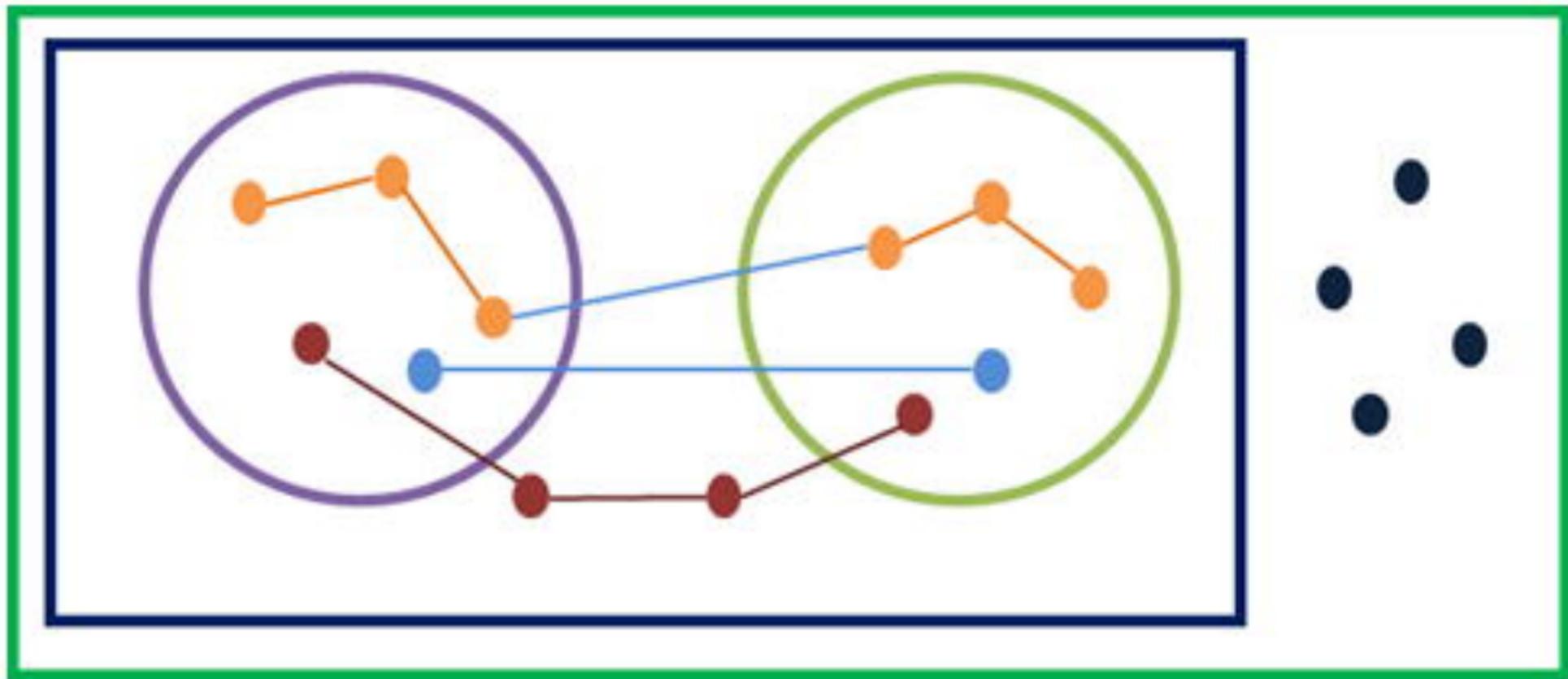
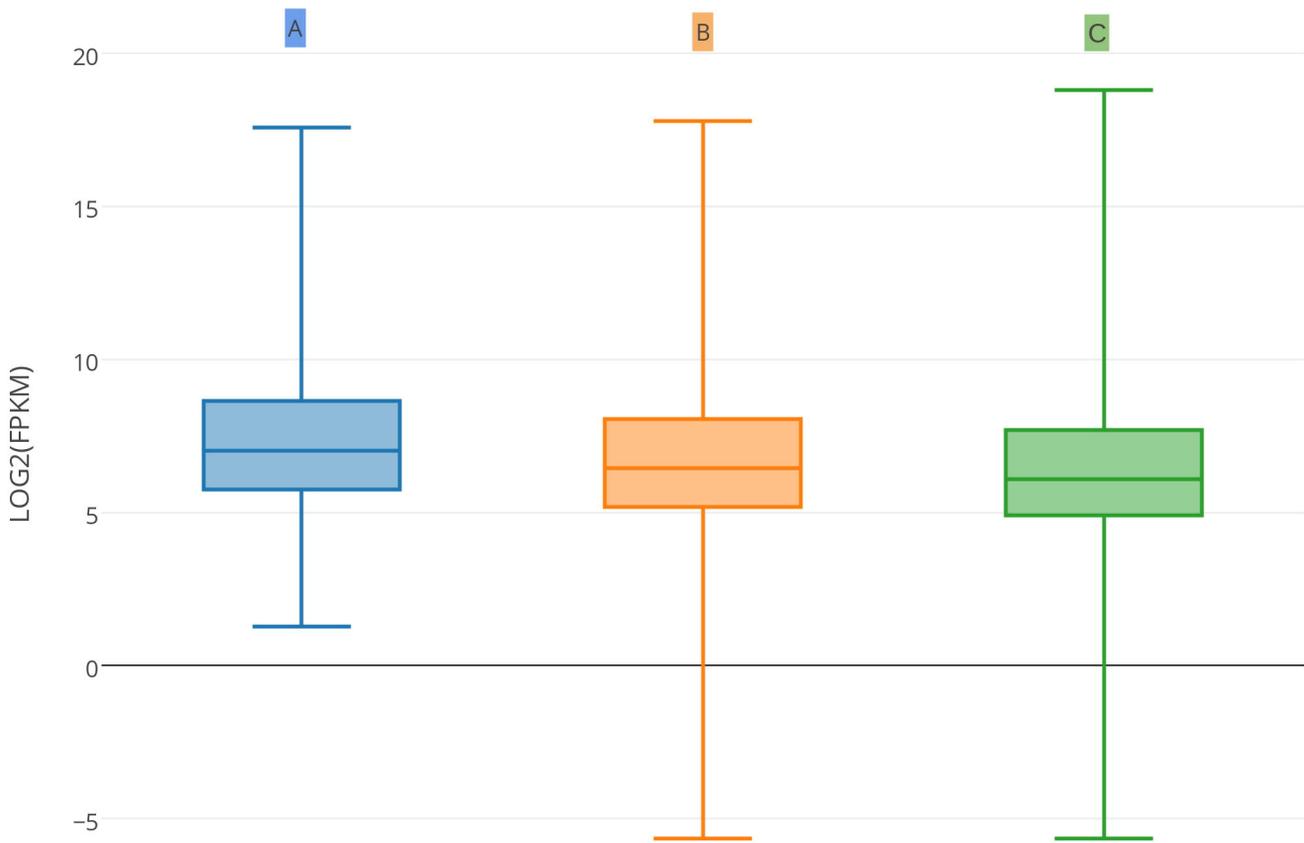


Figure 6.

POOLED 6 CELL LINES



- proximal intra-pathway
- proximal-inter-pathway AND proximal-generic
- non-proximal-generic

Figure 7.

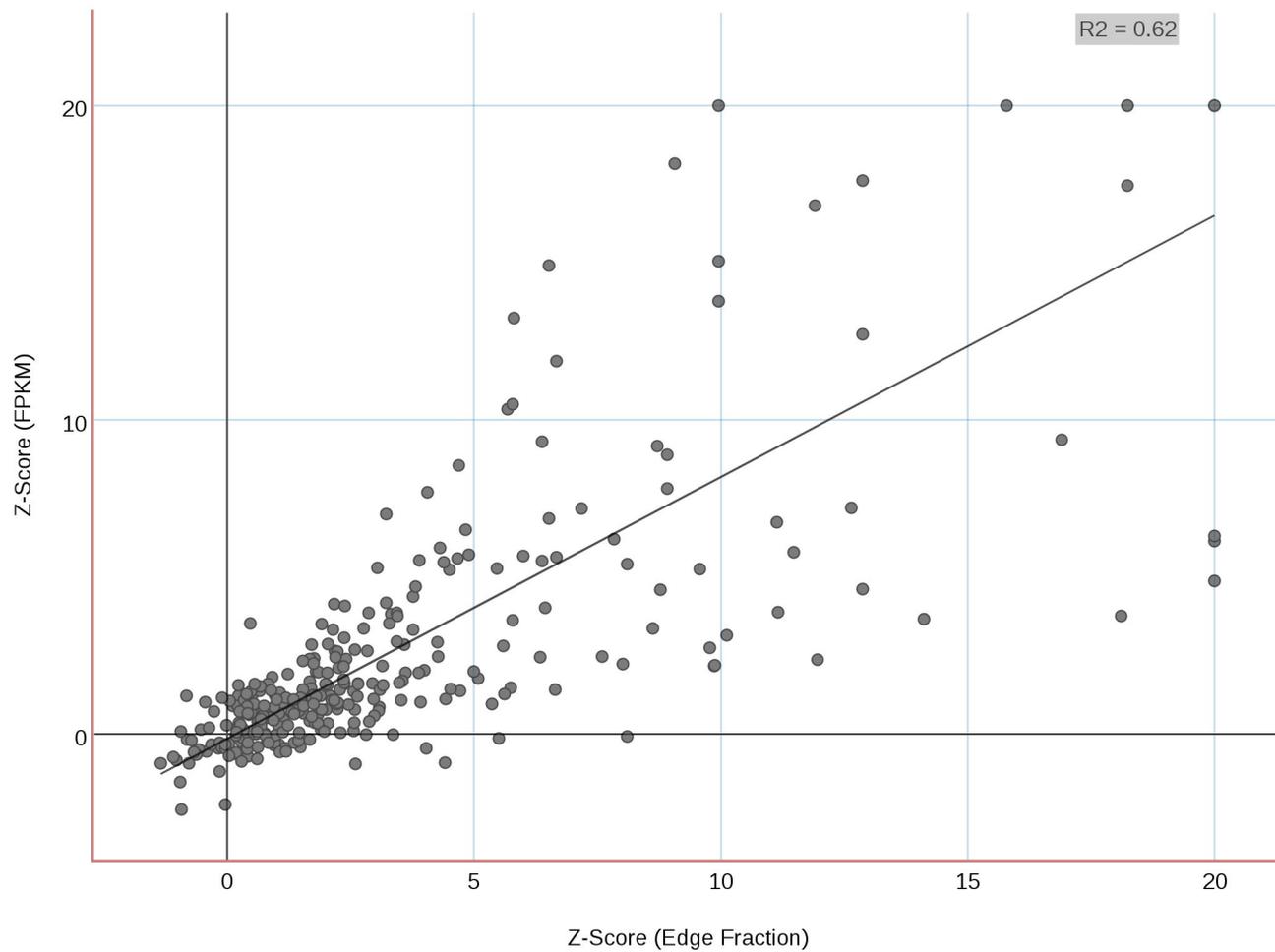


Figure 8.

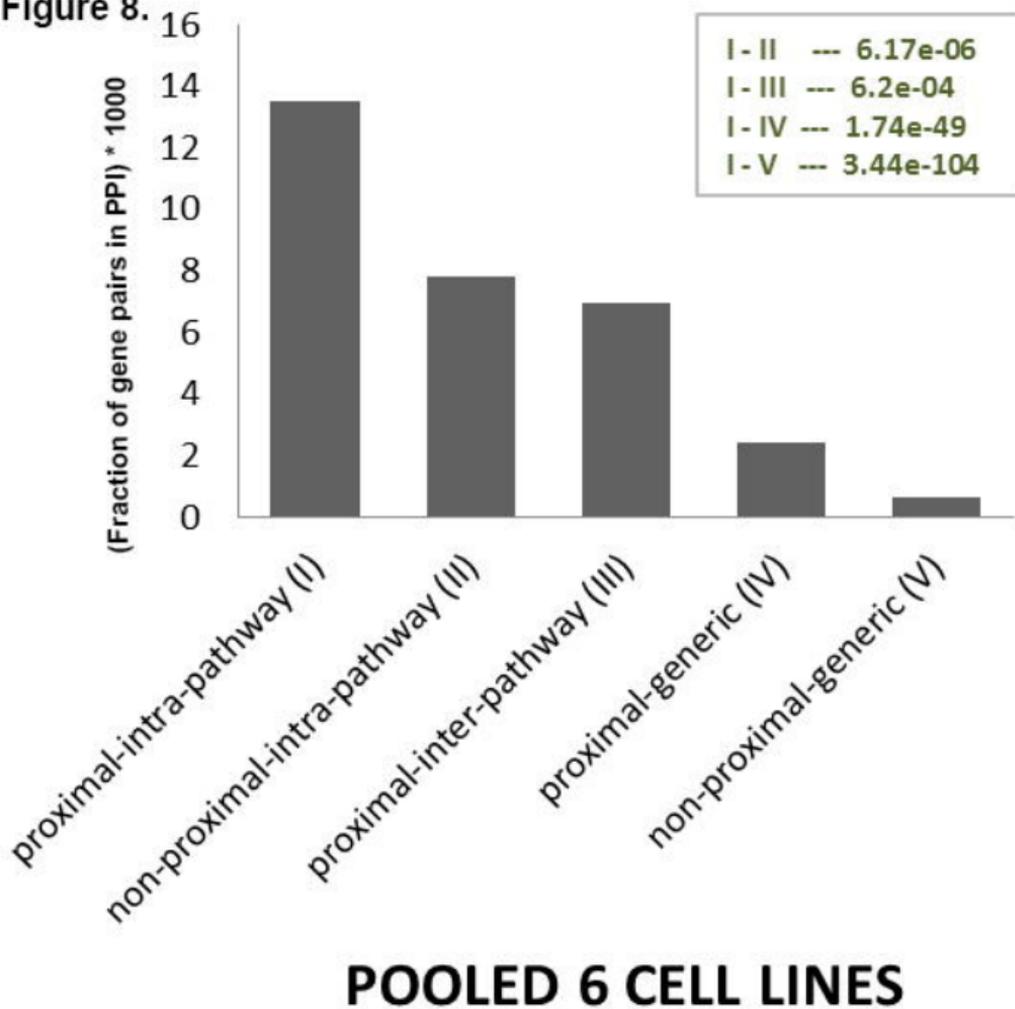


Figure 9

GO - Molecular Functions

