

Supporting Materials for “Core promoter T-blocks correlate with gene expression levels in *C. elegans*”

Uladzislau Hryshkevich, Tamar Hashimshony and Itai Yanai

Supplemental Information

Figure S1. T-blocks and A-blocks have similar effect on nucleosome occupancy.

Figure S2. Correlation between T-block occurrences and expression level in expression data from adult and L1 (first larvae) worms.

Figure S3. The correlation between T-block occurrences and expression level is unique to the core promoter.

Figure S4. 3'UTR length correlates with gene expression level.

Figure S5. The T/A rich tail of the SL1 motif is not sufficient to explain the correlation between the number of T-blocks and the presence of the SL1 motif.

Figure S6. The correlation between expression levels and the number of T-blocks is not explained by misannotated operons.

Figure S7. Thymines and T-blocks are enriched relative to adenines and A-blocks in the core promoters of many metazoans.

Figure S8. The distribution of the lengths of 5'UTRs.

Table S1. Motif frequencies across *Caenorhabditis* core promoters.

Table S2. Distribution of T-blocks and A-blocks in the core promoters.

Table S3. Combinatorial analysis.

Table S4. GO analysis.

Table S5. SNPs identified between the CB4856 (Hawaiian) and the reference N2 (Bristol) *C. elegans* strain.

Figure captions

Figure S1. T-blocks and A-blocks have similar effect on nucleosome occupancy. To compare the effect of T-blocks and A-blocks on nucleosome positioning, we composed 4 groups of genes: 1. genes with many (≥ 4.5) T-blocks and many (≥ 3.3) A-blocks, 2. genes with many (≥ 4.5) T-blocks, but few (≤ 1.4) A-blocks, 3. genes with few (≤ 1.9) T-blocks and many (≥ 3.3) A-blocks, and 4. genes with few (≤ 1.9) T-blocks and few (≤ 1.4) A-blocks. The nucleosome occupancy of different groups was compared by Kolmogorov-Smirnov test. The distribution of mean \log_2 nucleosome occupancy of the core promoters in the four gene groups. *P*-values indicate the result of the Komagorov-Smirnov test between each pair of distributions.

Figure S2. Correlation between T-block occurrences and expression level in expression data from adult and L1 (first larvae) worms. In a similar fashion to Figure 4, for each set of genes with a given level of expression, the distribution of the number of T-blocks is shown as a heat map. Twenty equally populated groups of genes were defined based upon their gene expression levels. Most genes have a low level of expression, and thus the ten groups with the lowest expression levels were merged into one, '1' group. The remaining ten are shown in the figures as groups 11 through 20. The distribution of T-blocks of each category was compared with that of the '1' group and the *P*-value is indicated to the right. B. same as A for A-blocks. C. and D. are the same as A and B but for expression data from L1 worms. Expression data were obtained from GEO DataSets (adult expression: GSM543285, GSM543286 and GSM543287, L1 expression: GSM146422, GSM146423 and GSM147330) (Irazoqui et al. 2010; Kirienko and Fay 2007) .

Figure S3. The correlation between T-block occurrences and expression level is unique to the core promoter. The presentation is the same as in Figure 5. For each of 15 gene regions, including the proximal promoter, starts and ends of the first, second and third exons and introns, 3'UTR, and the 3' adjacent and 3' distal end regions the T-block and A-block frequency is shown for different expression groups. For exons and introns, Ensembl annotation was used to identify the 1st, 2nd, or 3rd, ignoring alternatively spliced introns/exons.

Figure S4. 3'UTR length correlates with gene expression level. Genes were binned according to the length of the 3'UTR. The boxplots show the distribution of expression level for each 3'UTR length bin.

Figure S5. The T/A rich tail of the SL1 motif is not sufficient to explain the correlation between the number of T-blocks and the presence of the SL1 motif. We truncated the core promoter (Figure 1) by 40bp from the 3'end to remove the T/A rich region of the SL1 motif. We recomputed the SL1-fuzziness and again mapped it onto the heatmap of fuzziness vs. block numbers as in figure 7. Since the results are nearly identical we conclude that the T/A rich tail alone cannot account for the observed correlation.

Figure S6. The correlation between expression levels and the number of T-blocks is not explained by misannotated operons. We repeated the analysis shown in Figure 4 on strictly non-operon genes, defined as genes with an upstream and downstream intergenic region of at least 1kb that are not defined as operon genes in Wormbase. The correlation is preserved in this stringently filtered set of genes.

Figure S7. Thymines and T-blocks are enriched relative to adenines and A-blocks in the core promoters of many metazoans. **A.** The distribution of the thymine fraction within the forward strands of various metazoan core promoters. **B.** T-blocks are more frequent in the core promoters of most metazoans than A-blocks. Genomes were downloaded from Ensembl (<http://www.ensembl.org/>).

Figure S8. The distribution of the lengths of 5'UTRs. According to Wormbase release 210, 82% of the genes have shorter than or equal to 60bp.

Tables

Table S1. Motif frequencies across *Caenorhabditis* core promoters.

Frequency of different motifs across *Caenorhabditis* core promoters are shown as percentage of the core promoters containing each motif. Motif frequencies were calculated for motifs identified with $P < 0.005$, FDR-corrected.

Motif Species	SL1	Kozak	TATA box	Sp1	T-blocks
<i>C. remanei</i>	56.88	29.21	5.74	1.84	43.91
<i>C. briggsae</i>	63.85	33.50	4.03	2.43	41.85
<i>C. brenneri</i>	61.76	32.69	6.74	1.5	43.36
<i>C. elegans</i>	58.95	37.19	6.34	1.23	41.87
<i>C. japonica</i>	50.53	31.46	4.85	1.2	22.98

Table S2. Distribution of T-blocks and A-blocks in the core promoters.

Numbers of genes with different quantities of T- and A-blocks are indicated. Difference between two distributions ($P < 10^{-307}$) were computed using the Kolmagorov-Smirnov test.

# of blocks	T-blocks	A-blocks
0	670	1039
1	1517	2623
2	2324	3746
3	2820	3736
4	3129	3043
5	2712	1784
6	2137	892
7	1276	331
8	556	87
9	139	15
10	19	3

Table S3. Combinatorial analysis shows that SL1 and T-blocks are enriched for co-occurrences, while both are depleted in co-occurrences with TATA-box.

P-values are computed using the hypergeometric distribution. Occurrences were defined using the MEME-motifs and a MAST threshold of $P < 0.005$, FDR-corrected. Green *P*-values indicate enrichment in motif co-occurrence, while red *P*-values indicate depletion in motif co-occurrences.

Species \ Motifs	SL1 and T-blocks co-occurrences	SL1 and TATA-box co-occurrences	T-blocks and TATA-box co-occurrences
<i>C. remanei</i>	10^{-16}	10^{-72}	10^{-150}
<i>C. briggsae</i>	10^{-16}	10^{-52}	10^{-89}
<i>C. brenneri</i>	10^{-16}	10^{-115}	10^{-101}
<i>C. elegans</i>	10^{-16}	10^{-35}	10^{-87}
<i>C. japonica</i>	10^{-16}	10^{-22}	10^{-32}

Table S4. GO analysis. The table indicates the top 15 GO terms enriched in genes containing both SL1 and T-blocks, and genes containing a TATA-box.

Blue GO terms indicate association with development, growth and reproduction; while red GO terms indicate association with stress.

#	SL1-T-blocks Gene Ontology Terms	$P <$	TATA-box Gene Ontology Terms	$P <$
1	reproduction	10^{-125}	nucleosome assembly	10^{-125}
2	embryonic development ending in birth or egg hatching	10^{-11}	nucleosome	10^{-12}
3	nematode larval development	10^{-11}	extracellular region	10^{-10}
4	growth	10^{-10}	lysozyme activity	10^{-5}
5	cell division	10^{-10}	aminoacylase activity	10^{-4}
6	structural constituent of ribosome	10^{-10}	transferase activity, transferring hexosyl groups	10^{-4}
7	cytokinesis	10^{-10}	monooxygenase activity	10^{-4}
8	ribosome	10^{-10}	binding	10^{-4}
9	translation	10^{-10}	iron ion binding	10^{-4}
10	intracellular	10^{-8}	body morphogenesis	10^{-4}
11	hermaphrodite genitalia development	10^{-7}	2,3-dihydro-2,3-dihydroxybenzoate dehydrogenase activity	10^{-4}
12	positive regulation of growth rate	10^{-6}	enterobactin biosynthetic process	10^{-4}
13	nucleic acid binding	10^{-6}	DNA binding	10^{-4}
14	molting cycle, collagen and cuticulin-based cuticle	10^{-5}	sugar binding	10^{-4}
15	negative regulation of vulval development	10^{-4}	heme binding	10^{-3}

Table S5. SNPs identified between the CB4856 (Hawaiian) and the reference N2 (Bristol) *C. elegans* strain. 78,119 SNPs were identified with a coverage of at least 5 and no heterogeneity. Chromosomal location is relative to WormBase release 195. Table is included as an Excel spreadsheet.

Reference:

- Irazoqui, J.E., E.R. Troemel, R.L. Feinbaum, L.G. Luhachack, B.O. Cezairliyan, and F.M. Ausubel. 2010. Distinct pathogenesis and host responses during infection of *C. elegans* by *P. aeruginosa* and *S. aureus*. *PLoS Pathog* **6**: e1000982.
- Kirienko, N.V. and D.S. Fay. 2007. Transcriptome profiling of the *C. elegans* Rb ortholog reveals diverse developmental roles. *Dev Biol* **305**: 674-684.

Figure S1

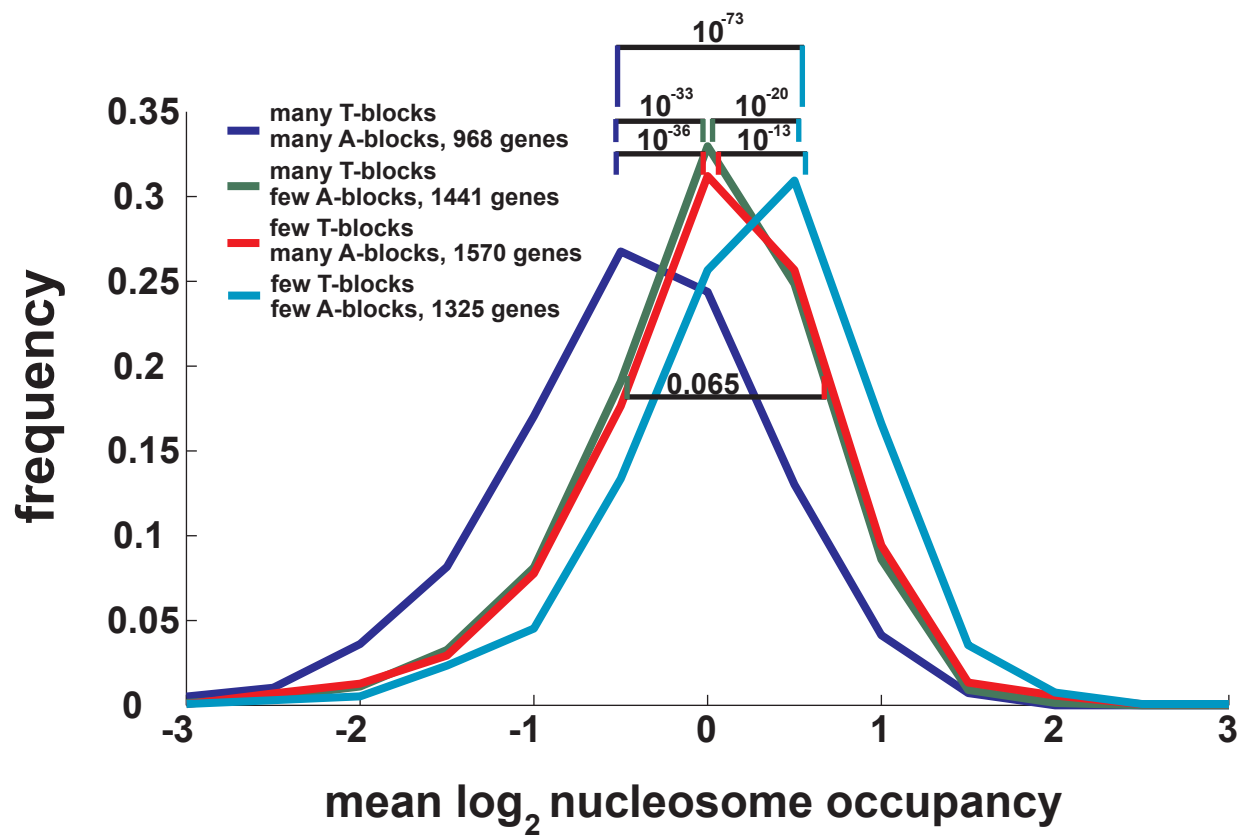


Figure S2

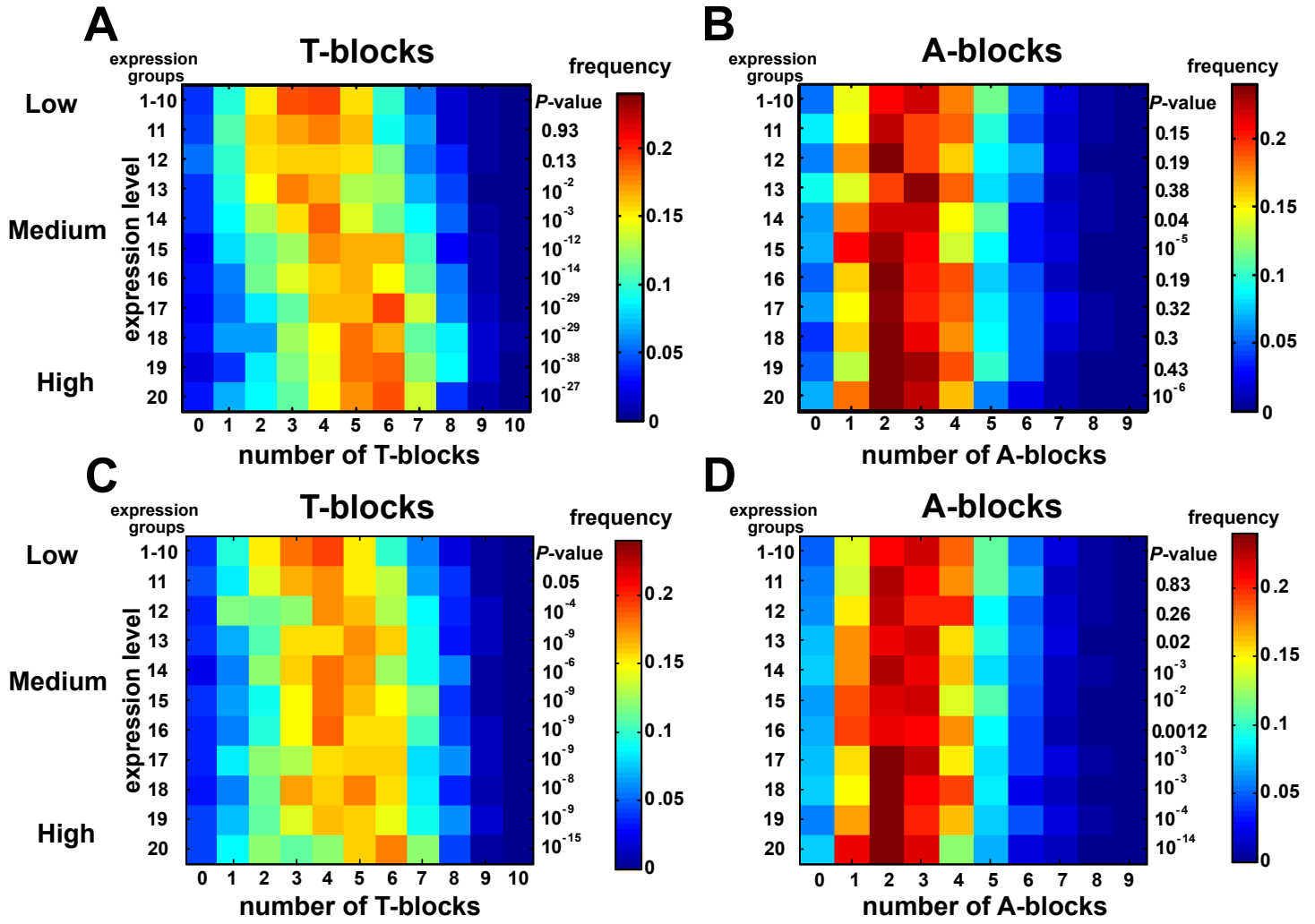


Figure S3

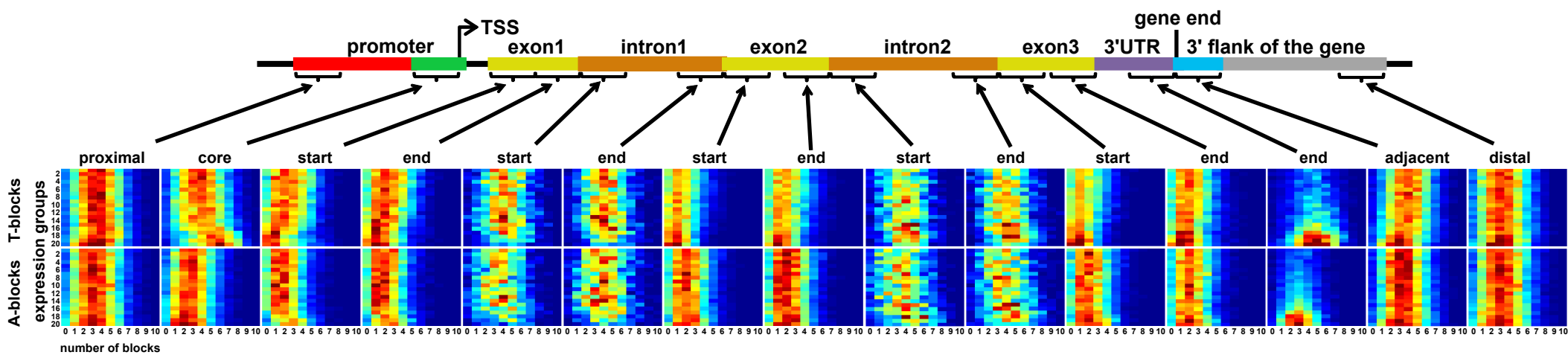


Figure S4

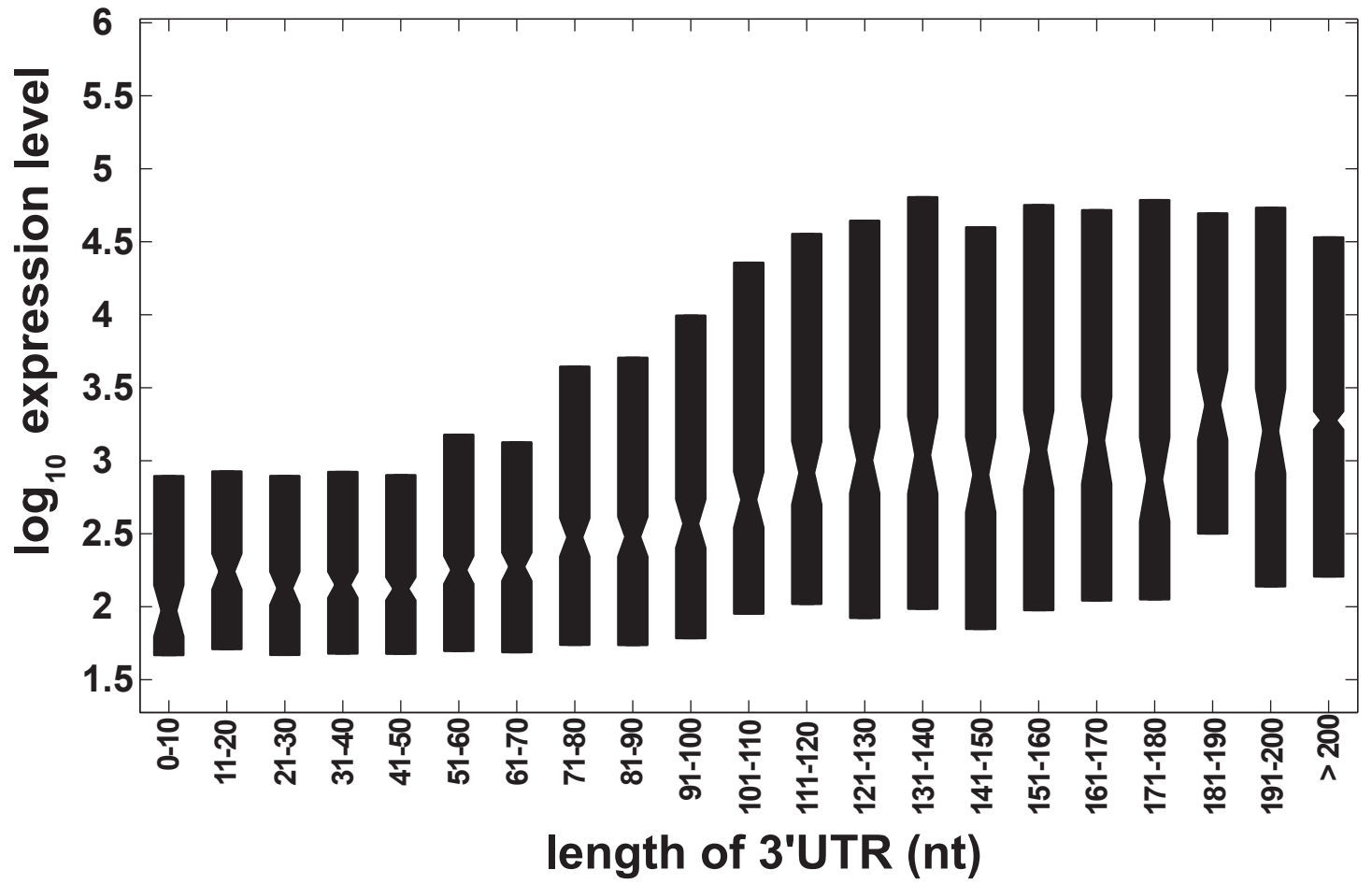


Figure S5

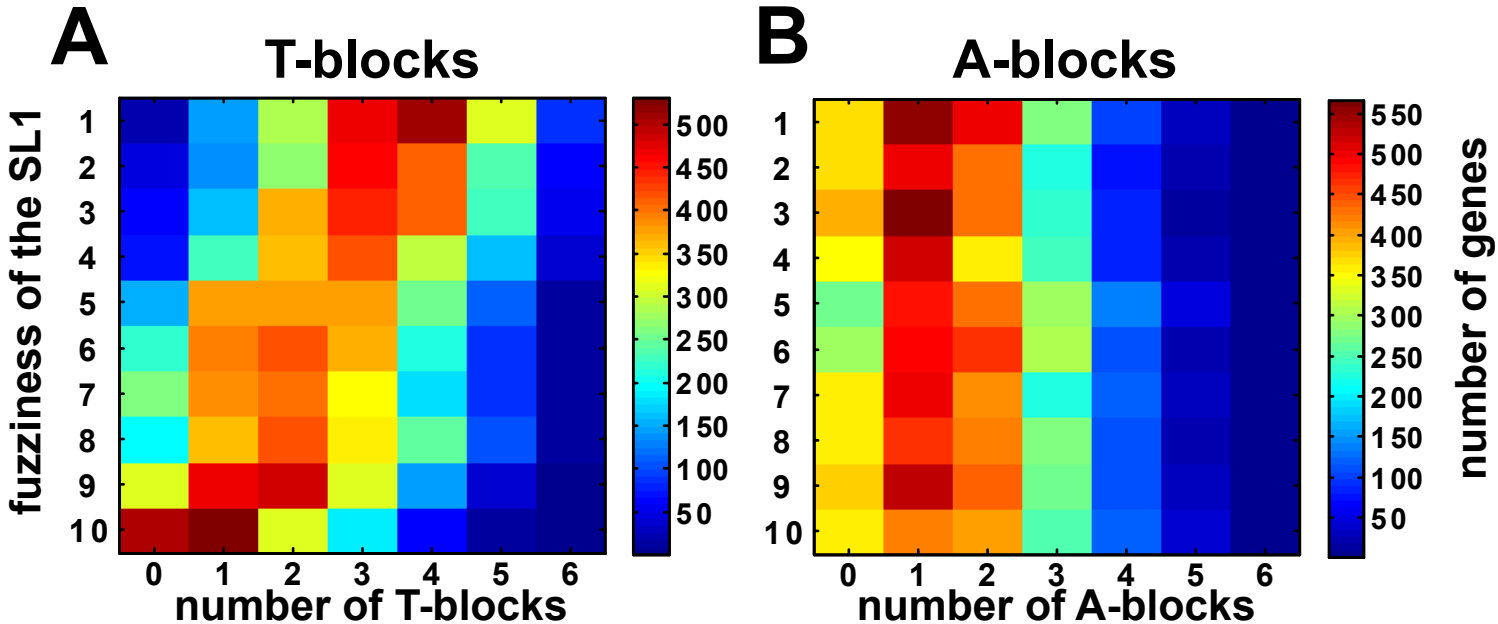


Figure S6

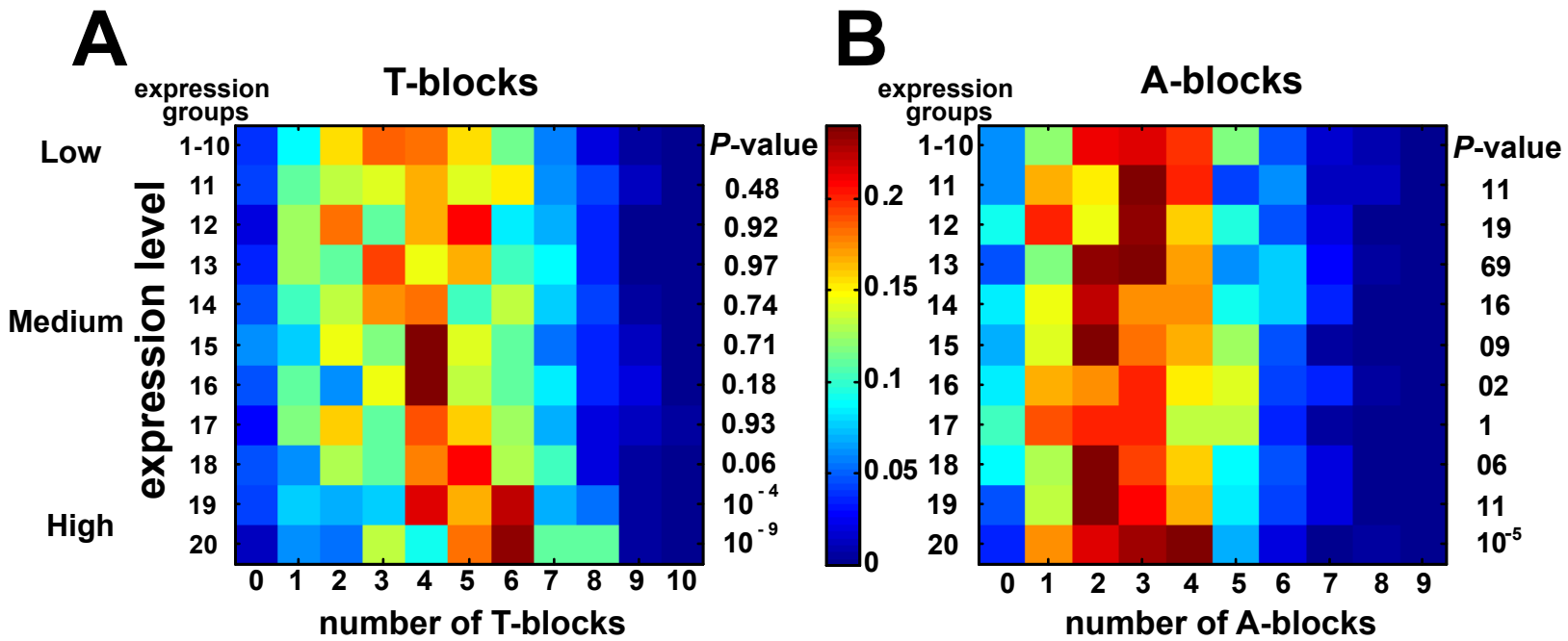
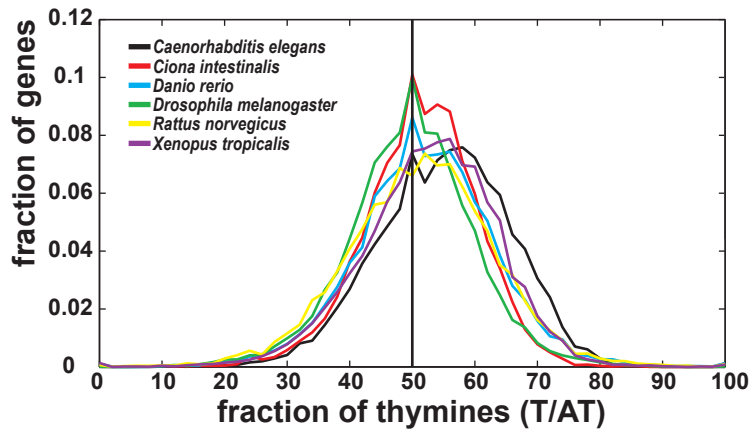


Figure S7

A



B

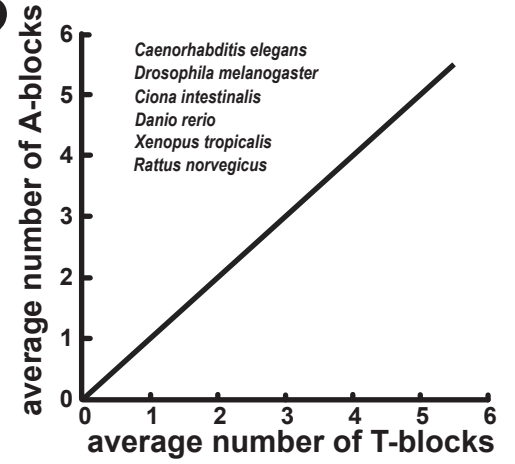


Figure S8

