

Supplementary Information

Predictive regulatory models in *Drosophila melanogaster* by integrative inference of transcriptional networks

Daniel Marbach[✉], Sushmita Roy[✉], Ferhat Ay, Patrick E. Meyer, Tamer Kahveci, Christopher A. Bristow, and Manolis Kellis

Companion Website

compbio.mit.edu/flynet Data and source code to reproduce networks and validation results (also included in Supplementary File 9).

Supplementary Figures

S1	Fraction of integrative network edges with and without physical support	4
S2	Degree distributions of integrative networks	5
S3	Structural properties of integrative networks	6
S4	Network motifs of integrative networks	7
S5	Prediction of novel functional annotations	8
S6	Expression prediction for eight representative target genes	9
S7	Squared errors for expression prediction using integrative and physical networks .	10
S8	Correlation of predictable and unpredictable genes across two time courses	11
S9	Definition of chromatin profiles	12
S10	Relative importance of input features for the unsupervised integrative network .	13

Supplementary Tables

S1	Performance for gene expressions prediction	14
----	---	----

Supplementary Notes

S1	Network motifs for integrative and mixed regulatory networks	15
----	--	----

Supplementary Files

- S1_integrative_regulatory_networks.zip** This zip file contains two tab-separated files representing the unsupervised and supervised integrative networks. The files have three columns: the first column is the name of the regulator, the second column is the name of the target gene, and the third column is the score of the edge connecting the regulator and the target. Genes are given by their FBgn numbers.
- S2_mixed_regulatory_networks.zip** The mixed transcriptional and post-transcriptional (TFs and microRNA) regulatory networks used for the network motif analysis described in Supplementary Note S1.
- S3_struct_properties_sup_network.zip** **Sheet 1:** Summary of the number of nodes and edges in the supervised integrative network and the mixed regulatory network. **Sheet 2:** Transcription factors in the supervised integrative and their in/out-degrees. **Sheet 3:** miRNAs in the mixed regulatory network and their in/out-degrees. **Sheet 4:** In-degrees of target genes and transcription factors of the supervised integrative network.
- S4_struct_properties_unsup_network.zip** **Sheet 1:** Summary of the number of nodes and edges in the unsupervised integrative network and the mixed regulatory network. **Sheet 2:** Transcription factors in the unsupervised integrative network and their in/out-degrees. **Sheet 3:** miRNAs in the mixed regulatory network and their in/out-degrees. **Sheet 4:** In-degrees of target genes and transcription factors of the unsupervised integrative network.
- S5_GO_enrichment_by_node_degree.xls** GO biological process enrichments for high (> 150) and low (= 1) in degree gene sets.
- S6_triads.xlsx** **Sheet 1:** List of all possible triads (subnetworks with only 3 nodes) for mixed regulatory networks. **Sheet 2:** List of all possible triads (subnetworks with only 3 nodes) for the integrative networks.
- S7_gene_function_prediction.zip** This archive has two files: One tgz file and sub-directories for the supervised and one for the unsupervised integrative network. Within each sub-directory is a file per GO process. Each line has four columns: first column is the name of the process, second column is the FBgn number of a gene, third column is the FDR, and fourth column is a label to identify the type of prediction. Specifically the labels take the values NEW (previously unannotated gene), KNOWN_PRED (additional new prediction for a previously annotated gene), KNOWN_HIT (previously annotated with the same function, that is, true positive). Only predictions that satisfy an FDR of

0.6 or lower and that are not in the FALSE set are included.

S8_ImaGO_enrichment_for_predicted_GO_terms.xls ImaGO enrichments for predicted GO terms for both unannotated genes and genes with previous annotations.

S9_datasets_and_source_code.zip Data and source code to reproduce networks and validation results.

Supplementary Figures

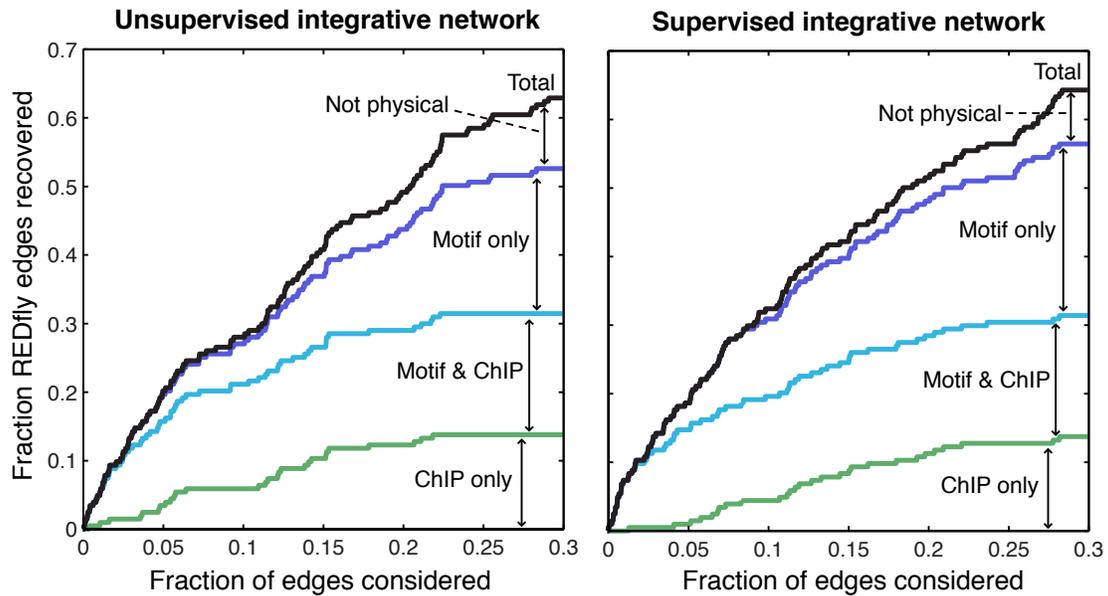


Figure S1: Fraction of integrative network edges with and without physical support
Recovery of known REDfly edges at different thresholds for the two integrative networks. The black curve (*total*) shows the fraction of the 204 REDfly edges recovered versus the fraction of predicted interactions considered for the integrative networks (cf. Fig. 2B). The area between the colored curves shows the fraction of these recovered edges that are supported by ChIP only, ChIP and conserved motifs, conserved motifs only, and functional datasets only (chromatin and expression). The majority of recovered edges are supported by ChIP and/or conserved motifs, especially the edges at the top of the prediction lists. However, at lower ranks some REDfly edges are also recovered based on functional data alone.

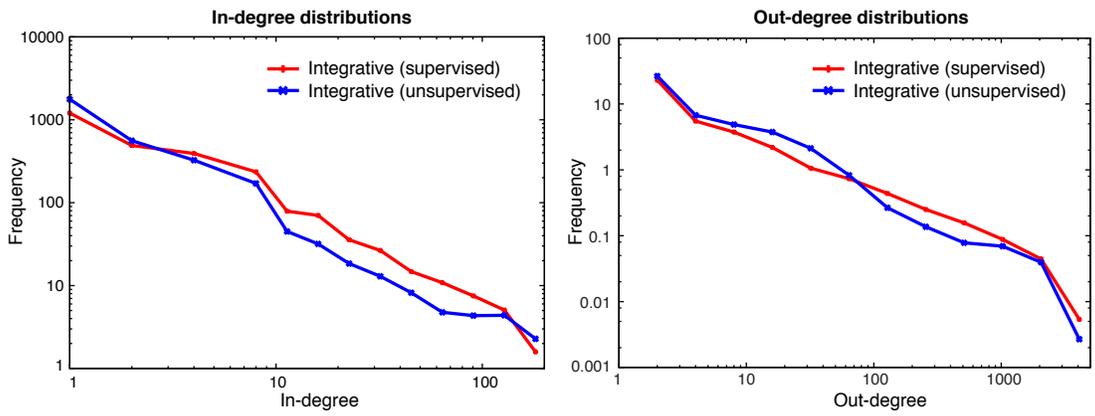


Figure S2: Degree distributions of integrative networks

Number of targets per TF (out-degree) and number of regulators per target (in-degree) for the supervised and unsupervised integrative networks, plotted on a log-log scale. Both networks show roughly a scale-free distribution for both in-degrees and out-degrees, with a small number of regulators showing many targets, and a small number of genes targeted by many regulators.

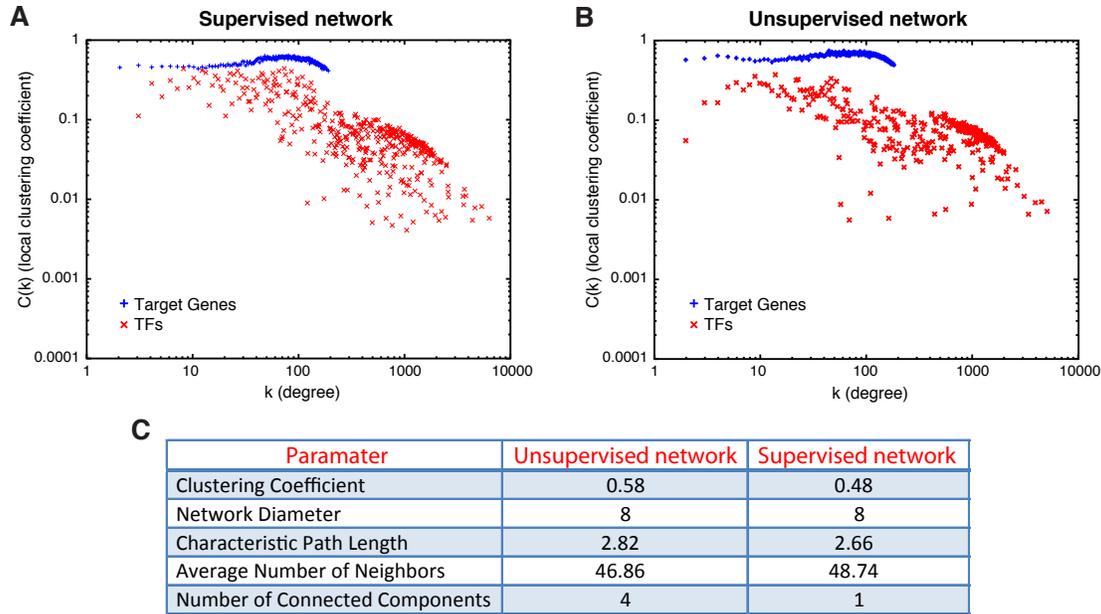


Figure S3: Structural properties of integrative networks

(A,B) Log-log plot of node degree (k) versus average clustering coefficient ($C(k)$) for the supervised and unsupervised integrative networks. Target genes show higher clustering coefficients as they are only connected to TFs, which themselves are highly interconnected leading to high $C(k)$ values for target genes. Target genes with the highest in-degrees (most heavily regulated) show decreased clustering coefficients, as their many regulators are less likely to be interconnected. In contrast, TFs are connected to both TFs and target-only genes (which themselves are not interconnected), resulting in overall lower coefficients. However, the same trend remains, where the most heavily-interconnected TFs show decreased clustering coefficients, as a higher proportion of their targets include predominantly target-only genes, which are not interconnected among themselves. (C) Clustering coefficients, networks diameters, characteristic path lengths, and number of connected components of the integrative networks.

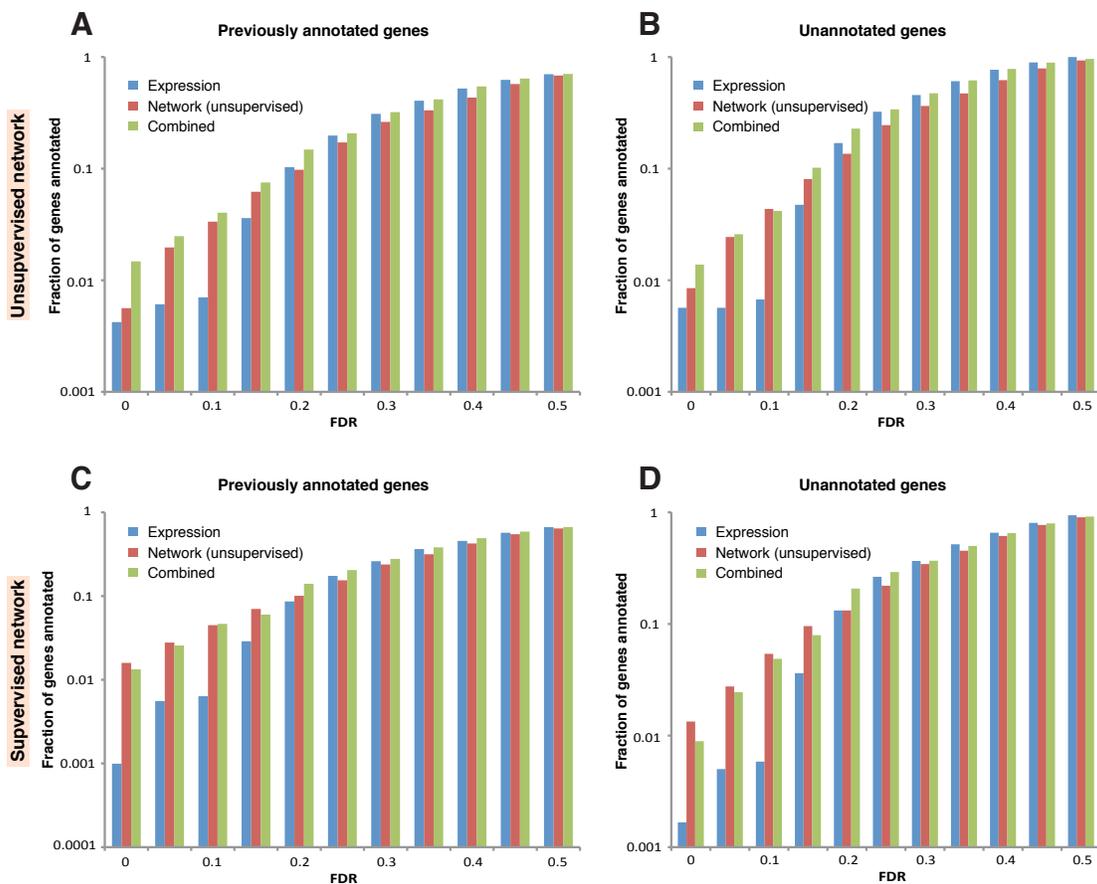


Figure S5: Prediction of novel functional annotations

The number of newly predicted GO process annotations at different FDR cutoffs for genes that are already annotated with other functions (A,C) and for genes that have no previous annotations (B,D). Predictions based on the unsupervised integrative network are shown in the first row (A,B), predictions based on the supervised integrative network are shown in the second row (C,D).

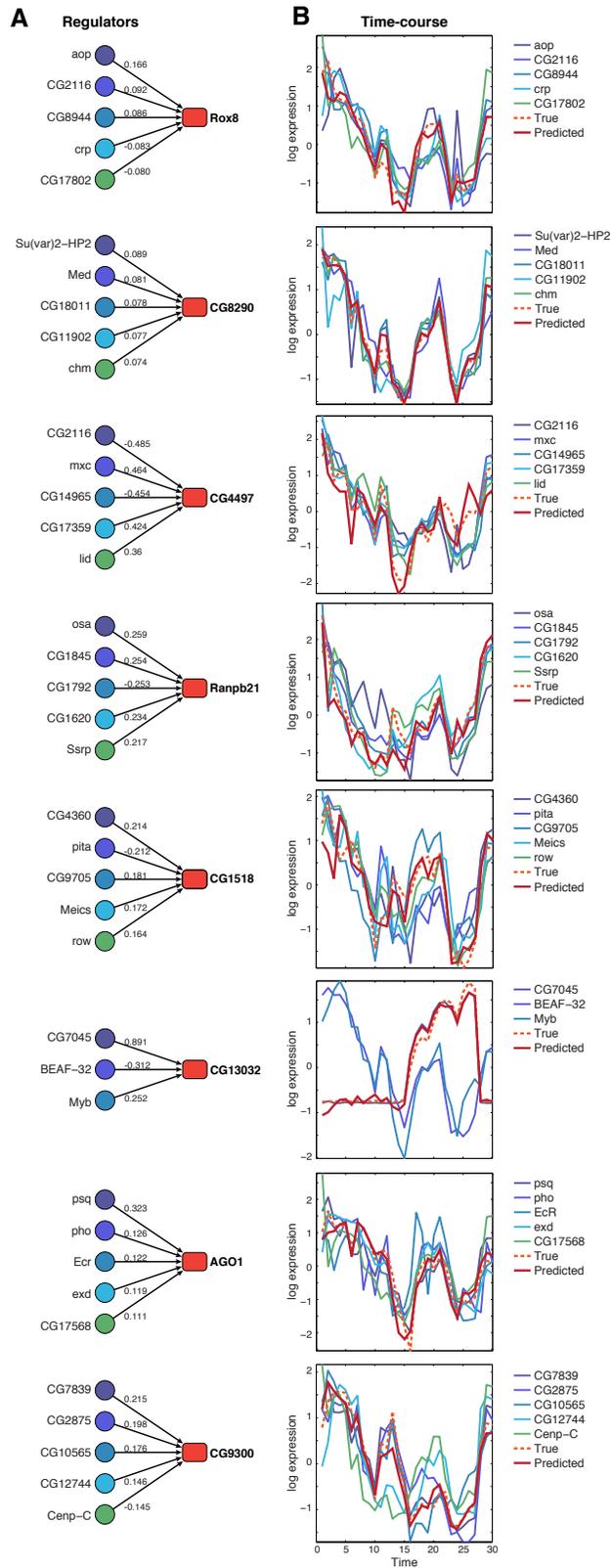


Figure S6: Expression prediction for eight representative target genes
 (A) The top 5 regulators and the magnitude of the regression coefficients (CG13032 has only three regulators). (B) True and predicted expression profiles for the target genes across the developmental time course. The expression profiles of the 5 regulators are also shown.

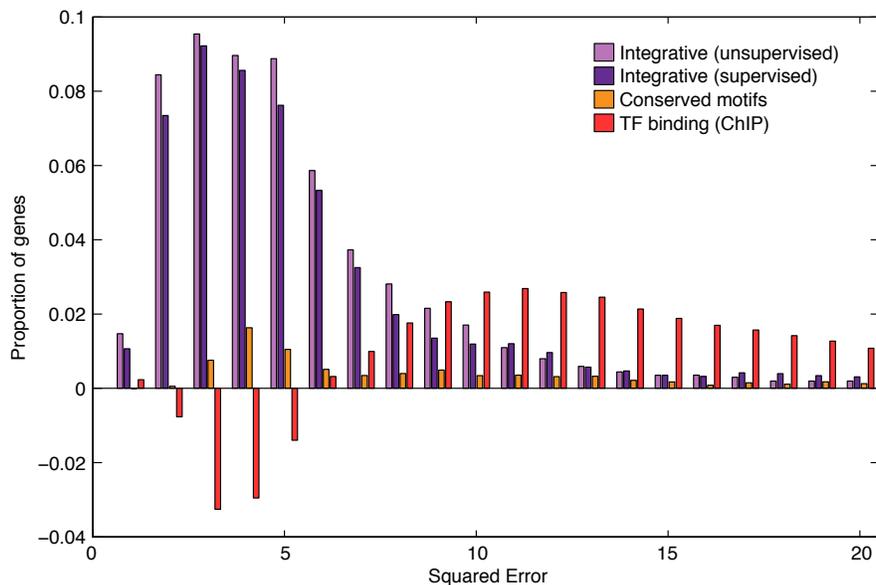


Figure S7: Squared errors for expression prediction using integrative and physical networks

Expression prediction was performed independently using the two integrative networks and the two physical networks based on TF binding (ChIP) and conserved motifs. For each of the four networks, we estimated the squared error between the true and predicted expression profiles, discretized the error into 20 bins, and estimated the proportion of genes with error less than the value corresponding to a particular bin. We did the same for 10 randomized networks, obtained an average proportion from randomized networks, and computed the difference between these two proportions (y-axis). In contrast to the two physical networks (ChIP and conserved motifs), the integrative networks show an excess of genes that are predicted accurately (low squared error) compared to predictions based on randomized networks.

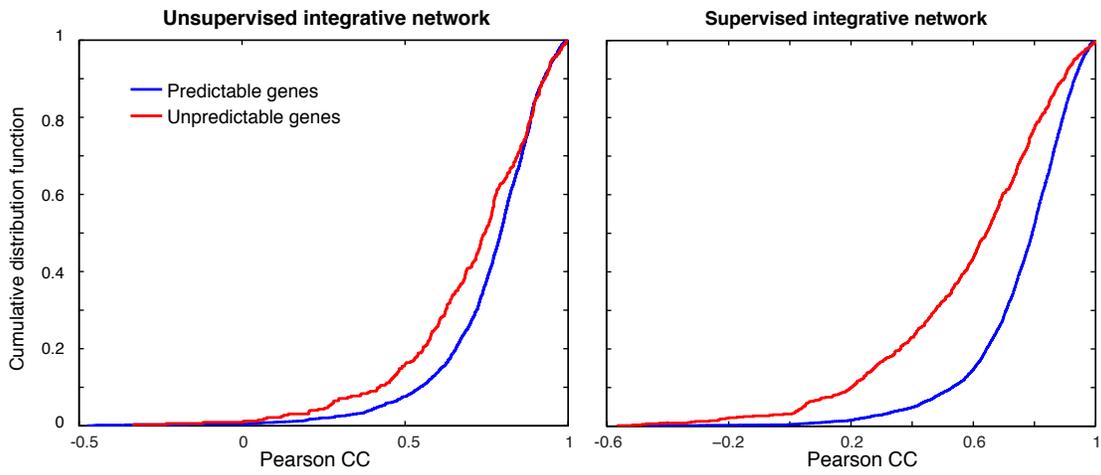


Figure S8: Correlation of predictable and unpredictable genes across two time courses

Distribution of Pearson's correlation of *predictable* and *unpredictable* genes across the two developmental time courses. Predictable genes showed a significantly higher correlation between the two datasets compared to the remaining genes (P-value < 1E-7), suggesting that predictable genes are also more reproducible.

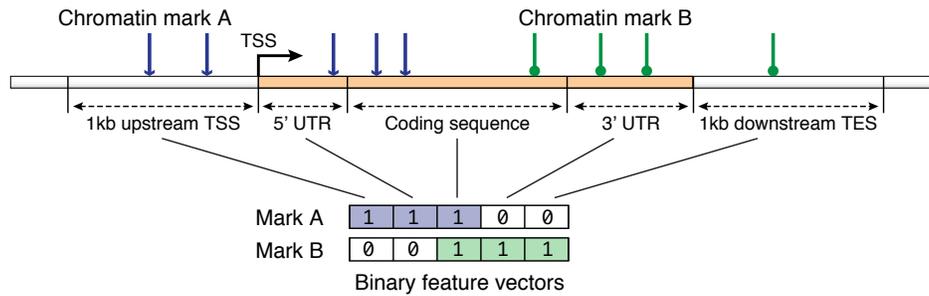


Figure S9: Definition of chromatin profiles

Chromatin profiles were defined by binning each gene into five regions and constructing binary vectors indicating the presence of each mark in these regions. The complete chromatin profile of a gene is obtained by concatenating the profiles of all marks for samples (time points or cell lines) into a single binary vector.

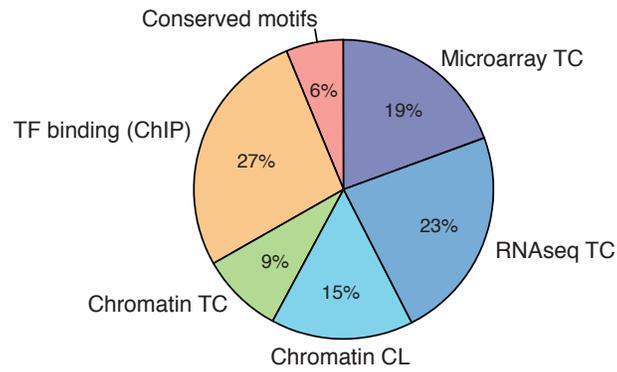


Figure S10: Relative importance of input features for the unsupervised integrative network

The pie chart shows the relative importance of the different input features in the *unsupervised* integrative network, analogous to Fig. 2C for the *supervised* network.

Supplementary Tables

Network	#TFs with expr.	#Targets with expr.	10-fold cross-validation in time course							Hold-out datasets in cell lines			
			#Pred. better than random	#Pred. worse than random	Enrichm. vs. random	Predictable genes		Unpredictable genes		Pred. genes (RNAseq)		Pred. genes (Microarray)	
Binding Motifs	72	10,248	1,707	1,642	1.04	656	6%	1,030	10%	296	45%		
REDfly	101	9,065	1,055	977	1.08	411	5%	214	2%	215	52%		
	74	82	13	3	3.91	12	15%	0	0%	4	33%		
Integrative (sup.)	541	8,433	1,988	787	2.53	1,502	18%	323	4%	992	66%	843	56%
Integrative (unsup.)	582	10,088	2,446	1,202	2.03	1,711	17%	629	6%	1,138	67%	1,118	65%

Table S1: Performance for gene expressions prediction

Shown are the number of TFs, targets, and predictable targets in the developmental time course and cell lines for the integrative, physical, and REDfly networks. The first two columns specify the number of TFs and targets that are in the expression datasets and used for predicting expression. The next four columns report cross-validation performance compared to randomized networks in the developmental time course. *#Predictions better/worse than random* specifies the average number of genes with significantly lower/higher error than randomized networks. The next two columns report the number of *predictable genes* and *unpredictable genes*, defined by genes that we can predict expression better (or worse) in 60% of the randomized networks. The last four columns report the number and percentage of predictable genes that remain predictable in the two cell-line hold-out datasets.

Supplementary Notes

S1 Network motifs for integrative and mixed regulatory networks

Results

We sought to reveal potentially meaningful network components by studying over-represented subgraphs known as network motifs¹ (see below for methodological details). Given the importance of microRNA (miRNA) regulation in feedback and repression in animal development, we complemented our transcriptional regulatory networks with post-transcriptional information on miRNA regulation. We included regulatory edges from miRNAs to their target genes from a previous study,³ where they were identified using evolutionary-conserved 7-mer miRNA seed motifs in three prime untranslated regions for 108 conserved microRNAs in 79 distinct families. Regulatory edges from TFs to miRNAs were identified based on binding of the TF in experimentally-inferred promoter regions for the precursor transcripts of miRNAs.^{2,3} The resulting combined pre- and post-transcriptional regulatory network contains 281 TF-to-miRNA edges and 599 miRNA-to-TF edges (out of a total of 6,141 miRNA-to-target edges).

We found 13 motifs in our mixed network that show significant fold enrichments compared to randomized networks (**Figure S4**), of which seven contain miRNA regulators. Four of the motifs in **Figure S4A** correspond to miRNA-mediated feedback loops (M2, M5, M6, M7), which may act as a noise buffer and prevent random switching events, as shown for *c-Myc*, *E2F1*, and *miR-17-20* in human.⁵ Three motifs constitute of only feed-forward loops (FFLs) (M1, M3, M4). Among these, the FFL in motif M4 is initiated by a miRNA and demonstrates co-operation between a TF and a miRNA to regulate a common target gene; motif M4 was also found to be enriched in human and mouse regulatory networks.⁴ Motif M3 corresponds to a double-FFL that can provide memory effect by converting transient signals of a common upstream regulator into states that persist after initial induction, which can be particularly important in developmental processes where decisions of cell fate occur^{1,7} (indeed, enriched GO terms for M3 target nodes include developmental processes, p-value < 10⁻¹⁰). Lastly, M1, which has two TFs cross-regulating each other and both targeting a common miRNA, was also significantly enriched in human regulatory networks.⁷

Overall, the inferred regulatory networks of *D. melanogaster* are enriched in similar network motifs as mammalian regulatory networks, suggestive of robustness, feedback, and co-operation between the transcriptional and post-transcriptional layers.

Methods

In order to identify building blocks of the regulatory networks, we considered all possible directed triads (i.e., connected three node sub-networks) of these networks. Among these triads, network motifs are the ones that appear significantly more compared to an ensemble of random networks with the same degree distribution as the original network (see *Network randomization* below). We used two different measures to assess the statistical significance of each triad, namely fold enrichment and z-score. We calculate the fold enrichment as the ratio of the percentage of each

triad in the original network to its average percentage in the randomized networks. Formally, for triad i the fold enrichment:

$$F_i = \frac{N * P_{or}}{\sum_{j=1}^N P_j} \quad (1)$$

where P_{or} is the percentage of triad i in original network, and P_1 to P_N are the percentages of triad i in N random networks. We compute these percentages and the z-scores for each triad using an existing network motif identification tool, FANMOD.⁶ We sort all possible triads according to their fold enrichment values calculated using 100 random networks generated by FANMOD. The top triads of this list are the most important network motifs in the original network. We also report the z-score, the number of appearances and the standard deviation of each network motif to better understand the portion of the original network that is explained by these conserved connectivity patterns (see [Fig. S4](#)).

We searched for network motifs in two different settings for each inferred network. First we focused on regulatory motifs in the integrative networks without miRNAs. We represented TFs and target genes with different colors to distinguish connectivity patterns between TFs from those including target genes. We generated 100 random networks using FANMOD in full enumeration mode, with three exchanges per edge and at most three attempts for each edge exchange. Among the 20 possible triads we picked the top six with highest fold enrichment in the integrative regulatory networks (Supplementary File S6-Sheet 2). The set of top six motifs is the same for both the unsupervised and supervised integrative network, and the motifs have a very similar ordering ([Figure S4B](#)).

To identify network motifs that include miRNAs, we colored miRNAs, TFs, and other genes with three different colors. Also, we distinguished between transcriptional and post-transcriptional regulations by assigning different colors to these different types of regulations. We used FANMOD to identify the motifs as described above. There are 66 possible triads with the node/edge coloring described above (Supplementary File S6-Sheet 1). We show the top seven motifs that contain at least one miRNA and that have the highest fold enrichment in the supervised integrative network ([Figure S4A](#)).

Network randomization. The random networks are generated from the original network by a series of edge switching operations (e.g. if there is an edge from a to b and c to d then an edge exchange operation changes these to an edge from a to d plus an edge from c to b). Edges are only exchanged if their endpoints have the same node type such as transcription factors or target genes. When randomizing the network, the edges are exchanged one after the other. We set the FANMOD options *exchanges per edge* and *exchange attempts* both to three, which means all the edges in the network are traversed three times and at each time at most three attempts were made to exchange an edge with another edge.

References

- [1] U. Alon. Network motifs: theory and experimental approaches. *Nat Rev Genet*, 8:450–461, 2007.
- [2] B. Graveley, A. Brooks, J. Carlson, M. Duff, J. Landolin, and et al. The developmental transcriptome of drosophila melanogaster. *Nature*, 471:473–479, 2010.
- [3] S. Roy, J. Ernst, P. Kharchenko, P. Kheradpour, N. Negre, and et al. Identification of functional elements and regulatory circuits by drosophila modencode. *Science*, 330:1787–97, 2010.
- [4] R. Shalgi, D. Lieber, M. Oren, and Y. Pilpel. Global and local architecture of the mammalian microrna-transcription factor regulatory network. *PLoS Comput Biol*, 3:e131, 2007.
- [5] J. Tsang, J. Zhu, and A. van Oudenaarden. Microrna-mediated feedback and feedforward loops are recurrent network motifs in mammals. *Mol Cell*, 26:753–767, 2007.
- [6] S. Wernicke and F. Rasche. Fanmod: a tool for fast network motif detection. *Bioinformatics*, 22:1152–1153, 2006.
- [7] X. Yu, J. Lin, D. Zack, J. Mendell, and J. Qian. Analysis of regulatory network topology reveals functionally distinct classes of micrnas. *Nucl Acids Res*, 36:6494–6503, 2008.