

Enhancer sharing promotes neighborhoods of transcriptional regulation across eukaryotes

Porfirio Quintero-Cadena¹ and Paul W. Sternberg^{1*}

¹*Division of Biology and Biological Engineering, California Institute of Technology; Howard Hughes Medical Institute, Pasadena, California, USA. *Corresponding author: pws@caltech.edu.*

1 **Enhancers physically interact with transcriptional promoters, looping over distances that**
2 **can span multiple regulatory elements. Given that enhancer-promoter (EP) interactions**
3 **generally occur via common protein complexes, it is unclear whether EP pairing is predom-**
4 **inantly deterministic or proximity guided. Here we present cross-organismic evidence sug-**
5 **gesting that most EP pairs are compatible, largely determined by physical proximity rather**
6 **than specific interactions. By re-analyzing transcriptome datasets, we find that the tran-**
7 **scription of gene neighbors is correlated over distances that scale with genome size. We**
8 **experimentally show that non-specific EP interactions can explain such correlation, and that**
9 **EP distance acts as a scaling factor for the transcriptional influence of an enhancer. We pro-**
10 **pose that enhancer sharing is commonplace among eukaryotes, and that EP distance is an**
11 **important layer of information in gene regulation.**

12 **Introduction**

13 Enhancers mediate the transcriptional regulation of gene expression, enabling isogenic cells to ex-
14 hibit remarkable phenotypic diversity (Davidson and Peter, 2015). In complex with transcription
15 factors, they interact with promoters via chromatin looping (Marsman and Horsfield, 2012), finely
16 regulating transcription in time and space. A prevailing view is that most enhancers have a mech-
17 anism to selectively loop to a target promoter (van Arensbergen et al., 2014). Examples in this
18 category usually require specific transcription factor binding at both enhancer and promoter sites
19 (Davidson and Peter, 2015), which could explain why some enhancers seem to influence different
20 promoters in varying degrees (Gehrig et al., 2009). On the other hand, EP looping is generally
21 mediated by common protein complexes (Kagey et al., 2010; Malik and Roeder, 2010), conflicting
22 with the specific molecular interactions required by such a model at a larger scale. Examples of
23 non-specific EP pairing seem to be common (Butler and Kadonaga, 2001), yet this model could
24 result in transcriptional crosstalk, which appears inconsistent with our current paradigm of gene
25 regulation. The predominant EP pairing scheme –specific or non-specific– and its determinants
26 are thus unclear. Here we ask to what extent are potential EP pairs compatible through a meta-
27 analysis of the genome-wide transcription of gene neighbors in five species. We propose that en-
28 hancer sharing occurs widely across eukaryotes, test key aspects of this hypothesis in *C. elegans*,
29 and analyze its implications in other genomic phenomena.

30 **Results and Discussion**

31 **Gene neighbors are transcriptionally correlated genome-wide**

32 To investigate the transcriptional relationship between gene neighbors, we paired every protein-
33 coding gene of five organisms (*Saccharomyces cerevisiae*, *Caenorhabditis elegans*, *Drosophila*
34 *melanogaster*, *Mus musculus* and *Homo sapiens*) with its 100 nearest neighbors within the same
35 chromosome, which yielded lists of around 600,000 (*S. cerevisiae*) and 2 million (each of the
36 rest) gene pairs. We then computed the Spearman correlation coefficient between paired genes
37 across multiple RNA-seq datasets (Attrill et al., 2016; Ellahi et al., 2015; Gerstein et al., 2010;
38 The ENCODE Project Consortium, 2012) and the intergenic distance between the the start of the
39 5' untranslated region of the first gene to the start of the second gene in each pair. We observed
40 that neighboring genes tend to be correlated in transcript abundance genome-wide in all analyzed
41 organisms, and that this correlation decays exponentially with increasing intergenic distance (Fig-
42 ure 1a). We thus fitted the data to an exponential decay function to compute the mean distance at
43 which a pair of genes remain correlated (d_{exp}). Consistent with the persistence of the correlation
44 pattern across organisms, d_{exp} scaled with genome size, to 1 kilobase in *S. cerevisiae*, ~10 kb in
45 *C. elegans* and *D. melanogaster*, and ~350 kb in *M. musculus* and in *H. sapiens* (Supplementary
46 Figure 1). This trend remained largely the same even after removing duplicated genes pairs (Sup-
47plementary Figure 2). Most genes had at least one neighbor closer than d_{exp} in all species (Figure
48 1b), and the representation of gene ontology annotations remained unbiased in correlated gene
49 pairs (Supplementary Figure 3), indicating that the average gene is correlated in expression with

50 its nearest neighbors beyond any particular gene class. In addition, sampled intergenic distances
51 go well beyond d_{exp} (Figure 1c), indicating that 100 gene neighbors is a sufficient number to study
52 this effect.

53 To examine the correlation of gene expression in the spatial domain, we analyzed RNA *in situ*
54 hybridization data for 6053 genes in *D. melanogaster* (Hammonds et al., 2013; Tomancak et al.,
55 2002, 2007). We computed the percentage overlap in tissue expression by dividing the number
56 of common tissues over the total number of unique tissues in which genes of any given pair are
57 expressed (Supplementary Figure 4a). This analysis revealed that close neighbors have a tendency
58 to be expressed in the same tissues, and that this overlap also decays exponentially with intergenic
59 distance (Supplementary Figure 4b). However, the correlation extends to a longer mean distance
60 ($d_{exp} = 22$ compared to 6 kb), suggesting that RNA-seq analysis, which included mostly whole-
61 organism transcriptome averages, resulted in a conservative estimate. Gene neighbors thus have a
62 spatio-temporal correlation in expression that is highly dependent upon the spacing between them.
63 Our meta-analysis extends the findings of previous reports (Michalak, 2008) genome-wide and
64 across organisms. In particular, pairing every gene with 100 proximal genes provides a complete
65 set of distance-dependent correlations between gene pairs.

66 **Enhancer sharing explains the transcriptional correlation of gene neighbors**

67 The pervasive nature of proximal gene co-expression suggested a common underlying mechanism.
68 We hypothesized enhancer sharing among nearby genes to be a plausible explanation, as it is in

69 agreement with several observations: i) enhancers regulate transcription by making contact with
70 promoters via chromatin looping (He et al., 2014), whose incidence also decays exponentially as
71 the distance between contacting sites increases (Rao et al., 2014; Ringrose et al., 1999), with the
72 same pattern as observed here (Supplementary Figure 5) ii) the average distance between studied
73 EP interactions scales with genome size in ranges consistent with d_{exp} for each analyzed organism
74 (Araya et al., 2014; He et al., 2014) iii) broad enhancer compatibility and sharing is consistent
75 with the idea that common protein complexes such as the mediator are widely utilized bridges in
76 EP looping (Kagey et al., 2010; Malik and Roeder, 2010) and iv) a high frequency of chromatin in-
77 teractions are observed within topologically associated domains identified through high-resolution
78 Chromosome Conformation Capture (Hi-C) (Rao et al., 2014). Consistent with this view, it seems
79 that enhancers often do not show promoter specificity (Butler and Kadonaga, 2001). This line of
80 reasoning suggests a model where, as opposed to only having a specific target gene (Figure 2a),
81 enhancers have a range of action in which they can influence any gene (Figure 2b), at least to
82 an extent that causes the correlation of nearby genes, likely within a topological domain. Tran-
83 scriptome analysis could thus provide indirect evidence of genome and condition-wide EP looping
84 that is difficult to access through Hi-C (Rao et al., 2014) due to the low signal-to-noise ratio of
85 short-range interactions.

86 Because of its compact genome, rapid development and availability of tissue specific en-
87 hancers (Corsi et al., 2015), we decided to use *C. elegans* to test this hypothesis. We first postu-
88 lated that unrelated enhancers should generally be compatible, showing qualitative additivity when
89 placed upstream of a single promoter. We thus paired the well characterized *myo-2* pharyngeal

90 enhancer with a body wall muscle (BWM) and a BWM plus intestine enhancer, placed them up-
91 stream a minimal promoter and a *gfp* reporter, and examined expression in transgenic animals. In
92 both cases, we observed fluorescence in the corresponding tissues (Figure 2c, d, e). This obser-
93 vation is consistent with typical enhancer studies in artificial constructs (Dupuy et al., 2004) and
94 agrees with some EP compatibility studies (Butler and Kadonaga, 2001).

95 Given that both chromatin looping and expression correlation decay exponentially, we rea-
96 soned that transcription of a given gene should also decay exponentially with increasing EP dis-
97 tance if the observed correlation is to be explained by enhancer sharing. To test this hypothesis,
98 we first built a series of genetic constructs with increasing neutral EP distances (0, 1, 1.5 and 2 kb)
99 for two different enhancers, *myo-2* and *unc-54*. We then integrated each construct in single copy
100 into the genome of *C. elegans* and used quantitative PCR to i) measure the influence of EP distance
101 on the reporter gene in native chromatin and ii) analyze the impact of the perturbation on the two
102 genes that natively flank the site of transgene insertion (*dpy-13* and *col-34*, Figure 2h), which we
103 reasoned should be affected in two counteracting ways. First, the ectopic enhancers should pro-
104 mote their expression. Second, the increased EP distance caused by the addition of spacers should
105 reduce their expression by scaling down the influence of both ectopic and native enhancers (the
106 latter of unknown identity and location) to the opposite side of the spacer.

107 We found that transcriptional levels of the reporter gene indeed fall rapidly with increasing
108 EP distance with both enhancers (Figure 2f, g); this occurred at a rate that seems congruent or
109 faster than the genome-wide correlation decay, likely because of the dramatic separation of every

110 regulatory element at once, as opposed to gradual separation from individual enhancers in a native
111 environment. Transcription was still well detected even when the enhancers were placed 2 kb away,
112 supporting the hypothesis that EP distance is a scaling factor on the enhancer's influence. Expres-
113 sion of *dpy-13* and *col-34* was reduced with the introduction of the 2 kb spacer when compared
114 to transgenic lines without it (Figure 2i, j). On the other hand, spacer-free lines were comparable
115 to wild-type, suggesting the incorporation of ectopic enhancers compensated for the EP distance
116 increase caused by the addition of the genetic construct itself. Hence, these observations fit the
117 corollaries of our model even amid the complexity of a native regulatory environment.

118 **Enhancer-promoter distance insulates neighboring genes**

119 We next wished to determine the extent to which enhancer sharing explains other genomic phe-
120 nomena. Previous reports have suggested that divergent, parallel and convergent gene pairs tend
121 to have distinct correlation profiles (e.g. (Chen and Stein, 2006)). To explore this hypothesis, we
122 first compared the distribution of intergenic distances of gene pairs oriented in parallel, divergent
123 and convergent orientations (Figure 3a, Supplementary Figure 6). As expected, divergent gene
124 pairs tend to be closest, followed by parallel and finally convergent genes. We then confirmed
125 that each group appears to have different distributions of correlations (Figure 3b, Supplementary
126 Figure 6). To consider the influence of EP distance, we sampled gene pairs from each orientation
127 controlling for intergenic size. This resulted in distributions of correlations that exactly overlap
128 (Figure 3c, Supplementary Figure 6), an observation that is supported by previous reports (Ghan-
129 barian and Hurst, 2015). We thus conclude that the apparent influence of gene orientation in the

130 transcriptional relationship of neighboring gene pairs is a consequence of enhancer sharing and EP
131 distance.

132 From the regulatory perspective, EP distance provides independence to most gene pairs, as
133 the vast majority have an intergenic distance that puts them in the baseline correlation regime (Fig-
134 ure 1c). To study the enhancer-blocking influence of insulators (Bushey et al., 2009) genome-wide,
135 we analyzed each group of genes flanked by insulator binding sites, which were previously deter-
136 mined using Chromatin ImmunoPrecipitation coupled with microarrays (ChIP-chip) for the six
137 known insulators in *D. melanogaster*: BEAF-32, CP190, CTCF, GAF, Mod(mdg4) and Su(Hw)
138 (Negre et al., 2010). We observed that gene neighbors closer than 10 kb bound by each of the insu-
139 lators tend to be less correlated in gene expression than gene pairs not bound by them (Figure 3e),
140 supporting their role in genome-wide insulation and agreeing with the observation that insulators
141 tend to separate differentially expressed genes (Negre et al., 2010; Xie et al., 2007). Nevertheless,
142 the same groups of gene pairs also tend to have much larger intergenic distances than genes that
143 are not flanked by insulator binding sites (Figure 3d). After controlling for the distribution of inter-
144 genic distances, we found very similar correlation distributions between insulator and not insulator
145 flanked gene pairs (Figure 3f). This agrees with previous reports suggesting that insulators do not
146 block enhancers everywhere they bind, but rather act only on very specific genomic contexts (Liu
147 et al., 2015; Schwartz et al., 2012); it also reconciles the lack of known insulator orthologs in *C. el-*
148 *egans* (Heger et al., 2009) in the context of local enhancer-blocking. In combination, these studies
149 strongly suggest that EP distance is the general source of transcriptional independence for close
150 gene neighbors.

151 Previous EP compatibility studies suggest that EP specificity is widespread (Gehrig et al.,
152 2009), while others that it is restricted to a smaller subset of enhancers (Butler and Kadonaga,
153 2001). Although our results support the latter, views arising from these studies might not be mu-
154 tually exclusive, as it is likely that enhancers have specificity to promoter classes (Danino et al.,
155 2015), whose limited number could result in general EP compatibility. Nevertheless, the implica-
156 tions from considering our observations are broadly applicable to gene regulation and likely act
157 in conjunction with other regulatory mechanisms, such as chromatin accessibility. For example,
158 position effects, in which transgene expression levels are influenced by the insertion site (Gierman
159 et al., 2007), are naturally expected from enhancer sharing. Chromosomal translocations and mu-
160 tations involving regulatory elements likely impact genetic contexts rather than individual genes.
161 Furthermore, enhancer scaling by EP distance contributes to the weight of any given EP interaction
162 and thus enriches both the complexity and tunability of genomic logic. Our analysis provides a
163 clarifying perspective of gene regulation consistent with both mechanistic and genome-wide stud-
164 ies.

165 **Experimental procedures**

166 **Computational biology**

167 For each analyzed organism, Ensembl (Flicek et al., 2014) protein-coding genes were grouped
168 by chromosome, sorted by position, and paired with the 100 nearest neighbors within the same
169 chromosome. A list of duplicated gene pairs for *H. sapiens* and *M. musculus* was obtained from
170 the Duplicated Genes Database (Ouedraogo et al., 2012) (<http://dgd.genouest.org>). A list of *C.*
171 *elegans* genes predicted to be in operons was obtained from Allen et al. (2011), and gene pairs
172 present in the same operon were removed from the analysis to prevent co-transcriptional bias.
173 Processed RNA-seq data was obtained from multiple sources (Attrill et al., 2016; Ellahi et al.,
174 2015; Gerstein et al., 2010; The ENCODE Project Consortium, 2012) and converted to transcripts
175 per million (TPM) (Wagner et al., 2012) when necessary. Formatted datasets are available upon
176 request. Genes detected in less than 80% of experiments were discarded. To compute the Spearman
177 correlation coefficient, TPM values were ranked in each RNA-seq experiment and the pairwise
178 Pearson correlation coefficient was computed on the ranked values according to the following
179 equation:

$$\rho = \frac{cov(\text{gene}_1, \text{gene}_2)}{\sigma_{\text{gene}_1} \sigma_{\text{gene}_2}}$$

180 where gene_1 and gene_2 are the corresponding ranks of each paired gene in a given RNA-seq experi-
181 ment, cov their covariance and σ their standard deviation. The list of gene pairs with intergenic dis-
182 tances and correlation coefficients was sorted by increasing intergenic distance, and subsequently
183 smoothed using a sliding median with window size of 1000 gene pairs. The result was then fitted

184 to an exponential decay function of the form:

$$\rho(d) = \rho_0 e^{-\lambda d} + c$$

185 where ρ_0 is the median Spearman correlation coefficient of the closest neighboring genes, d the
186 intergenic distance and c the baseline correlation. The mean distance at which a pair of genes
187 remain correlated was then computed as:

$$d_{exp} = 1/\lambda$$

188 A list of genes annotated with RNA *in situ* hybridization data (Hammonds et al., 2013; Tomancak
189 et al., 2002, 2007) was obtained from the Berkeley Drosophila Genome Project ([http://insitu.fruit-](http://insitu.fruit-fly.org)
190 [fly.org](http://insitu.fruit-fly.org)). Insulator ChIP-chip data was obtained from Negre et al. (2010) (GSE16245); the intersec-
191 tion of replicates was used. HiC data was obtained from Rao et al. (2014) (GSE63525, GM12878
192 primary replicate HiCCUPS looplevelist). Functional protein classification was conducted using Pan-
193 ther (Mi et al., 2016). Genomic manipulations were conducted using Bedtools v2.24.0 (Quinlan
194 and Hall, 2010). Data analysis was conducted using Python 2.7.9 and the Scipy library (McKinney,
195 2010).

196 **Molecular biology**

197 *C. elegans* was cultured under standard laboratory conditions (Stiernagle, 2006). For enhancer
198 additivity experiments, transgenic *C. elegans* lines carrying extra-chromosomal arrays were gen-
199 erated by injecting each plasmid at 50 ng/ μ l into *unc-119* mutant animals. The minimal $\Delta pes-10$
200 promoter (Fire et al., 1990) was used in all constructs. Minimal regions of the *myo-2* and *unc-54*

201 enhancers (Okkema et al., 1993) able to drive tissue specific expression were used. The BWM
202 enhancer was obtained from the upstream region of *F44B9.2*; the BWM plus intestine enhancer
203 was obtained from the upstream region of *rps-1*. Animals were imaged at 40x using a GFP filter
204 on a Zeiss Axioskop microscope. For the enhancer promoter distance and ectopic enhancer exper-
205 iments, we defined an EP distance of 0 to be the enhancer placed just upstream of the $\Delta pes-10$
206 promoter, which is ~ 350 bp. To ensure neutrality yet maintain a similar GC content as non-coding
207 sequences in *C. elegans*, we used non-overlapping AT-rich DNA spacers obtained from the genome
208 of *Escherichia coli*. Constructs were integrated in single-copy into chromosome IV via CRISPR-
209 Cas9 using plasmids provided as gifts by Dr. Zhiping Wang and Dr. Yishi Jin (unpublished re-
210 sults). Briefly, plasmids containing the following expression cassettes were co-injected: reporter
211 and hygromycin resistance genes flanked by homologous arms for recombination-directed repair
212 (10 ng/ μ l), single-guide RNA (10 ng/ μ l), Cas9 (10 ng/ μ l), and extra-chromosomal array reporter
213 for expression of either *rfp* or *gfp* outside the tissue of interest (10 ng/ μ l). Transformants were
214 selected for using hygromycin at 10 μ g/ μ l, and those not carrying extra-chromosomal transgenes,
215 lacking of *gfp* or *rfp* expression outside the tissue of interest, were subsequently isolated. Animals
216 homozygotic for the insertion were identified by polymerase-chain reaction (PCR) and Sanger se-
217 quencing. Quantitative PCR was carried out as previously described (Ly et al., 2015) using *pmp-3*
218 as a reference gene (Zhang et al., 2012). Briefly, third-stage larval (L3) worms, when expression
219 from the test enhancers is maximal according to RNA-seq data, were synchronized at 20°C via
220 egg-laying. 15 animals were lysed in 1.5 μ l of Lysis Buffer (5 mM Tris pH 8.0 (MP Biomed-
221 icals), 0.5% Triton X-100, 0.5% Tween 20, 0.25 mM EDTA (Sigma-Aldrich)) with proteinase-K

222 (Roche) at 1.5 µg/µl, incubated at 65°C for 10 minutes followed by 85°C for one minute. Reverse
223 transcription was carried out using the Maxima H Minus cDNA synthesis kit (Thermo Fisher) by
224 adding 0.3 µl H₂O, 0.6 µl 5x enzyme buffer, 0.15 µl 10mM dNTP mix, 0.15 µl 0.5 µg/µl oligo dT
225 primer, 0.15 µl enzyme mix and 0.15 µl DNase, and incubated for 2 minutes at 37°C, followed
226 by 30 minutes at 50°C and finally 2 minutes at 85°C. The cDNA solution was diluted to 15 µl
227 and 1 µl was used for each qPCR reaction, so that on average each well contained the amount of
228 RNA from a single worm. All qPCR reactions were performed with three technical replicates and
229 at least three biological replicates using the Roche LightCycler[®] 480 SYBR Green I Master in
230 the LightCycler[®] 480 System. Crossing point-PCR-cycle (Cp) averages were computed for each
231 group of three technical replicates; these values were then subtracted from the respective average
232 Cp value of the reference gene.

233 **Data and reagent availability**

234 Strains are available upon request. Relevant DNA sequences, including spacers, enhancers, primers,
235 sgRNA and homology arms are available in the supplementary table 1. Correlation datasets are
236 available in the supplementary file 1. Python source code, and links to all expression datasets used
237 in this study, are available for download on the following github repository: [https://github.com/](https://github.com/WormLabCaltech/QuinteroSternberg2016.git)
238 [WormLabCaltech/QuinteroSternberg2016.git](https://github.com/WormLabCaltech/QuinteroSternberg2016.git).

239 **Acknowledgments**

240 Our work was supported by the Howard Hughes Medical Institute, of which P.W.S is an investi-
241 gator. We thank Zhiping Wang and Yishi Jin for plasmids for Crispr-Cas9 single copy insertion,
242 Carmie Robinson for discussions, experimental suggestions and comments on the manuscript, Han
243 Wang for discussions, technical advise and comments on the manuscript, Hillel Schwartz, Mitchell
244 Guttman, Mihoko Kato, David Angeles-Albores, Jonathan Liu, Barbara Wold, Isabelle Peter and
245 Angelike Stathopoulos for discussions and comments on the manuscript, the Encode and ModEn-
246 code consortiums, Flybase, Wormbase and Ensembl databases, the Wold Lab and the Guigo Lab
247 for data accessibility. This paper is dedicated to the memory of Eric H. Davidson.

248 **Contributions**

249 P.Q.C. performed the experiments and analyzed the data. P.Q.C. and P.W.S. designed the experi-
250 ments and wrote the paper.

251 **Competing financial interests**

252 The authors declare no competing financial interests.

253

254 Allen, M. A., Hillier, L. W., Waterston, R. H., and Blumenthal, T. (2011). A global analysis of *C.*
255 *elegans* trans-splicing. *Genome Res.*, 21:255–264.

256 Araya, C. L., Kawli, T., Kundaje, A., Jiang, L., Wu, B., Vafeados, D., Terrell, R., Weissdepp, P.,
257 Gevirtzman, L., Mace, D., Niu, W., Boyle, A. P., Xie, D., Ma, L., Murray, J. I., Reinke, V.,
258 Waterston, R. H., and Snyder, M. (2014). Regulatory analysis of the *C. elegans* genome with
259 spatiotemporal resolution. *Nature*, 512(7515):400–405.

260 Attrill, H., Falls, K., Goodman, J. L., Millburn, G. H., Antonazzo, G., Rey, A. J., Marygold, S. J.,
261 and the FlyBase consortium (2016). Flybase: establishing a gene group resource for drosophila
262 melanogaster. *Nucleic Acids Res.*, 44(D1):D786–D792.

263 Bushey, A. M., Dorman, E. R., and Corces, V. G. (2009). Chromatin insulators:regulatory mecha-
264 nisms and epigenetic inheritance. *Mol. Cell*, 32(404):1–9.

265 Butler, J. E. and Kadonaga, J. T. (2001). Enhancer–promoter specificity mediated by dpe or tata
266 core promoter motifs. *Genes Development*, 15(19):2515–2519.

267 Chen, N. and Stein, L. D. (2006). Conservation and functional significance of gene topology in the
268 genome of *Caenorhabditis elegans*. *Genome Res.*, 16(5):606–617.

269 Corsi, A. K., Wightman, B., and Chalfie, M. (2015). A transparent window into biology: A primer
270 on *Caenorhabditis elegans*. *Genetics*, 200(2):387–407.

271 Danino, Y. M., Even, D., Ideses, D., and Juven-Gershon, T. (2015). The core promoter: At the

- 272 heart of gene expression. *Biochimica et Biophysica Acta (BBA) - Gene Regulatory Mechanisms*,
273 1849(8):1116 – 1131.
- 274 Davidson, E. H. and Peter, I. S. (2015). Chapter 1 - the genome in development. In Davidson,
275 E. H. and Peter, I. S., editors, *Genomic Control Process*, pages 1–40. Academic Press, Oxford.
- 276 Dupuy, D., Li, Q.-R., Deplancke, B., Boxem, M., Hao, T., Lamesch, P., Sequerra, R., Bosak, S.,
277 Doucette-Stamm, L., Hope, I. A., Hill, D. E., Walhout, A. J., and Vidal, M. (2004). A first
278 version of the caenorhabditis elegans promoterome. *Genome Res.*, 14(10b):2169–2175.
- 279 Ellahi, A., Thurtle, D. M., and Rine, J. (2015). The chromatin and transcriptional landscape of
280 native saccharomyces cerevisiae telomeres and subtelomeric domains. *Genetics*, 2:505–521.
- 281 Fire, A., Harrison, S. W., and Dixon, D. (1990). A modular set of lacZ fusion vectors for studying
282 gene expression in caenorhabditis elegans. *Gene*, 93(2):189 – 198.
- 283 Flicek, P., Amode, M. R., Barrell, D., Beal, K., Billis, K., Brent, S., Carvalho-Silva, D., Clapham,
284 P., Coates, G., Fitzgerald, S., Gil, L., Girón, C. G., Gordon, L., Hourlier, T., Hunt, S., Johnson,
285 N., Juettemann, T., Kähäri, A. K., Keenan, S., Kulesha, E., Martin, F. J., Maurel, T., McLaren,
286 W. M., Murphy, D. N., Nag, R., Overduin, B., Pignatelli, M., Pritchard, B., Pritchard, E., Riat,
287 H. S., Ruffier, M., Sheppard, D., Taylor, K., Thormann, A., Trevanion, S. J., Vullo, A., Wilder,
288 S. P., Wilson, M., Zadissa, A., Aken, B. L., Birney, E., Cunningham, F., Harrow, J., Herrero, J.,
289 Hubbard, T. J., Kinsella, R., Muffato, M., Parker, A., Spudich, G., Yates, A., Zerbino, D. R., and
290 Searle, S. M. (2014). Ensembl 2014. *Nucleic Acids Res.*, 42(D1):D749–D755.
- 291 Gehrig, J., Reischl, M., Kalmar, E., Ferg, M., Hadzhiev, Y., Zaucker, A., Song, C., Schindler, S.,

292 Liebel, U., and Muller, F. (2009). Automated high-throughput mapping of promoter-enhancer
293 interactions in zebrafish embryos. *Nat Meth*, 6(12):911–916.

294 Gerstein, M. B., Lu, Z. J., Van Nostrand, E. L., Cheng, C., Arshinoff, B. I., Liu, T., Yip, K. Y., Ro-
295 bilotto, R., Rechtsteiner, A., Ikegami, K., Alves, P., Chateigner, A., Perry, M., Morris, M., Auer-
296 bach, R. K., Feng, X., Leng, J., Vielle, A., Niu, W., Rhrissorrakrai, K., Agarwal, A., Alexander,
297 R. P., Barber, G., Brdlik, C. M., Brennan, J., Brouillet, J. J., Carr, A., Cheung, M.-S., Clawson,
298 H., Contrino, S., Dannenberg, L. O., Dernburg, A. F., Desai, A., Dick, L., Dosé, A. C., Du, J.,
299 Egelhofer, T., Ercan, S., Euskirchen, G., Ewing, B., Feingold, E. A., Gassmann, R., Good, P. J.,
300 Green, P., Gullier, F., Gutwein, M., Guyer, M. S., Habegger, L., Han, T., Henikoff, J. G., Henz,
301 S. R., Hinrichs, A., Holster, H., Hyman, T., Iniguez, A. L., Janette, J., Jensen, M., Kato, M.,
302 Kent, W. J., Kephart, E., Khivansara, V., Khurana, E., Kim, J. K., Kolasinska-Zwierz, P., Lai,
303 E. C., Latorre, I., Leahey, A., Lewis, S., Lloyd, P., Lochovsky, L., Lowdon, R. F., Lubling, Y.,
304 Lyne, R., MacCoss, M., Mackowiak, S. D., Mangone, M., McKay, S., Mecnas, D., Merrihew,
305 G., Miller, D. M., Muroyama, A., Murray, J. I., Ooi, S.-L., Pham, H., Phippen, T., Preston,
306 E. A., Rajewsky, N., Rättsch, G., Rosenbaum, H., Rozowsky, J., Rutherford, K., Ruzanov, P.,
307 Sarov, M., Sasidharan, R., Sboner, A., Scheid, P., Segal, E., Shin, H., Shou, C., Slack, F. J.,
308 Slightam, C., Smith, R., Spencer, W. C., Stinson, E. O., Taing, S., Takasaki, T., Vafeados, D.,
309 Voronina, K., Wang, G., Washington, N. L., Whittle, C. M., Wu, B., Yan, K.-K., Zeller, G., Zha,
310 Z., Zhong, M., Zhou, X., , Ahringer, J., Strome, S., Gunsalus, K. C., Micklem, G., Liu, X. S.,
311 Reinke, V., Kim, S. K., Hillier, L. W., Henikoff, S., Piano, F., Snyder, M., Stein, L., Lieb, J. D.,
312 and Waterston, R. H. (2010). Integrative analysis of the *Caenorhabditis elegans* genome by the

- 313 modencode project. *Science*, 330(6012):1775–1787.
- 314 Ghanbarian, A. T. and Hurst, L. D. (2015). Neighboring genes show correlated evolution in gene
315 expression. *Mol. Biol. Evol.*, 32(7):1748–1766.
- 316 Gierman, H. J., Indemans, M. H., Koster, J., Goetze, S., Seppen, J., Geerts, D., van Driel, R., and
317 Versteeg, R. (2007). Domain-wide regulation of gene expression in the human genome. *Genome*
318 *Res.*, 17(9):1286–1295.
- 319 Hammonds, A. S., Bristow, C. A., Fisher, W. W., Weiszmann, R., Wu, S., Hartenstein, V., Kellis,
320 M., Yu, B., Frise, E., and Celniker, S. E. (2013). Spatial expression of transcription factors in
321 drosophila embryonic organ development. *Genome Biol.*, 14(12):R140.
- 322 He, B., Chen, C., Teng, L., and Tan, K. (2014). Global view of enhancer–promoter interactome in
323 human cells. *Proc. Natl. Acad. Sci. USA*, 111(21):E2191–E2199.
- 324 Heger, P., Marin, B., and Schierenberg, E. (2009). Loss of the insulator protein CTCF during
325 nematode evolution. *BMC Mol. Biol.*, 5:1–14.
- 326 Kagey, M. H., Newman, J. J., Bilodeau, S., Zhan, Y., Orlando, D. A., van Berkum, N. L., Ebmeier,
327 C. C., Goossens, J., Rahl, P. B., Levine, S. S., Taatjes, D. J., Dekker, J., and Young, R. A.
328 (2010). Mediator and cohesin connect gene expression and chromatin architecture. *Nature*,
329 467(7314):430–435.
- 330 Liu, M., Maurano, M. T., Wang, H., Qi, H., Song, C.-z., Navas, P. A., Emery, D. W., Stamatoy-
331 annopoulos, J. A., and Stamatoyannopoulos, G. (2015). Genomic discovery of potent chromatin
332 insulators for human gene therapy. *Nat. biotechnol.*, 33(2):198—203.

- 333 Ly, K., Reid, S. J., and Snell, R. G. (2015). Rapid RNA analysis of individual *Caenorhabditis*
334 *elegans*. *MethodsX*, 2:59–63.
- 335 Malik, S. and Roeder, R. G. (2010). The metazoan mediator co-activator complex as an integrative
336 hub for transcriptional regulation. *Nat. Rev. Genet.*, 11(11):761–772.
- 337 Marsman, J. and Horsfield, J. A. (2012). Long distance relationships: Enhancer–promoter com-
338 munication and dynamic gene transcription. *Biochimica et Biophysica Acta (BBA) - Gene Reg-*
339 *ulatory Mechanisms*, 1819(11–12):1217 – 1227.
- 340 McKinney, W. (2010). Data structures for statistical computing in python. In van der Walt, S. and
341 Millman, J., editors, *Proceedings of the 9th Python in Science Conference*, pages 51–56.
- 342 Mi, H., Poudel, S., Muruganujan, A., Casagrande, J. T., and Thomas, P. D. (2016). Panther
343 version 10: expanded protein families and functions, and analysis tools. *Nucleic Acids Res.*,
344 44(D1):D336–D342.
- 345 Michalak, P. (2008). Coexpression, coregulation, and cofunctionality of neighboring genes in
346 eukaryotic genomes. *Genomics*, 91(3):243 – 248.
- 347 Negre, N., Brown, C. D., Shah, P. K., Kheradpour, P., Morrison, C. A., Henikoff, S., Kellis, M., and
348 White, K. P. (2010). A Comprehensive Map of Insulator Elements for the *Drosophila* Genome.
349 *Plos Genet.*, 6(1):e1000814.
- 350 Okkema, P. G., Harrison, S. W., Plunger, V., Aryana, A., and Fire, A. (1993). Sequence Re-
351 quirements for Myosin Gene Expression and Regulation in *C. elegans*. *Genetics*, 135(Waterston
352 1988):385–404.

- 353 Ouedraogo, M., Bettembourg, C., Bretaudeau, A., Sallou, O., Diot, C., Demeure, O., and Lecerf, F.
354 (2012). The duplicated genes database: Identification and functional annotation of co-localised
355 duplicated genes across genomes. *PLoS ONE*, 7(11):1–8.
- 356 Quinlan, A. R. and Hall, I. M. (2010). Bedtools: a flexible suite of utilities for comparing genomic
357 features. *Bioinformatics*, 26(6):841–842.
- 358 Rao, S. S. P., Huntley, M. H., Durand, N. C., Stamenova, E. K., Bochkov, I. D., Robinson, J. T.,
359 Sanborn, A. L., Machol, I., Omer, A. D., Lander, E. S., and Aiden, E. L. (2014). A 3d map
360 of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell*,
361 159(7):1665–1680.
- 362 Ringrose, L., Chabanis, S., Angrand, P. O., Woodroffe, C., and Stewart, A. F. (1999). Quantitative
363 comparison of dna looping in vitro and in vivo: chromatin increases effective dna flexibility at
364 short distances. *EMBO J.*, 18(23):6630–6641.
- 365 Schwartz, Y. B., Linder-basso, D., Kharchenko, P. V., Tolstorukov, M. Y., Kim, M., Li, H.-b.,
366 Gorchakov, A. A., Minoda, A., Shanower, G., Alekseyenko, A. A., Riddle, N. C., Jung, Y. L.,
367 Gu, T., Plachetka, A., Elgin, S. C. R., Kuroda, M. I., Park, P. J., Savitsky, M., and Karpen,
368 G. H. (2012). Nature and function of insulator protein binding sites in the *Drosophila* genome.
369 *Genome Res.*, 11:2188–2198.
- 370 Stiernagle, T. (2006). Maintenance of *c. elegans*. In *C. elegans Research Community*, editor,
371 *WormBook*. WormBook.

- 372 The ENCODE Project Consortium (2012). An integrated encyclopedia of dna elements in the
373 human genome. *Nature*, 489(7414):57–74.
- 374 Tomancak, P., Beaton, A., Weizmann, R., Kwan, E., Shu, S., Lewis, S. E., Richards, S., Ash-
375 burner, M., Hartenstein, V., Celniker, S. E., and Rubin, G. M. (2002). Systematic determina-
376 tion of patterns of gene expression during drosophila embryogenesis. *Genome Biol.*, 3(12):re-
377 search0088.1–88.14.
- 378 Tomancak, P., Berman, B. P., Beaton, A., Weizmann, R., Kwan, E., Hartenstein, V., Celniker,
379 S. E., and Rubin, G. M. (2007). Global analysis of patterns of gene expression during drosophila
380 embryogenesis. *Genome Biol.*, 8(7):R145.
- 381 van Arensbergen, J., van Steensel, B., and Bussemaker, H. J. (2014). In search of the determinants
382 of enhancer–promoter interaction specificity. *Trends in cell biology*, 24(11):695–702.
- 383 Wagner, G. P., Kin, K., and Lynch, V. J. (2012). Measurement of mrna abundance using rna-seq
384 data: Rpkm measure is inconsistent among samples. *Theory Biosci.*, 131(4):281–285.
- 385 Xie, X., Mikkelsen, T. S., Gnirke, A., Lindblad-toh, K., Kellis, M., and Lander, E. S. (2007).
386 Systematic discovery of regulatory motifs in conserved regions of the human genome , including
387 thousands of CTCF insulator sites. *Proc. Natl. Acad. Sci. USA*, 104(17):7145–7150.
- 388 Zhang, Y., Chen, D., Smith, M. A., Zhang, B., and Pan, X. (2012). Selection of reliable reference
389 genes in /textitCaenorhabditis elegans for analysis of nanotoxicity. *PLoS ONE*, 7(3):1–7.

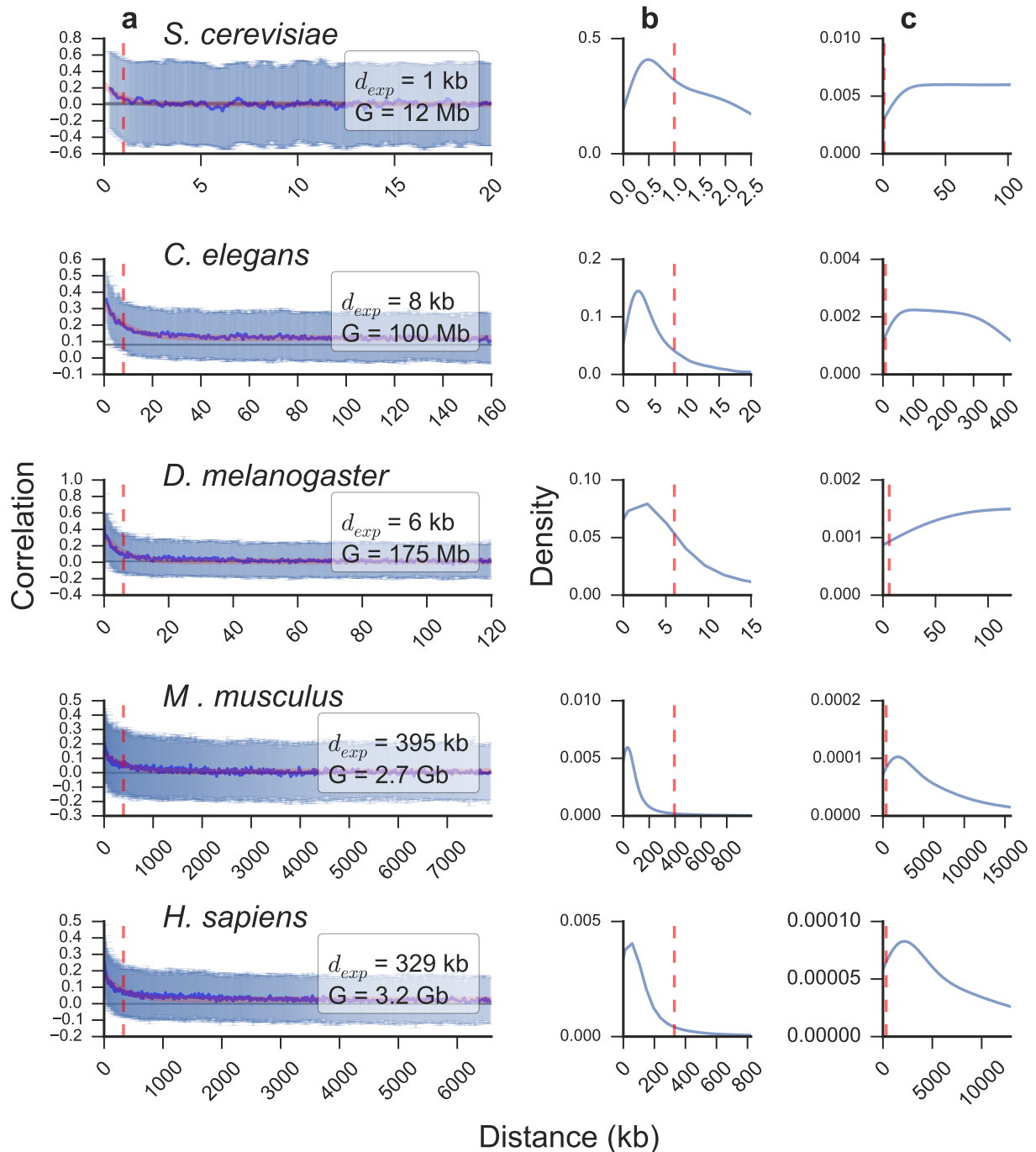


Figure 1: Neighboring genes are transcriptionally correlated genome-wide across eukaryotes. a) Sliding median of correlations between paired neighbors (blue line) and interquartile range (pale blue) with increasing intergenic distance. Median \pm 95% confidence interval of randomly paired genes is shown as a horizontal gray line. Fit to an exponential decay function (red line) was used to compute the mean distance at which gene neighbors remain correlated (d_{exp} , vertical red dashed line). The genome size (G) is displayed for each organism. Distribution of intergenic distances between each gene and its nearest neighbor (b) and all paired genes (c). The organism analyzed in each case is indicated for each group of three subplots.

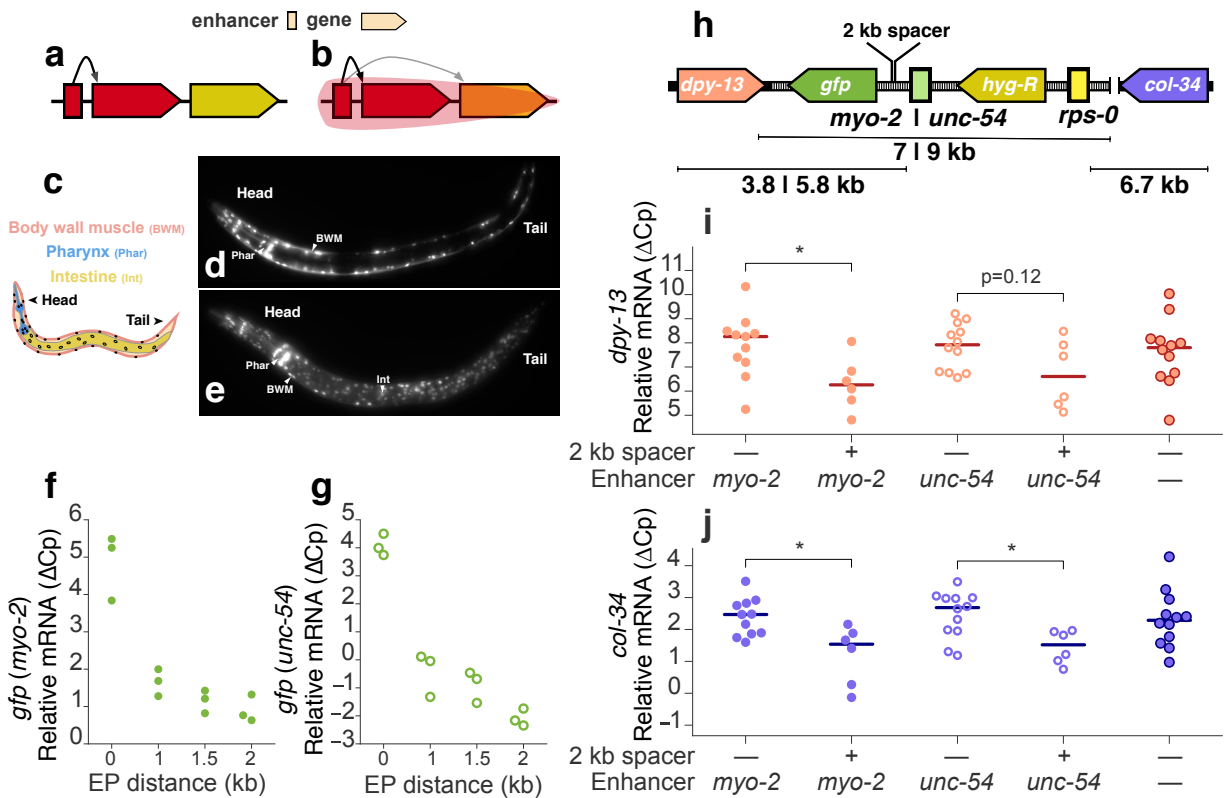


Figure 2: Enhancer sharing explains the transcriptional correlation of gene neighbors. Two possible models for EP relationship: a) Enhancers have specific target genes and b) enhancers have a range of action in which they influence genes by physical proximity. Tissue specific enhancers (c) are generally compatible. Pharynx and body wall muscle (d) and pharynx, body wall muscle and intestine (e) enhancers driving nuclear *gfp* expression. mRNA levels of *gfp* with increasing EP distance for lines with *myo-2* (filled circles, f) and *unc-54* (hollow circles, g) enhancers. h) Genomic context of the integration site. The inserted construct is shown over a dashed black line and includes a hygromycin resistance gene (*hyg-R*) regulated by a ribosomal enhancer (*rps-0*) and promoter in addition to the reporter (*gfp*) regulated by either the *myo-2* or *unc-54* enhancers; the native genes *dpy-13* and *col-34* flank the insertion site. Relative mRNA levels of *dpy-13* (i) and *col-34* (j) in wild-type and lines with and without the 2 kb spacer (*two tailed P-val<0.05, Mann-Whitney U test). The difference in crossing point-PCR-cycle (Δ Cp) with the reference gene *pmp-3* and the corresponding median for each group of biological replicates is shown for every qPCR experiment.

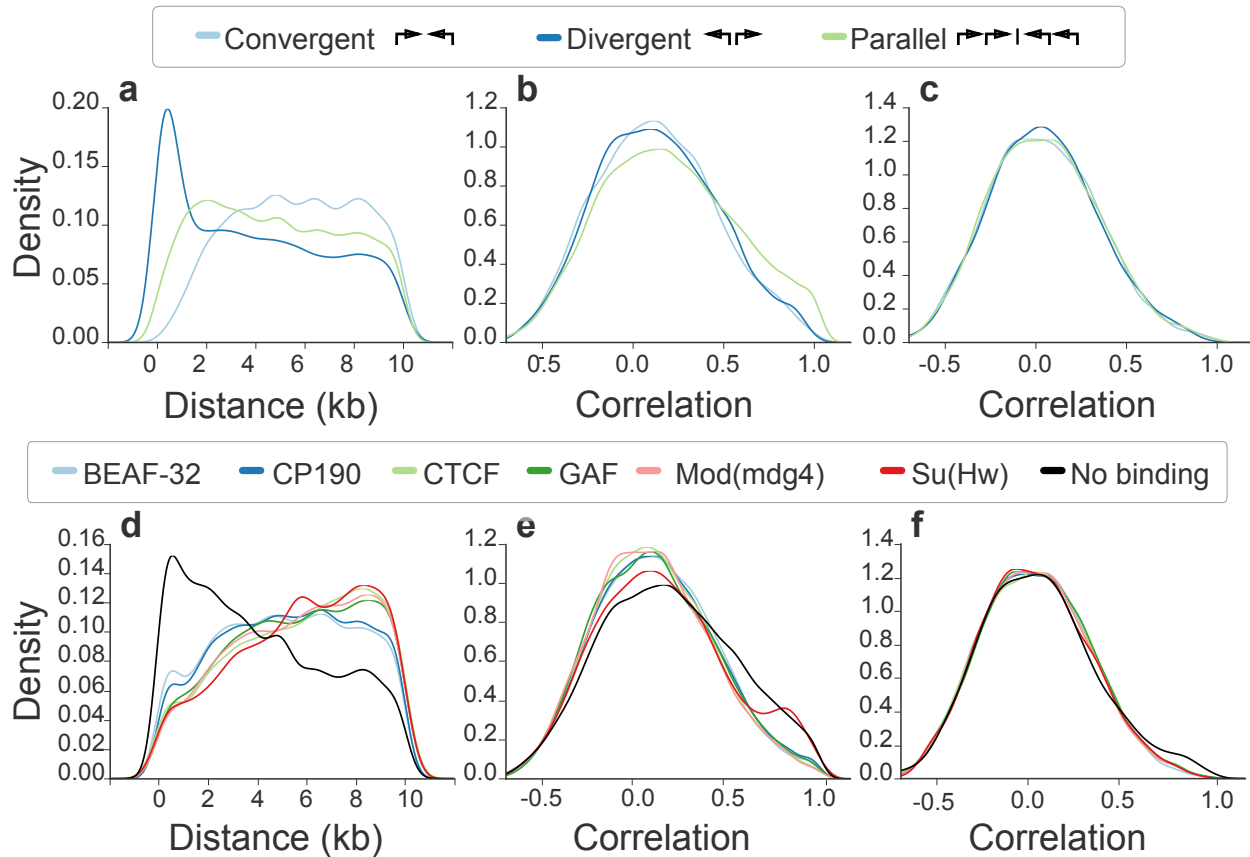
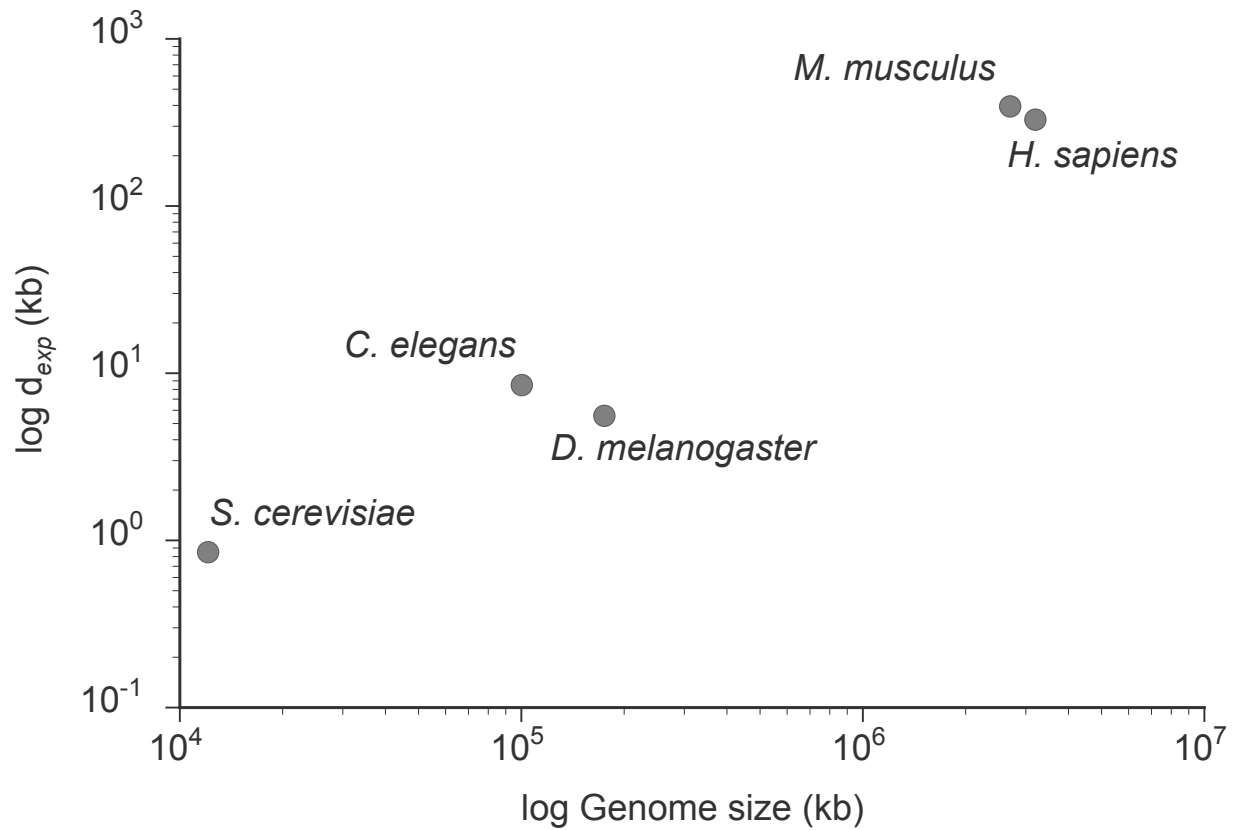
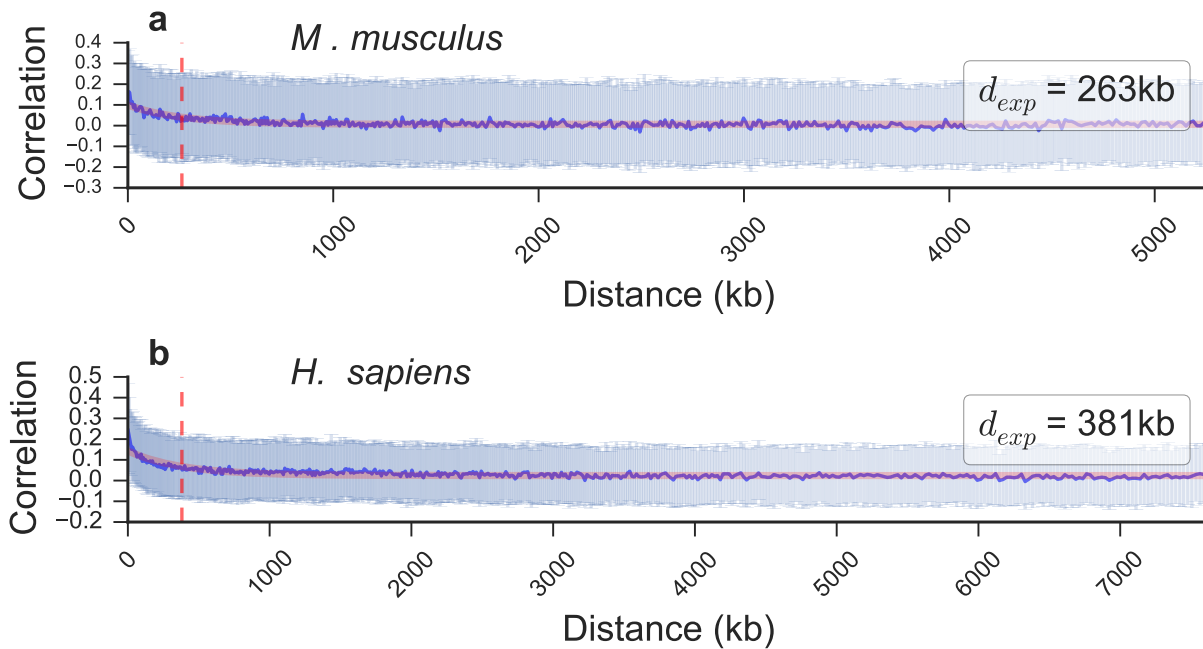


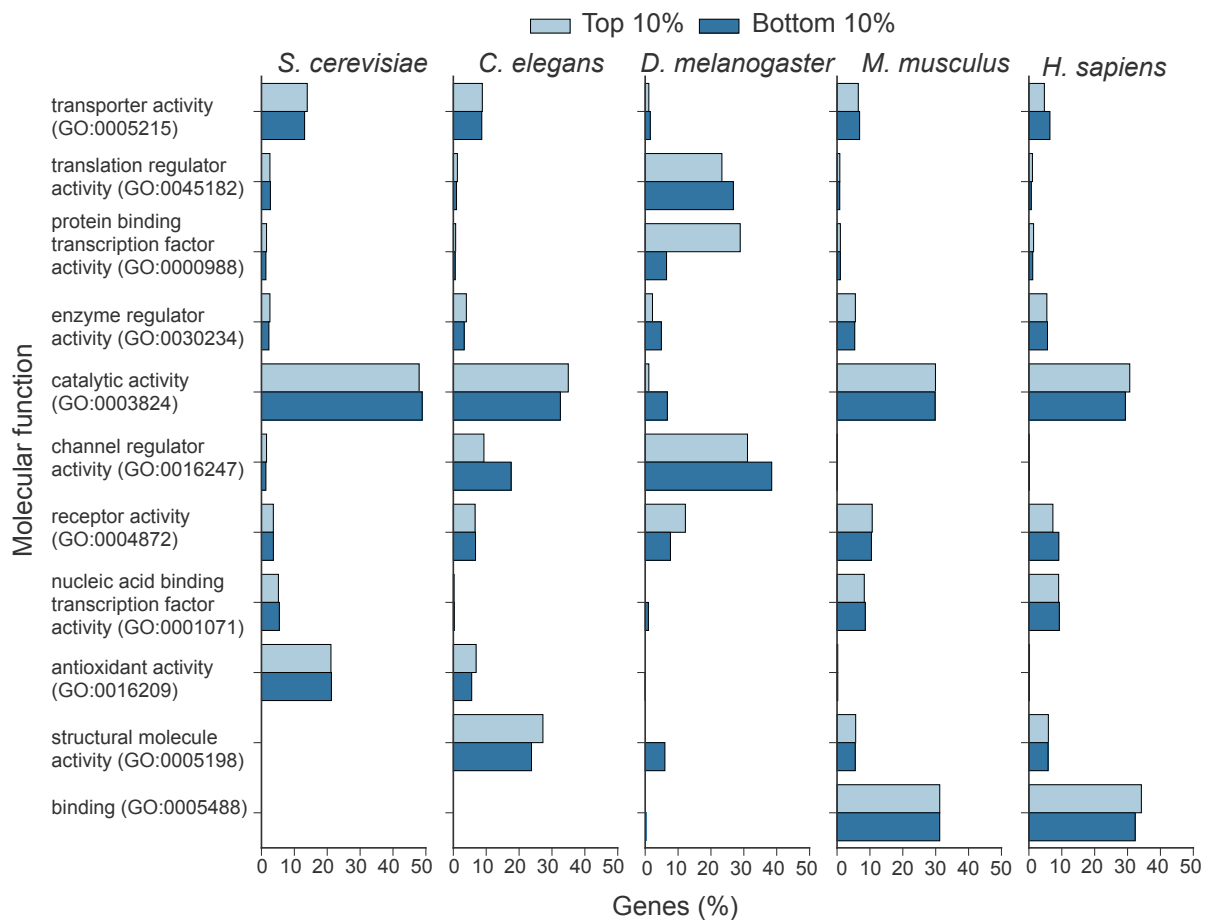
Figure 3: EP distance causes gene orientation-dependent correlation and provides regulatory independence to gene neighbors. Distribution of intergenic distances below 10 kb of gene pairs in *D. melanogaster* by configuration (~5 to 18 thousand gene pairs for each group, a) and flanking insulator binding sites identified through ChIP-chip (Negre et al., 2010) (~5 to 15 thousand pairs for each group, d). The corresponding distribution of correlations is shown for the same gene pairs (b, e) and pairs with controlled distributions of intergenic distances between 30 and 40 kb (~7 to 14 thousand pairs for gene orientation groups, ~10 to 18 thousand for insulator groups, c, f).



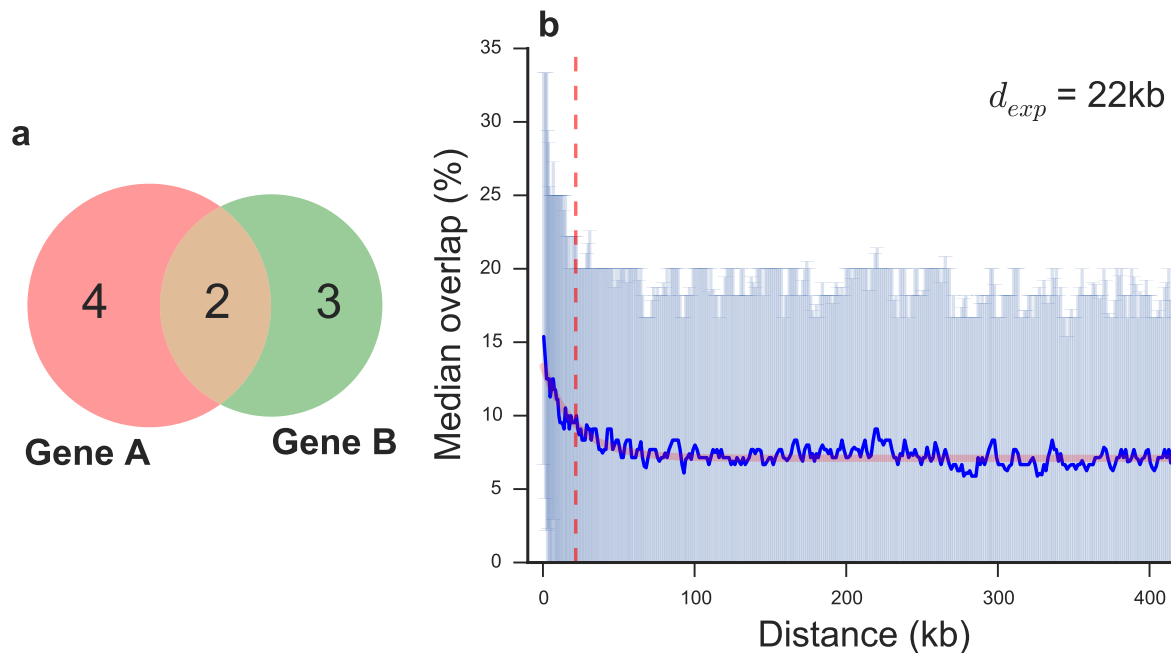
Supplementary Figure 1: The distance at which a pair of genes remain correlated (d_{exp}) scales with genome size.



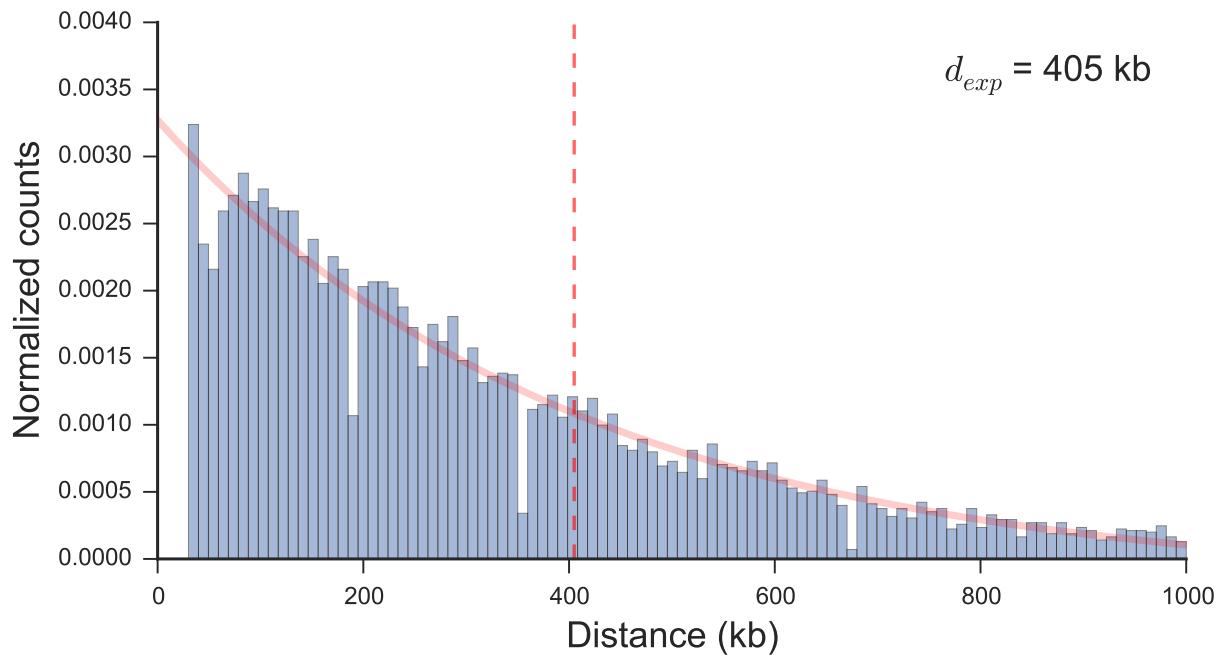
Supplementary Figure 2: Removing duplicated genes does not affect the overall correlation of gene neighbors. Sliding median of correlations between paired neighbors (blue line) and interquartile range (pale blue) with increasing intergenic distance in *M. musculus* (a) and *H. sapiens* (b) after removing duplicated gene pairs. Fit to exponential decay function (red line) and corresponding d_{exp} (red dashed line) are shown.



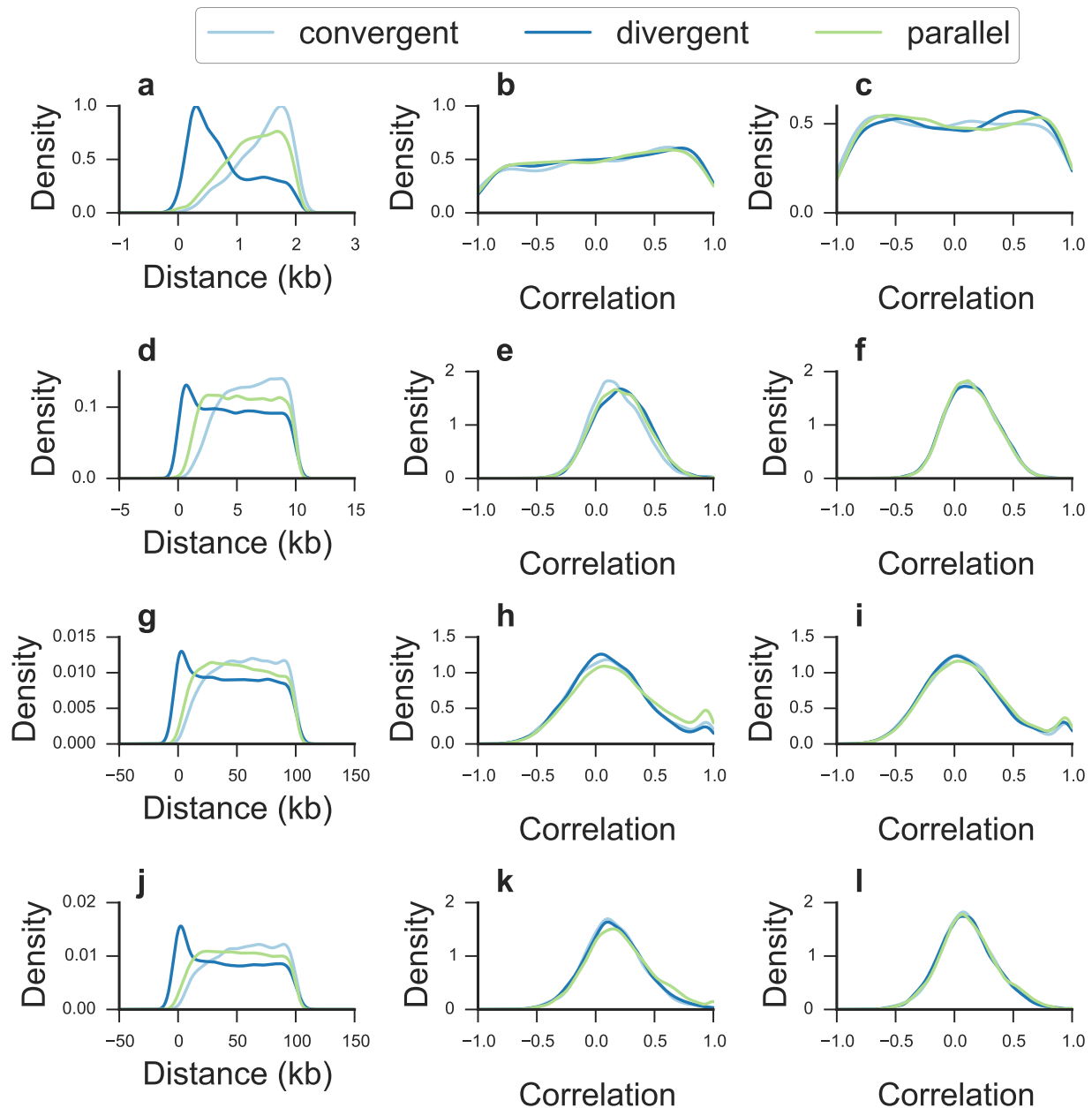
Supplementary Figure 3: Representation of gene ontology annotations remains unbiased in correlated gene pairs. The molecular function classification of top and bottom 10% correlated gene pairs with intergenic distance below d_{exp} is shown for each organism.



Supplementary Figure 4: Gene pairs are correlated in spatial expression in *D. melanogaster*. The size of the intersection between the set of tissues in which each gene of a given pair is expressed was divided over the size of the union of the same sets. An example is shown in (a), where the percentage overlap is $2/(4+3-2)=0.4$. b) Sliding median of the percentage overlap in tissue specific expression (blue line) and interquartile range (pale blue) with increasing intergenic distance. Fit to exponential decay function (red line) and corresponding d_{exp} (red dashed line) are shown.



Supplementary Figure 5: Chromatin looping decreases exponentially with distance in human cell lines. Normalized count of loops identified through HiC by Rao et al. (2014) were fit to exponential decay function (red line); the resulting d_{exp} (red dashed line) is shown.



Supplementary Figure 6: Gene orientation effect in correlation of gene pairs is explained by EP distance. Distribution of intergenic distances and the corresponding distribution of correlations of gene pairs is shown in the first and second columns, respectively; correlations after controlling for intergenic distance are shown in the third column. The range of distances between paired genes for each plot is as follows: *S. cerevisiae* below 2 kb (a,b) and between 2 and 4 kb (c). *C. elegans* below 10 kb (d,e) and between 10 and 20 kb (f). *H. sapiens* and *M. musculus* below 100 kb (g, h, j, k) and between 100 and 200 kb (i, l).