

## Supporting Online Material for

### **Identification of Functional Elements and Regulatory Circuits by *Drosophila* modENCODE**

The modENCODE Consortium, Sushmita Roy, Jason Ernst, Peter V. Kharchenko, Pouya Kheradpour, Nicolas Negre, Matthew L. Eaton, Jane M. Landolin, Christopher A. Bristow, Lijia Ma, Michael F. Lin, Stefan Washietl, Bradley I. Arshinoff, Ferhat Ay, Patrick E. Meyer, Nicolas Robine, Nicole L. Washington, Luisa Di Stefano, Eugene Berezikov, Christopher D. Brown, Rogerio Candeias, Joseph W. Carlson, Adrian Carr, Irwin Jungreis, Daniel Marbach, Rachel Sealfon, Michael Y. Tolstorukov, Sebastian Will, Artyom Alekseyenko, Carlo Artieri, Benjamin W. Booth, Angela N. Brooks, Qi Dai, Carrie A. Davis, Michael O. Duff, Xin Feng, Andrey Gorchakov, Tingting Gu, Jorja G. Henikoff, Philipp Kapranov, Renhua Li, Heather MacAlpine, John Malone, Aki Minoda, Jared Nordman, Katsutomo Okamura, Marc Perry, Sara Powell, Nicole C. Riddle, Akiko Sakai, Anastasia Samsonova, Jeremy E. Sandler, Yuri Schwartz, Noa Sher, Rebecca Spokony, David Sturgill, Marijke van Baren, Kenneth H. Wan, Li Yang, Charles Yu, Elise Feingold, Peter Good, Mark Guyer, Rebecca Lowdon, Kami Ahmad, Justen Andrews, Bonnie Berger, Steven E. Brenner, Michael R. Brent, Lucy Cherbas, Sarah C. R. Elgin, Thomas R. Gingeras, Robert Grossman, Roger A. Hoskins, Thomas C. Kaufman, William Kent, Mitzi Kuroda, Terry Orr-Weaver, Norbert Perrimon, Vincenzo Pirrotta, James W. Posakony, Bing Ren, Steven Russell, Peter Cherbas, Brenton R. Graveley, Suzanna Lewis, Gos Micklem, Brian Oliver, Peter J. Park, Susan E. Celniker,\* Steven Henikoff,\* Gary H. Karpen,\* Eric C. Lai,\* David M. MacAlpine,\* Lincoln D. Stein,\* Kevin P. White,\* Manolis Kellis\*

\*To whom correspondence should be addressed. E-mail: manoli@mit.edu (M.K.) (integrative analysis); celniker@fruitfly.org (S.E.C.) (transcripts); karpen@fruitfly.org (G.H.K.) (chromatin); kpwhite@uchicago.edu (K.P.W.) (transcription factors); david.macalpine@duke.edu (D.M.M.) (replication); laie@mskcc.org (E.C.L.) (small RNAs); steveh@fhcrc.org (S.H.) (nucleosomes); lincoln.stein@gmail.com (L.D.S.) (data availability)

Published 22 December 2010 on *Science* Express  
DOI: 10.1126/science.1198374

#### This PDF file includes:

Materials and Methods  
SOM Text  
Figs. S1 to S23  
Tables S1 to S8

## References

**Other Supporting Online Material for this manuscript includes the following:** (available at [www.modencode.org/publications/integrative\\_fly\\_2010/](http://www.modencode.org/publications/integrative_fly_2010/))

Data Sets S1 to S17

## ***Drosophila* modENCODE - Supporting Online Material**

This document includes five parts:

- Supporting methods and text
- Supporting figure legends
- Supporting table legends
- List of supporting data sets (data sets are available at  
[http://www.modencode.org/publications/integrative\\_fly\\_2010/](http://www.modencode.org/publications/integrative_fly_2010/))
- Supporting references

### **Supporting methods and text**

#### **Data Production.**

The data set includes a total of nearly 3 billion high-throughput whole-transcriptome next-generation sequencing (RNA-seq) short poly(A)+ RNA-seq (Illumina) (2.3B Developmental Stage RNA-seq + 126M 4-cell line RNA-Seq + 69M other species RNA-Seq + 211M small ncRNAs), long poly(A)+ RNA-seq (454) (1.75M RNA-seq), and total RNA-seq (SOLiD) (92M RNA-seq) sequencing platform reads for RNA-seq and ChIP-seq experiments. Together, these were used for informing gene structure, profiling RNA expression, identifying copy number variations, and ascertaining chromatin binding sites, and comprise more than 200 billion bases (more than 1700 genome equivalents). The raw data set also includes more than 300 replicated genome tiling arrays used for chromatin IP (ChIP-chip), RNA expression profiling, DNA replication timing, nucleosome density, and comparative genomic hybridization (CGH). The ChIP experiments assayed a total of 103 transcription factors (TFs) and other chromatin-binding proteins and 23 histone modifications and variants.

#### **Data Processing and Community Access.**

All data sets were released to the community upon verification and vetting by the Data Collection Center (DCC). Each data set is accompanied by structured metadata, submitted to the DCC at the same time as the data set, and subjected to hand-curation prior to automated processing. The metadata describe reagents used in the experiment, experimental protocols (including bioinformatics protocols), and the experimental design (including biological and technical replicates). All metadata, interpreted results, and raw data can be searched, displayed, and downloaded via modMine (<http://intermine.modencode.org>), the DCC's interface to the modENCODE corpus.

The primary site for access to verified modENCODE data is <http://www.modencode.org>, where readers can find complete sets of both raw and analyzed data in downloadable, browsable, and searchable formats; the data sets are also available in a cloud computing environment via the Bionimbus system. Arrays and sequencing reads are archived at the Short Read Archive (SRA; <http://www.ncbi.nlm.nih.gov/sra>) and Gene Expression Omnibus (GEO;

<http://www.ncbi.nlm.nih.gov/geo/>) (S1). Results from modENCODE are being incorporated into the community database FlyBase (<http://www.flybase.org>) (S2). These data sets are summarized in Table S1, and additional files presented in this paper are available at [http://www.modencode.org/publications/integrative\\_fly\\_2010/](http://www.modencode.org/publications/integrative_fly_2010/).

Most assays were conducted in the sequenced strain *y; cn bw sp* (S3), with some additional RNA expression and copy number variation studies conducted in Oregon-R, Canton-S, and others. Four commonly used cell lines [ML-DmBG3 (BG3), Kc167, CME W1 cl.8+ (Cl.8+), and S2-DRSC] (Table S2) were used for analysis of transcription, transcription factor (TF) binding sites, and histone mark distributions. Each element was tested in a variety of developmental time points, with very granular dissection of RNA expression (30 time points), and less granular for TF and histone studies (typically 8 time points). To date, only short ncRNAs have been assayed in isolated tissues, though future experiments are planned to measure patterns total RNA expression in isolated tissues.

#### **Methods for assessing data reproducibility and quality.**

For embryo ChIP-seq assays, the data standards were as follows: we generated at least two biological replicates, where at least 80% of the top 40% of peaks called on the shorter of the two peak lists is present in the longer of the two. For embryo ChIP-chip assays, we conducted at least three biological replicates, and insisted on a probe-level Pearson correlation of at least 0.7 between each of the three biological replicates. In some cases, ChIP data that did not meet these standards has been included in this study at the discretion of data producers on the grounds that, although of lesser quality, the data still include valuable biological information. For RNA-seq assays, we conducted extensive browsing of thousands of genomic locations.

For cell line ChIP-chip assays, we generated at least two biological replicates. The data set was deemed consistent if the enriched regions from the individual biological replicates had at least 75% reciprocal overlap, or if at least 80% of the top 40% of the regions identified by each replicate were identified by the other replicate (before comparison the replicates were size-equalized by increasing the significance threshold for a replicate with more enriched sequence).

#### **Methods for coding and non-coding transcripts.**

Transcript models were included if they were detected with  $\text{FPKM} \geq 1$  and low-confidence  $\text{FPKM} > 0$  in polyA(+) RNA-seq samples. Additional transcript models were included by manual curation of cDNAs and genomic regions detected by total RNA-seq samples. While these transcripts were detected in at least one developmental time point, the connectivity and precise structures such as exon start and stop positions remain to be validated by long-sequencing methods such as cDNA screening.

For validation, we targeted 3887 new predicted transcript fragments derived from the sequence reads using sim4 (S4), SPA (S5), TopHat (S6), and Cufflinks (S7). For all transcript fragments at least one exon did not overlap a FB5.12 exon, but a majority

showed evidence of being part of an alternate transcript of an annotated gene. We designed inverse PCR primers and performed directed cDNA screening and captured and sequenced 970 cDNA clones. The gene expression data also allowed us to select the appropriate library to screen 2,114 FB5.12 gene models with no validating cDNA. We sequenced 1,021 cDNAs for 757 genes that had eluded prior detection. We validated an additional 1,937 genes and transcripts using ESTs. Finally, the cDNA sequences support 136 FB5.12 gene merges and validate 485 of 1938 new predicted genes (25%). A total of 1,498 genes were not detected by RNA sequencing. This set is highly enriched in rRNAs, histones, olfactory and gustatory receptor genes, as well as computationally predicted genes (data not shown).

To provide additional cross-validation of the predicted exons, we have compared a subset of exons active in S2 cells ( $\text{FPKM} > 0$  with 95% confidence, 21071 exons) with putative promoter and transcribed regions predicted based on the chromatin signatures (see main manuscript). We find that 81% of the novel S2 exons is overlapped (by at least 10% of its length) by the predicted transcribed regions, 44% by promoter predictions (score  $\geq 0.8$ ). Overall, 89% of novel S2 exons are covered by either putative promoter or transcribed regions.

### **Coding potential of novel genes.**

We used a comparative genomics approach to assess coding potential in the novel genes. We first extracted alignments of each transcript model from MULTIZ whole-genome alignments of 12 *Drosophila* species (S8, S9), and then used PhyloCSF (S10) to search for protein-coding evolutionary signatures in each alignment. For each transcript alignment, we ran PhyloCSF in two modes: first, searching for start-to-stop ORFs on the annotated strand (`--frames=3 --orf=ATGStop`), and second, searching for high-scoring windows ending with a stop codon on either strand (`- frames=6 --orf=StopStop3`). In both cases, we used a minimum ORF/window length of 20 codons.

We selected transcripts with a PhyloCSF window score of at least 80 decibans, a cutoff that accounts for the multiple testing inherent in this analysis for further examination. We excluded as likely antisense transcripts any with a higher window score on the opposite strand than ORF or window score on the annotated strand (the vast majority of these are indeed antisense to known coding genes). We then manually examined the alignments of all transcripts with higher ORF score than window score to select those likely to represent complete open reading frames ("Complete ORF") based on conservation of the start and stop sites. The remaining examples, as well as those with higher window score than ORF score, were categorized as "Incomplete ORF". Lastly, to exclude the possibility that substantial numbers of coding genes were missed by PhyloCSF, we analyzed all transcripts lacking high PhyloCSF scores using BLASTX homology searches against Swiss-Prot. Only 15 of the 1,800 remaining genes (<1%) have a transcript with any BLAST hit on its annotated strand with an E-value of 10-6 or less (`-p blastx -S 1 -e 1e-6`). It is therefore unlikely that these remaining transcripts include many functional coding genes, since the vast majority show neither similarity to known protein sequences, nor evidence of selection on a protein-coding sequence across closely related species.

## **Methods and details for predicted secondary structures.**

Since many functional noncoding RNAs (ncRNAs) rely on a defined secondary structure, we searched for significant RNA structures in the newly discovered transcribed regions. We also analyzed transcripts listed in Flybase that were annotated as "ncRNA" but do not belong to well known classes of structural RNAs and for which in most cases their function is unknown.

Although functional RNA structures are in most cases more stable than random sequences, thermodynamic stability alone is generally not statistically significant enough to detect functional RNAs. However, comparative approaches that incorporate information from homologous sequences can detect functional RNA structures with reasonable accuracy. We therefore used the RNAz algorithm that combines evidence from thermodynamics and evolutionary structure conservation (S11) and applied it to alignments of 12 drosophilid genomes.

We analyzed 1740 novel transcripts that did not show coding potential according to our evolutionary analysis. In addition we added 140 noncoding transcripts from Flybase. This list included all transcripts from FB annotated as "noncoding" with transcripts removed that were explicitly annotated as miRNAs, snoRNAs, tRNAs, rRNAs, and all other small noncoding RNAs (smnRNAs).

For each region, we extracted genome wide Multiz alignments of 12 drosophilid genomes. We then scanned the alignments in overlapping windows of size 120 and slide 40. Sequences with more than 25% gaps and alignment blocks with fewer than 5 aligned sequences were discarded. The remaining windows were analyzed using RNAz version 2.0pre [[www.tbi.univie.ac.at/~wash/RNAz](http://www.tbi.univie.ac.at/~wash/RNAz) (S12)]. For the CuffLinks predictions without strand information, we scanned both strands separately. We considered structure predictions with RNAz classification score  $P>0.5$ . Alignments containing a predicted RNA structure were re-aligned using LocARNA version 1.5.2b (S13) (<http://www.bioinf.uni-freiburg.de/Software/LocARNA/>) and mutational patterns were color-annotated using tools from the RNAz package.

Our pipeline predicts significant structures in 11 of the 140 Flybase transcripts. Six could be assigned to well known structures such as the RNase P (FBtr0085775), SRP RNA (FBtr0081624, FBtr0302398) and miRNAs (FBtr0083362, FBtr0085391) and snoRNAs (FBtr0300454) not explicitly annotated as such in Flybase. These RNAs were not excluded from the initial set but—other than serving as a positive control—are of limited interest here. We only could find novel structures in five other long ncRNAs of Flybase: 3 transcripts do not have functional annotation (FBtr0073592, FBtr0073625, FBtr0080943), 2 structures were found in the well studied roX2 RNA (FBtr0073514) and the HSR $\omega$  (FBtr0084059) transcripts and are therefore of special interest (fig. S1).

roX2 is a major determinant of sex chromosome dosage compensation in flies. Our analysis recovered a previously described region containing a conserved sequence motif predicted to form a stem loop and shown to be sufficient for the MSL complex recruitment and H4-Lys16 acetylation functions of roX2 (S14). Our consensus

structure model shows a striking number of compensatory mutations and thus unambiguously confirms this stem-loop. We also found a second stem-loop of similar size and equally strong evolutionarily support immediately downstream of the known stem-loop, as well as a somewhat weaker structural signal in the 3' end of the roX2 transcript (fig. S1).

We also found a strong structure signal within the HSR $\omega$ , a ncRNA involved in the heat-shock response (S15). This ncRNA is functionally analogous to the mammalian SatIII ncRNA, that has roles in protein folding and post-transcriptional regulation (S16). Although to our knowledge no structure had been previously reported for this ncRNA, we find a highly-structured region in the 5'-end of the transcript that contains at least two conserved stem-loops and an exceptional number of compensatory changes (fig. S1), suggesting an important role in the function of HSR $\omega$ .

Of the 1740 newly identified transcripts without any coding potential, we found 132 (7.6%) with at least one predicted structure. To assess the significance of these results, we repeated the analysis on randomized alignments. We used a shuffling algorithm described previously (S17) that removes any correlations from secondary structure while keeping crucial alignment characteristics like conservation and GC-content unchanged. 3.3% of the randomized regions was predicted contain at least one structure, we thus observe a 2.3-fold enrichment of structures in the novel transcripts over the random control. To assess the sensitivity of the approach and to rule out the possibility that the relatively small fraction of detected structured RNAs is due to overly stringent cut-offs, we used the same parameters on alignments of known RNA structures that are at least conserved in five species. We observed sensitivities of 84.7% for tRNAs, 75.8% for spliceosomal RNAs (snRNAs), 24.7% for small nucleolar (snoRNAs) RNAs, and 93.0% for miRNA precursors. These results are consistent with previous evaluations of the RNAz algorithm and show that the protocol reliably recovers conserved structures.

#### **Annotation of Argonaute-mediated small RNAs.**

To annotate various classes of small RNAs relevant to Argonaute pathways, we aggregated 187 data sets of ~18-28 nt RNAs, including 76 data sets produced for the modENCODE project and additional published data available from NCBI GEO. We annotated new canonical miRNAs (S18) and mirtrons (S19). Both efforts generally relied upon the recovery of ~21-24 small RNA reads with dominant 5' ends that comprised a duplex with 3' overhangs within a hairpin precursor. In rare cases we accepted hairpin loci with only one cloned arm as genuine miRNAs, if the reads exhibited a preferred 5' end and accumulated in AGO-IP library data. The annotation of primary piRNA production from mRNAs, and particularly from 3' UTRs, was reported earlier (S20).

To annotate novel cis-NAT-siRNA loci, we used a strategy described by Okamura (S21). We collected bidirectionally transcribed exonic regions >20 nt from the MB8 gene models, defining 2532 genomic regions representing 389 kB. We then used Bowtie (S6, S7) to map perfectly matching reads from the collected small RNA data. Finally, we selected those regions with substantial enrichment for 21 nt reads (2/3 of all reads

mapped to a given region) and a minimum read depth (>50 nt 21 nt reads) as cis- NAT-siRNA loci. These comprised 245 loci covering 48kb and 195303 reads, which divide into 4 5' UTR cis-NAT regions; 4 "ambiguous" cis-NAT regions including overlapping internal exons, and 237 3' UTR cis- NAT regions (Table S4).

### **Chromatin signatures for different classes of functional elements in S2 cells.**

To calculate the average enrichment patterns associated with expressed genes, individual enrichment profiles (log<sub>2</sub> intensity ratio) were calculated for a set of expressed genes (FPKM ≥ 0.65) that are least 1Kbp long (using 150 spatial bins, see scaling diagram in Fig. 3A). The profiles of individual genes were then averaged to obtain the final enrichment profile for a given mark. The average profiles associated with the replication origins, insulator binding positions and TF binding hotspots were similarly calculated by averaging the ±2Kbp profiles (100 spatial bins) of individual sites. The TSS-proximal and distal positions were classified based on a 2Kbp proximity to an annotated TSS.

### **Integrative prediction of active promoters in S2 cells.**

We implemented a supervised learning approach to distinguish active promoters from inactive promoters based on their chromatin modification, transcription factor binding, and nucleosome physical property profiles in S2 cells. We focused on transcripts with validated promoters (evidence from two or more of the CAGE, RACE, RE-EST assays) and used RNAseq data to define active and inactive promoters based on a FPKM threshold, only transcripts with a FPKM value greater than 0.5 were considered actively transcribed.

The features used to distinguish active from inactive promoters correspond to 18 chromatin modifications (H3K9ac, H3K18ac, H3K23ac, H3K27ac, H4K16ac, H4K5ac, H4K8ac, H4acTetra, H3K4me1, H3K4me2, H3K4me3, H3K9me2, H3K9me3, H3K27me3, H3K36me3, H3K79me1, H3K79me2, H2B.ubiq), binding site profiles for PolII and ORC2, and physical properties of components of nucleosomes (salt fractionation data, CATCH-IT data, mononucleosome profiles) profiled by a combination of ChIP-chip and ChIP-seq assays. A 1Kb region was centered on each validated promoter and split into 100bp bins. The ChIP-chip/seq values were assigned to each bin producing a vector for each data set at each valid promoter. If multiple values were present across a bin, the values were weighted by the fraction of overlapping signal and averaged. We required that at least 5 of the 10 bins had values and filtered out regions that did not have consistent signal across all the data sets.

We implemented a strategy developed for integrating chromatin signatures into promoter prediction models in order to represent the values as features for the classifier (S22). For the set of positives, we first computed an average vector for each mark and then for each individual promoter and calculated the Pearson correlation and the dot product with the average of the positives. The Pearson correlation was selected to represent the shape of the signal and the dot product was used to represent the intensity. We labeled the promoters as positive if they are adjacent to an active transcript and negative if they are adjacent to an inactive transcript. This resulted in

5803 active promoters and 1890 inactive promoters. We also select a set of 7693 random regions that serve as additional negative examples used in training that are matched based on the chromosomal distribution of valid promoter in the heterochromatic and euchromatin regions of the genome.

A training set was assembled from 70% of the labeled regions, leaving the remaining 30% for model evaluation. Chromosomal distribution and FPKM values were balanced across the training and test sets. We learned a classifier for the S2 cell data based on the training set using a logistic regression model implemented in the WEKA machine learning software suite (S23). We then used overlapping 1000 bp windows (100bp shift) to generate features across the genome and ran the model to identify regions that have similar experimental profiles as the active valid promoters on both strands of the genome. For this analysis chrU and chrM were excluded.

We evaluated the performance of the model using a distance-based evaluation strategy developed for promoter predictions (S24), where a promoter prediction that falls within 500 bp of a valid active promoter is classified as correct and each active promoter that has a promoter prediction is classified as discovered. The total set of predictions is those predictions that fall within the genic regions of the genes in the test set. For this approach, all intergenic regions are ignored and bins that overlap a known TSS within the gene body are also filtered. The distribution of annotated promoter size indicates that 75% of the promoters are less than 203 bps and half have some spacing between the promoter annotation and the TSS. For these promoters that do not overlap the TSS, 75% have an 80bp or less spacing and thus we added a 500bp margin to the TSS, when removing TSS overlapping bins (fig. S4A, B). Precision is defined as the number of correct predictions divided by the total number of predictions and the recall is the number of discovered active promoters normalized by the total number of active promoters.

We then made a comparison with the chromatin state models described below. To generate a precision-recall curve for the chromatin-states, we start with the most enriched state and subsequently add additional states. For the 9-state model the state1 is the most enriched followed by state 3. For the 30-state model, the following order of states was used: 1, 5, 6, 3. We find that based on the evaluation protocol described above, the 9-state model performs slightly better than the supervised approach (fig. S4C), however, this evaluation protocol does not account for the size of the promoter predictions and larger predictions could result in a better performance.

An alternative approach for evaluation is to consider the enrichment of a promoter annotation set within the predicted regions. Here the enrichment is determined to be the number of bps that are shared between the annotation set and the predicted regions normalize by the total size of the background. This quantity is the normalized by the product of the number of bps in the annotation set and the number of bp in the prediction set. We first investigated the enrichment for the active valid promoters of the test set, using the genic regions of the test genes as the background, and find that the supervised predictions show the strongest enrichment, followed by the 30-state model and then the 9-state model (fig. S4D). We also used an annotation set that uses the

TSS for all active genes in S2 cells by adding a 100 bp window around the start position and for this analysis the entire genome is the background. Here we find a similar trend with the supervised predictions showing the best performance (fig. S4E).

Using mRNA expression to define the annotation set may be problematic as we're selecting for poly-A transcripts in the RNAseq prep, non-poly-A transcripts could be missed. These promoters could be associated with non-coding non-poly-A transcripts. In addition, promoters bound by PolII that are paused may not show transcriptional support by the current RNAseq protocol. Therefore, in order to provide additional validation to the promoter predictions, we incorporated a published data set for short capped RNA that are derived from PolII carried out in S2 cells (S25). For this annotation set, we added a 100 bp to each region of the genome with more than one unique read mapped from the short capped RNA library and find the supervised approach shows the strongest enrichment for short RNA reads (fig. S4F).

We selected an operating point with 93.7% sensitivity at 21.5% FDR and used promoter activity predictions to validate the 52914 new or modified exons and new promoters discovered by transcriptional profiling efforts. There are 7405 first exons that are associated with transcripts expressed with FPKM values of 0.5 or greater and 5459 (73.1%) have an active promoter prediction directly overlapping the TSS. This reduced sensitivity likely reflects the observation that there is a difference between the location of the annotated promoter and the annotated TSS. We find promoter predictions overlapping 310 of the 428 (72.4%) new promoters adjacent (with 300bp) of an active transcript in S2 cells. Furthermore, 82% (264 of 321) of the promoter annotations that have support from 2 of the 3 promoter validation assays were verified as active and only 43% (46/107) of the promoters with primarily RACE support are annotated as active. Overall, we find 2203 predicted promoters that are not within 500 bp of an annotated start site or overlap a valid promoter. For this set, 67% have support from the short capped RNA data set (at least one unique read) indicating that many of these new promoter locations are indeed valid.

The putative transcribed regions were determined as continuous regions of significant enrichment for H2B-ubiquitination or H3K36me3 at least 3 kbp in size, positioned at least 2 kbp away from any annotated genes. The enriched regions were identified using a 3-state HMM model with fixed emission signal distributions ( $M=0.5$  for enriched state, 0 for neutral, -1 for depleted). See Figure S3 for characterization of predictive power.

#### **De novo discovery of chromatin states.**

The derivation and analysis of the 9-state continuous and 30-state discrete chromatin models is described in Kharchenko *et al.* (S26). Both models utilized multivariate Hidden Markov Models (HMM) to jointly model the chromatin patterns observed across BG3 and S2 genomes, by effectively treating the data from the two cell types as if they were different sets of chromosomes in the same cell type.

The discrete (30-state) model was derived using a method described in (S27). Briefly, the modification data was binarized locally at a 200-bp resolution, based on the local signal level at a 5% FDR compared to 10-genome wide randomizations. For H3K23ac,

H1, and H4 depletion significant depletion was used as the input feature instead of enrichment. Model learning and training was applied as described in (S27), except in the nested initialization procedure models with up to sixty states were learned, and states from the highest scoring based on the BIC score was pruned based on euclidean distance, and from inspection a 30-state model was selected to provide a desired level of biological interpretation. State enrichments for elements were based on the full posterior distribution over state assignments. Hard assignments based on the maximum posterior were also computed. Fold enrichments for elements were computed based on enrichment at the 200-bp bin overlap resolution. Expression level of chromatin states were computed based on the average expression of 200bp intervals of the chromatin state that overlapped genes. The expression level of a 200bp interval was determined as the average expression of overlapping genes after adding a pseudo count of one and applying a log base 10 transformation. The nucleosome solubility and turnover values associated with states was the average signal in the state.

The intensity-based (9-state) model modeled the average log intensity ratio within the 200bp bins using multivariate normal emission distributions. Briefly, the 200bp bin average M values were quantilenormalized between the two cell lines, and an HMM model was determined on the joint data using Baum-Welch algorithm, using 30 seeding configurations determined from K-means clustering. States with minor intensity variations (Euclidian distance of mean emission values < 0.15) were merged.

### **TF binding reveals HOT regions.**

In total of 25 factors from 41 data sets generated by modENCODE project and 16 factors from BDTNP (Berkeley Drosophila Transcription Network Project) (S28, S29) (25, 26) were combined to investigate the co-localization of transcript factors. Data sets generated for the same factor were merged, and in total 41 transcription factors' binding sites were used for the following analysis. For each pair of TFs, a Fisher's exact test was employed to present the significant of co-localization. And – log10 transformed p\_values were used to make a heatmap of all TFs co-localization. To define the HOT regions we performed a Gaussian kernel density estimation across the genome with a bandwidth of 300 bp, using the centers of each of the TF binding peaks as points. We then scanned this density for peaks, and denoted each peak a HOT region. To determine the complexity of the HOT region, we summed the Gaussian kernalized distance from the peak to each transcription factor that contributed at least 0.1 to its strength. The reported window around eat HOT peak was derived by finding the maximum distance (in bp) from the HOT peak to a contributing TF, and then adding 150 bp (one half of the bandwidth). Each window is centered on the HOT peak. In those analyses in this paper which required a complexity cutoff, those binding sites with TFBS complexity more than 8 were defined as HOT regions.

### **Predicting combinations of TF binding from combinations of sequence motif.**

In order to investigate whether combinations of motifs improve the prediction of transcription factor binding sites, we constructed a logistic regression classifier. The classifier was trained on transcription factor binding sites, and classifier accuracy was

assessed using ten-fold cross-validation. We studied the combinations of features that were most predictive of transcription factor binding. We reasoned that if combinations of motifs increased predictive power, than accuracy of the predictions should be improved by including the most enriched (cooperative) or the most depleted (inhibitory) motifs. This was tested by comparing performance using: the single most enriched feature for each transcription factor, the five most enriched motifs, the ten most enriched motifs, the five most enriched as well as the five most depleted motifs (Fig. 5D). We found that including depleted as well as enriched motifs improved the power of the classifier, suggesting that this approach provides insight into antagonistic as well as synergistic relationships among transcription factors.

#### **Motif enrichments at varying complexity classes.**

To determine if the recruitment of factors to the HOT regions is specific, we examined the enrichment of known motifs for the corresponding factors in the regions bound at different complexity levels (Fig. 5D). When multiple known motifs exist for a given factor, the one whose enrichment is highest in the bound regions compared to intergenic regions is taken. Shuffled motifs are used to generate a baseline level of enrichment and correct for potential biases in composition. Using the number of motif instances inside the lower complexity and higher complexity regions and comparing that to the corresponding counts for the shuffled motifs, we compute a bias corrected log odds ratio. We then use the standard error of this statistic to compute a 95% ( $Z=1.96$ ) confidence interval around each value and use the most conservative (i.e. closest to 0) enrichment value. Known DNA binding motifs tend to be more enriched in low complexity regions. This suggests that either low complexity regions are bound more strongly, or that the binding in low complexity regions is greater due to binding to the known motif. It might also suggest that the binding of factors to HOT regions is less specific and raises the question of their function. Indeed, we found that enhancers [as defined by the CAD database (S30)] are more represented in low complexity categories (S31), suggesting that HOT regions might have a different function than tight transcriptional regulation. We also examined the enrichment of the corresponding known motif in each of the experimental regions at varying levels of complexity (Fig. 5D). We do not find significantly different levels of enrichment in low vs. high complexity peaks for most factors. However, for the factors where we do see a difference, low complexity peaks tend to be more enriched. This supports the hypothesis that the low complexity peaks are more driven by sequence-specific complementarity. Conversely, the few factors for which higher enrichment is seen in high complexity regions may be due to these factors being responsible for establishing these high complexity regions.

#### **Methods for estimating genome coverage.**

Estimates of genome coverage where based on the resolution available from the data sets; these can be expected to evolve as new and higher resolution experiments are performed. The portion of the genome under evolutionary constrain was identified using a 14-way alignment of *D. melanogaster* with 11 additional *Drosophila* species (*D. simulans*, *D. sechellia*, *D. yakuba*, *D. erecta*, *D. ananassae*, *D. pseudoobscura*, *D. persimilis*, *D. willistoni*, *D. virilis*, *D. mojavensis*, *D. grimshawi*), mosquito (*A. gambiae*),

honeybee (*A. mellifera*) and beetle (*T. castaneum*) and applying the phastCons algorithm as described in the text. For Fig. 6A, data sets were ordered such that those most incontrovertibly considered as annotation are accounted for first (e.g. exons) and those which are more disputably considered an annotation are accounted for last (e.g. introns). We identified large blocks of repeats and other low complexity sequences by taking all elements determined by RepeatMasker (A.F.A. Smit, R. Hubley & P. Green, unpublished) with a score of at least 2000. For Fig. 6B data sets were grouped into three broad domains: transcribed elements, bound regulators and chromatin domains. Transcribed elements consist of coding exons, Ago smRNAs (including piRNAs, cis-NAT-siRNAs, hpRNAs, miRNA and mirtrons), UTRs (3'/5'), non-coding RNAs (including pri-mirs) and introns. Protein isoforms observed through RNA-seq positively contribute to transcribed element density. Bound regulators consist of transcription factors (48 data sets), insulators (10 data sets), and other bound proteins (14 BG3 data sets and 22 S2 data sets). Chromatin domains data sets consist of Polymerase II, Polycomb (PC) domain proteins and Origin Recognition Complex (ORC) proteins as well as the derived chromatin states (enhancer/promoter states, transcribed chromatin states and heterochromatic states). In Fig. 6C, each component of 6A was plotted on a heatmap to display the amount of overlap occurring amongst elements of the set at specific coordinates. Heatmap tracks were drawn using custom modENCODE software. In fig. S11 the genome was linearly divided into partially overlapping blocks such that each base pair in the genome was sampled exactly twice in two distinct blocks and the corresponding blocks between data groups were compared and plotted.

#### Coverage of remaining conserved regions.

After accounting for all modENCODE data sets just under 9% of the identified conserved regions remain unaccounted for, or just over 19% of conserved regions if derived chromatin states and introns and not included as coverage. While it is possible some of the remaining conserved regions are false positives from the conservation identification algorithm (S32), conserved but non functional regions, or are not amenable to identification with modENCODE experiments (including a tiny fraction of conserved bases found within long repeats), it is likely that much of the remaining portion could be accounted for by examining new factors and cell lines.

When excluding introns and derived chromatin states from the analysis, exons (3.05%), insulators (6.09%), other chromatin binding proteins (6.16%) and transcription factors (4.57%) cover the highest percentages of conserved bases overall which are not covered by any other modENCODE data set group (based on groups defined in Fig. 6A). When counting insulators, transcription factors and other chromatin binding proteins together as single group of chromatin interacting proteins, the percentage of conserved bases uniquely covered by that group alone increases to 24.23%. We took all identified conserved bases in the genome and masked out those covered by exons, UTRs, Ago smRNAs, non-coding RNAs, PolII, Polycomb domain proteins and ORC peaks and looked for saturation within the chromatin interacting protein group. We removed the 6 smallest data sets that did not contribute much

overall or unique coverage and tested the amount of unique coverage of conserved bases gained from each of the remaining data sets. We found that they added an average of 44389 unique conserved bases ( $\sigma = 79682$ ) not covered by any of the other modENCODE data set with 82% adding at least 1KB of unique coverage (the equivalent of 77% of data sets overall when including the smallest 6). Although the rate of unique conserved bases added slows substantially from the first factor added to the last, it does not saturate after examining n-1 factors in addition to the other data sets. With higher resolution experiments, the amount of current coverage would likely appear reduced and it is improbable we would be any closer to reaching saturation.

We believe a portion of the remaining conserved region is likely accounted for by chromatin interacting factors which we have not performed ChIP on or which we need to examine in more cell lines. Fully accounting for all conserved bases in the genome, even with just a single level of annotation, will likely take some combination of examining new factors, testing the factors we have looked at in additional cell lines and under more conditions, identification of other novel functional regions and possibly require refined identification of conserved bases in the genome.

### **Physical regulatory network inference from regulator binding and motif instance conservation.**

We created the physical regulatory network (pRN) of TF binding and microRNA targeting, as inferred from the experimentally derived binding profiles of 76 TFs and genome-wide occurrences of evolutionarily conserved microRNA seed motifs for 105 microRNAs (77 distinct seeds). The TF binding profiles were obtained by integrating ChIP binding data sets from the modENCODE project and a literature survey that we conducted (Table S5). Note that for simplicity, we refer to these 76 regulators collectively as TFs even though some of them are co-factors or chromatin remodeling factors. To infer TF edges, we used TF binding in experimentally confirmed promoter regions for each target gene, and to infer microRNA edges, we used evolutionarily conserved motif occurrences for the 7-mer seeds of microRNA families in experimentally determined 3'UTRs. We further included regulatory edges from TFs to miRNAs based on binding of the TF in experimentally-inferred promoter regions for the precursor transcripts of miRNAs based on RNAseq (S5) and CAGE data (S33).

The resulting physical regulatory network contains 76 TF regulators with 158,558 edges targeting 12,482 targets, and 77 miRNA regulators with 5,805 edges targeting 3,099 targets. Of these regulatory interconnections, 1,847 TF edges specifically target other TFs, 89 TF edges target miRNAs, and 106 miRNA edges target transcription factors.

### **Structural and functional properties of the physical regulatory network.**

**Combinatorial regulation and pleiotropy:** The combined network shows high connectivity and rapid spread of regulatory information, requiring traversal of only ~2 edges on average between any two genes, and no more than 5 edges between any pair of genes. Target genes are regulated by ~12 TFs on average and can have up to 54 regulatory TFs, potentially leading to extensive combinatorial control. The most heavily-

targeted genes are typically associated with increased pleiotropy, as measured by the number of distinct functional processes they are associated with, with 4.76 specific GO terms for the top 10% most targeted genes (*p*-value <10e-5) and 3.18 for the top 50% (*p*-value <10e-5), compared to 1.81 terms for the bottom 50% least targeted genes (*p*-value <10e-5). For instance, *spi*, the most heavily regulated gene, is annotated based on experimental evidence with 11 different GO processes ranging from negative regulation of apoptosis to peripheral nervous system development and is known to play a role in at least 17 distinct developmental contexts (S34). Similarly, *shn* is regulated by 48 TFs, is annotated with 9 different biological processes and is a component of the BMP signaling pathway that is involved in a wide range of tissue patterning events during development (S35). Additionally, we validated this observation by looking at the number of tissue annotations per gene. This number is 10.66 for the top 10% most targeted genes with at least one tissue annotation (*p*-value <10e-5), 9.03 for top 50% (*p*-value <10e-5), 4.88 for the bottom 50% least targeted genes (*p*-value <10e-5), and 2.96 for the bottom 10% (*p*-value <10e-5).

**Hierarchical structure of pRN:** We examined structural and functional properties of this network through global and local topology metrics. In order to understand the hierarchical organization of regulation, we assigned the TFs to different levels based on the relative proportion of TF targets versus TF regulators for each TF. We augmented this network with the miRNA regulators most closely interacting with each level, separated into miRs solely targeting TFs and those also targeted by TFs. miRNAs were added to the network in two groups: miRNAs regulated by TFs (left in Fig. 7A) and miRNAs that only regulate TFs (right in Fig. 7A). We positioned miRNAs vertically according to their interactions with different TF levels. Also, the horizontal positions of TFs at each level show three different properties: the TFs on the left regulate miRNAs; the TFs in the middle have no regulatory edges to or from miRNAs; the TFs on the right are targeted by miRNAs but do not regulate them. TFs at various levels in the network are associated with different expression patterns in the fly developmental time course (S5). In general, TFs are expressed earlier in the time course than the target genes. Starting from top level to bottom, 7 out of 9, 11 out of 13, 17 out of 22, 9 out of 12 and 14 out of 20 TFs reach their peak expression values within the first 6 hours of embryonic development. Additionally although top, second level and third level regulators have no significant tissue enrichment, the bottom most regulators are enriched for certain tissues in later stages (e.g. Dorsal ectoderm epidermis stages 11-12, 13-16, Pharynx stages 11-12, Head Epidermis stages 13-16).

**Local connectivity patterns (network motifs):** To study the regulatory feedback, cooperation, and cross-regulation properties of TFs and miRNAs, we studied three-node network motifs that are significantly over-represented in the physical regulatory network (Fig. 7B). Previous studies in other model organisms have suggested that such over represented connectivity patterns correspond to basic circuit primitives associated with specific functional roles (S36). Of the 61 possible 3-node interconnection patterns between TFs, miRNAs, and target genes, we found eight network motifs that are significantly over-represented in the physical regulatory

network compared to their frequency in 100 random networks with the same in- and out-degree distributions (obtained by randomizing true networks while preserving the number of regulators and targets for each gene) (S37). Five of these motifs exhibit patterns of TF co operation (motifs 1, 2, 4, 7, 8), confirming our previous observations of extensive co-binding and co-targeting.

Moreover, for each of these five motifs, at least two TFs bind each other's promoter regions, suggesting extensive positive and negative feedback. Two other motifs correspond to mixed feed-forward loops. They involve cooperation of TFs and miRNAs (motifs 3, 6), which can lead to sign-sensitive delay properties in the expression of target gene depending on the activating or repressive action of the TF. Lastly, one motif (motif 5) corresponds to a feedback-loop of a downstream TF targeting an upstream TF through an miRNA, as previously suggested by observing the topological properties of the hierarchical network (See data S11 for complete lists of motif instances for these eight motifs). One notable motif instance involved the transcription factors SNAIL, TWIST and the miRNA mir-1, all three of which are co-expressed mesodermally. A number of regulatory connections of this motif instance is confirmed experimentally [activations of mir-1 and SNAIL by TWIST (S38, S39)].

### **Integration of expression, chip, motif and chromatin data sets to predict a functional regulatory network.**

**Network inference:** In addition to the physical regulatory network (pRN) based on ChIP binding data sets of TFs described earlier and to the motif regulatory network (mRN) (S40), we constructed five feature specific networks that were used to infer an integrated, functional regulatory network (fRN). Before explaining how each of these networks were inferred, we describe some attributes that are common to all of them. First, each edge has a score between 0 and 1 indicating the confidence of an interaction between a TF and a target gene (TG). These scores have been computed for all possible interactions between 707 TFs and 14,444 TGs (the number of genes for which we have at least one measurement in one data set). If a data set had missing values preventing the correct computation of a score, the default value of 0.1 was assigned to the edge.

We used three expression data sets: our RNAseq fly developmental time-course (S5), time course expression from Hooper *et al.* (S41), and expression data from FlyAtlas (S42). The three corresponding co-expression networks were built using squared Spearman correlation between TF expressions and TG expressions across all experiments.

Two chromatin data sets are also used in this analysis, the first one is a time course data set (S31) and the other is a cell-lines data set (S26). We first binned into five regions, a gene's DNA sequence ranging from 1Kb upstream of the transcription start site (TSS) to 1Kb downstream the transcription end site (TES). For each mark, we had five binary values corresponding to the presence/absence of the mark in that region. As a result, each gene has a binary string profile of (5xm) bits (with m the number of marks) as its chromatin profile. The six regions delimiters are TSS 1kb, TSS, start of CDS, end

of CDS, TES, TES+1kb respectively. Finally, we computed the similarity between the binary profiles of the TF and TG using squared Spearman correlation.

In order to generate the functional regulatory network (fRN), we used two different approaches, namely a supervised logistic regression trained on a set of known regulatory interactions from the REDfly database (S43), and an unsupervised average-rule method (S44). The unsupervised method computes a confidence for each edge, which is the average weight of that edge across the seven feature specific networks. In other words, the unsupervised approach is a linear model where the feature weights are all equal. We found that the two methods performed comparably with respect to recovery of REDfly edges.

Although, the supervised approach recovered a similar number of REDfly interactions (based on a 10-fold cross-validation), the unsupervised average-rule outperformed the logistic regressor on other validation metrics, such as enrichment (for each pair of genes sharing at least 50% of their regulators) for common gene ontology terms, correlation of gene expression profiles across the developmental time course, and co-expression in similar tissues as observed by in situ hybridization experiments (see validation metrics below). Thus, we selected the unsupervised network for further study in the remainder of the paper, given its unbiased nature. Out of the ~10M interactions that constitute the complete network, only 231,181 edges (top 2%) highest confidence edges were selected from the network generated by the unsupervised method. We selected 231,181 edges in order to have the same number of edges identified by the supervised method (number of edges that satisfy a probability of 0.5 or higher). The final fRN connects 576 TFs and 9,436 targets.

**Validation metrics of networks:** In addition to Precision-Recall curves using REDfly as gold standard, we evaluated the different networks based on different functional data of co-regulated pairs of genes: (a) common Gene Ontology (GO) terms, (b) correlation of gene expression profiles across the developmental time-course (S5), (c) common Imago tissues terms (S45), (d) frequency of protein-protein interactions (PPI). Each of these data are indicative of co regulatory relationships. For (a), (b) and (c) we considered co-regulatory relationships implied by multiple TFs. For (d) we focused only on co-regulatory relationships per TF. We estimated the enrichment of each metric in the true networks to randomized versions of these networks. Random networks were generated by selecting pairs of edges from the original network and swapping the targets, thus preserving both in and out degree distributions within the random networks.

To consider co-regulatory relationships from multiple TFs, we used the Jaccard coefficient of regulator sets of pairs of genes. Jaccard coefficient between two sets is defined as the ratio of the intersection to the union of the two sets. We selected those pairs of genes with Jaccard coefficient >0.5. Co-expression was measured by squared Spearman correlation. Similarity in GO and Imago terms was measured based on the Jaccard coefficient of the set of terms associated with the gene pair under consideration. For PPI, we calculate the proportion of the number of PPIs to the number of possible PPIs among the targets of a TF. For each of these metrics, we calculate an

average from the true networks and compare this to the average from the random networks. Fig.7C reports the positive predictive value. The value reported in Fig.7C is expected positive predictive value,  $(x-1)/x$ , with  $x$ =enrichment in network/enrichment in random network.

**Intersection of pairs of networks:** To estimate the set of common edges between a pair of networks we focused only on TFs and targets that are common between the two networks. Then we selected the subset of the edges connecting these common TFs and targets for each network. Overlap proportion between two networks expected by random,  $p$ , is the ratio of the size of the smaller edge set to the total number of edges between the common TFs and targets. The proportion of edges actually overlapped,  $q$ , is the ratio of the number of common edges to the size of the large edge set. Fold enrichment for intersection of two networks is the ratio  $q/p$ .

#### **Inferring functional gene modules from clusters of activity patterns along developmental time course.**

We clustered the expression profiles of 10,733 genes from the fly developmental time course using a Gaussian mixture model, which assumes that genes within each cluster are distributed according to a Gaussian distribution. The parameters of each cluster, a mean and a covariance matrix are estimated based on expectation maximization. The data were pre-processed as follows: add 1 to each expression value, log transform and standardize. The number of clusters,  $k$ , to maximize the penalized log likelihood of hold out data. The result of clustering is a probability that a gene belongs to a particular cluster. A gene is assigned to the cluster with the highest probability value. This probabilistic assignment is used in the down-stream analysis of function prediction described below. GO enrichment and motif enrichment of each cluster is present in fig. S15.

#### **Predicting function of genes using the functional regulatory network and co-expression.**

We developed a “guilt-by-association” approach to predict the biological process categories of genes with no previous GO annotations as well as predicting additional functions of genes with known annotations. The inputs to this approach are: (a) a set of probabilistic assignments of a gene for each cluster (obtained from above) (i.e. similarity of expression profiles), (b) normalized regulatory similarity for each pair of gene  $i$  and  $j$ , based on the Jaccard coefficient of the set of regulators associated with each gene  $i$  (i.e. similarity of regulatory program). Regulators for each gene were obtained from the functional regulatory network.

To predict the function of a gene  $g$  we define a score which combines the cluster membership probability as well as the regulatory similarity of the  $g$  with all other genes. We first estimate the contribution of every cluster to a  $g$ 's score for a functional category by summing the regulatory similarity of all other genes that are known to belong to the category. We then sum the per-cluster scores across all clusters, with the contribution of each cluster weighted by the cluster membership probability of  $g$ . The regulatory similarity is normalized on a per-gene basis using the sum of all regulatory similarities

obtained from all genes that are co-regulated with a gene. This allows genes that are co-regulated of lots of other genes to have a smaller contribution, as compared to genes that are co-regulated with fewer genes. For each gene, this procedure results in a score for each functional category, which we then transform into a confidence value based on the score distribution for previously annotated genes for each functional category.

We selected all the GO process terms for predicting process categories of 2832 genes. The total number of genes on which we could test our predictions was 5606, which was the intersection of genes that could be clustered (10733), had an edge in the functional regulatory network (9436) and was annotated with one of the GO processes.

For each GO process term we created a negative set based on the union of members of all GO terms that shared less than x% of its members with the term in question. We used the Jaccard coefficient of the member gene sets of the two terms to estimate the proportion of shared genes. We experimented with different values of the Jaccard coefficient (0.01, 0.05, 0.1, 0.15, 0.2) and found 0.1 (10%) to give the best performance. This negative set allowed us to define “false positives”, and thus compute an FDR as the ratio of the number of false positives to the sum of the number true positives and false positives. While estimating the number of genes for which we could predict function we excluded terms that were very general: *metabolic process*, *cellular process*, *biological process*, *biological regulation*, *cellular metabolic process*, *primary metabolic process*, *biosynthetic process*, *cellular biosynthetic process*.

### **Predicting regulators associated with developmental gene expression changes.**

The regulatory event map was constructed using the Dynamic Regulatory Events Miner (DREM) (S46). The expression data was transformed to the log base 2 ratio with the first time point, the Embryo 0-2h sample. A value of two was added to all expression values before taking the ratio in order to smooth the influence of low expression values on the ratio. To focus on genes which showed greater dynamic ranges of expression levels we filtered genes which had expression in the range of -3 to 3 at all time points leaving 6129 genes. The model search was constrained to consider only two and three way split events in the time series with DREM’s regularization parameter on model complexity set to 100. A minimum standard deviation of 0.5 was enforced on the gaussian distributions in order for DREM to avoid over-fitting the no or very low variance of sets of genes in RNA-seq data. DREM was run with the option to use only expression data to learn the model opposed to the option to learn the map also using TF-gene interaction information, thus enabling the same set of splits to be analyzed for enrichment based on multiple networks. The TFs predicted to be associated with a path out of split, were those with a significant number of targets genes on the path conditioned on the set of genes going into the split. An uncorrected p-value of <0.001 based on the hypergeometric distribution was used to determine the significance threshold. We note that unlike the physical network, the functional network was partially learned on a different developmental expression data so is not a completely independent validation source.

In evaluating if the expression change of regulators predicted to be associated with a split were greater than those not, we computed for each regulator the absolute value of the expression change at the time point immediately following the split with the expression at the split. The p-value significance for this average change being different reported in the main text was based on a one sided t-test.

#### **Motif enrichment in chromatin defined regions specific to S2/BG3 cells.**

We collected motifs from the literature (S29, S27–S54) and placed into a cluster all those that either shared a transcription factor or had a Pearson correlation of at least 0.75. The motifs were then matched to the genome (excluding coding exons, repeats, transposons, 3' untranslated regions and non-coding RNAs from FlyBase version 5.28) using a PWM threshold corresponding to a p-value of  $4^{-7}$  as determined by TFM-Pvalue (S55). We used conservation on the 12 fly species to filter motif instances to only those that reach a 30% confidence level (S31, S40). We compared the enrichment of these motifs in chromatin regions that were unique to either S2 (red) to BG3 (blue) cells using conservative, bias corrected log-odds with shuffled motifs as a baseline (see Supplemental Methods). In cases where multiple motifs existed (from either variants for the same factor or from groups of factors that were grouped due to motif similarity), we took the motif that had the highest maximal enrichment in some mark.

#### **Predicting gene expression from combinations of chromatin states.**

To relate the overall chromatin state of a given gene with the presence and magnitude of expression, we've utilized a Support Vector Machines (SVM) approach. The chromatin state of a given gene was described by a set of factors giving the mean enrichment level of each mark, top level, and top 10% average within entire gene body, TSS-proximal (within 2Kbp), and 3'-end proximal regions. The nuclassification was used to distinguish expressed and silent genes (defined based on the RNA-seq FPKM thresholds). 10-fold bootstrap validation was used in assessing the predictive power of the classifier. Similarly, nu-regression was used to build models relating the combined set of chromatin features of each gene to the expression magnitude (RNA-seq FPKM on a log10 scale).

#### **Predicting gene expression from expression of regulators.**

To predict expression of a target gene from the expression of its regulators, we obtained the set of regulators of every target gene from a regulatory network. We assume that the expression of a target gene is a linear combination of the expression level of the TFs. The weights of the linear combination are obtained from the training expression profiles of both TFs and targets. We learned these models for four networks: ChIP-based binding (bRN), motif-based network (mRN), the REDfly network (kRN) and the functional network (fRN). The number of genes for which we learned these models using different networks is described in Table S7.

We tested each model for its ability to predict the expression of the three time points that it had never observed, and compared the squared errors of the predicted gene expression profiles to the errors found by applying the same prediction scheme to 10

randomized control networks of same size and degree distribution, generated by permuting the labels of TFs and targets. We performed two types of comparisons against random networks. First, we estimated the average number of genes on which models learned from the original networks had a significantly lower error (squared error) than random networks (One-sided t-test p-value < 0.05). Second, we generated a list of “predictable” genes and a list of “unpredictable genes” that had consistently lower or higher errors than random networks. Specifically, a gene was considered predictable if it had lower error in the true network than in 6 out of the random networks, and was considered unpredictable if it had higher error in the true network than in 6 out of 10 random networks.

To understand why we were not able to predict expression of some genes, we estimated the correlations of gene time courses across two developmental time courses, which were generated by the modENCODE project, one of which is a 30 time points time course (S5) and the other is a 12 time points time course (S31). The shorter time course has a subset of sampling points of the longer time course, allowing us to collapse the longer time course by taking a mean of the extra time points. We found that the correlation distribution of genes that could be predicted was significantly shifted towards higher positive correlation (KS T-test 1E-7) compared to genes for which we could not predict expression (fig. S21). This suggests that the genes that we cannot predict expression for tend to have variable, random fluctuations, which cannot be captured by the variations in expression levels of its regulators.

To assess how well our predictable set of genes were able to predict expression in a completely new setting, we used the learned models from the time course to predict expression in 4 cell lines. For each network we made a prediction for cell-lines for each gene using the models learned from each fold of cross-validation. We did the same for each random network and estimated the number of genes on which the errors from the true networks were significantly lower than the random network, and repeated this for all random networks. We report the average number of genes on which the models learned from the true networks were better than random (Table S7). We repeated this procedure for the unpredictable set.

***Drosophila* modENCODE Consortium membership list.**

Listed authors: Sushmita Roy, Jason Ernst, Peter V Kharchenko, Pouya Kheradpour, Nicolas Nègre, Matthew L Eaton, Jane M Landolin, Christopher A Bristow, Lijia Ma, Michael F Lin, Stefan Washietl, Bradley I Arshinoff, Ferhat Ay, Patrick E Meyer, Nicolas Robine, Nicole L Washington, Luisa Di Stefano, Eugene Berezikov, Christopher D. Brown, Rogerio Candeias, Joseph W Carlson, Adrian Carr, Irwin Jungreis, Daniel Marbach, Rachel Sealfon, Michael Y Tolstorukov, Sebastian Will, Artyom A Alekseyenko, Carlo G Artieri, Benjamin W Booth, Angela N Brooks, Qi Dai, Carrie A Davis, Michael O. Duff, Xin Feng, Andrey A Gorchakov, Tingting Gu, Jorja G Henikoff, Philipp Kapranov, Renhua Li, Heather K MacAlpine, John H. Malone, Aki Minoda, Jared T Nordman, Katsutomo Okamura, Marc D. Perry, Sara K Powell, Nicole C. Riddle, Akiko Sakai, Anastasia A Samsonova, Jeremy E Sandler, Yuri B Schwartz, Noa Sher, Rebecca F. Spokony, David Sturgill, Marijke J van Baren, Kenneth H. Wan, Li Yang, Charles Yu, Elise A Feingold, Peter J Good, Mark S Guyer, Rebecca F Lowdon, Kami Ahmad, Justen Andrews, Bonnie Berger, Steven E Brenner, Michael R Brent, Lucy Cherbas, Sarah C. R. Elgin, Thomas R Gingera, Robert L Grossman, Roger A Hoskins, Thomas C. Kaufman, William J Kent, Mitzi I Kuroda, Terry L Orr-Weaver, Norbert Perrimon, Vincenzo Pirrotta, James W. Posakony, Bing Ren, Steven Russell, Peter Cherbas, Brenton R Graveley, Suzanna E Lewis, Gos Micklem, Brian Oliver, Peter J Park, Susan E. Celniker, Steven Henikoff, Gary H. Karpen, Eric C Lai, David M MacAlpine, Lincoln D Stein, Kevin P White, Manolis Kellis

Additional investigators: David Acevedo, Richard Auburn, Galt Barber, Hugo J. Bellen, Eric P Bishop, Terri D. Bryson, Aurelien Chateigner, Jia Chen, Hiram Clawson, Charles L. G. Comstock, Sergio Contrino, Leyna C. DeNapoli, Queying Ding, Alex Dobin, Marc H. Domanus, Jorg Drenkow, Sandrine Dudoit, Jackie Dumais, Thomas Eng, Delphine Fagegaltier, Sarah E Gadel, Srinka Ghosh, Francois Guillier, David Hanley, Gregory J Hannon, Kasper D. Hansen, Elizabeth Heinz, Angie S Hinrichs, Martin Hirst, Sonali Jha, Lichun Jiang, Youngsook L. Jung, Helena Kashevsky, Cameron D. Kennedy, Ellen T Kephart, Laura Langton, Ok-Kyung Lee, Sharon Li, Zirong Li, Wei Lin, Daniela Linder-Basso, Paul Lloyd, Rachel Lyne, Sarah E Marchetti, Marco Marra, Nicolas R Mattiuzzo, Sheldon McKay, Folker Meyer, David Miller, Steven W Miller, Richard A Moore, Carolyn A. Morrison, Joseph A. Prinz, Michelle Rooks, Richard Moore, Kim M Rutherford, Peter Ruzanov, Douglas A Scheftner, Lionel Senderowicz, Parantu K. Shah, Gregory Shanower, Richard Smith, E.O. Stinson, Sarah Suchy, Aaron E Tenney, Feng Tian, Koen J.T. Venken, Huaijen Wang, Robert White, Jared Wilkening, Aarron T Willingham, Chris Zaleski, Zheng Zha, Dayu Zhang, Yongjun Zhao, Jennifer T Zieba

**Author contributions:** **Integrative analysis (PI: Kellis)**. Coding potential/RNA structures: Lin, Washietl, Jungreis, Will, Berger, Carlson, Celniker, Kellis. Transcript models: Landolin, Carlson, Carr, Oliver, Graveley, Celniker. **microRNAs**: Berezikov, Robine, Lin, Kellis, Lai. **Chromatin integration**: Ernst, Kharchenko, Eaton, Bristow, Tolstorukov, Henikoff, Park, MacAlpine, Kellis, Karpen. **TF integration**: Kheradpour, Ma, Brown, Nègre, Sealoff, Bristow, Eaton, MacAlpine, White, Kellis. **Genome coverage analysis**: Arshinoff, Karpen, Stein, Kellis. **Regulatory networks**: Roy, Ay, Meyer, Candeias, Marbach, Bristow, White, Kellis. **Predictive models**: Roy, Ernst, Kheradpour, Kellis. **Scientific management**: Di Stefano, **Data sets**: Washington, **Program management**: Feingold, Good, Guyer, Lowdon. **Manuscript preparation**: Stein, Park, Henikoff, Lai, White, Celniker, MacAlpine, Karpen, Kellis. **Primary data sets**. **Transcription (PI: Celniker)**: Andrews, Booth, Brooks, Cherbas, Davis, Duff, Kapranov, Samsonova, Sandler, van Baren, Wan, Yang, Yu. **coPIs**: Brenner, Brent, Cherbas, Gingeras, Hoskins, Kaufman, Perrimon, Graveley, **Other contributors**: Charles L. G. Comstock, Alex Dobin, Jorg Drenkow, Sandrine Dudoit, Jackie Dumais, Delphine Fagegaltier, Srinka Ghosh, Kasper D. Hansen, Sonali Jha, Laura Langton, Wei Lin, David Miller, Aaron E. Tenney, Huaien Wang, Aarron T. Willingham, Chris Zaleski, Dayu Zhang. **Chromatin (PI: Karpen)**: Alekseyenko, Gorchakov, Gu, Minoda, Riddle, Schwartz, **coPIs**: Elgin, Kuroda, Park, Pirrotta, **Other**: David Acevedo, Eric P Bishop, Sarah E Gadel, Youngsook L. Jung, Cameron D. Kennedy, Ok-Kyung Lee, Daniela Linder-Basso, Sarah E Marchetti, Gregory Shanower. **TFs (PI: White)**: Spokony. **coPIs**: Grossman, Posakony, Ren, Russell, **Other**: Richard Auburn, Hugo J. Bellen, Jia Chen, Marc H. Domanus, David Hanley, Elizabeth Heinz, Zirong Li, Folker Meyer, Steven W. Miller, Carolyn A. Morrison, Douglas A. Scheftner, Lionel Senderowicz, Parantu K. Shah, Sarah Suchy, Feng Tian, Koen J. T. Venken, Robert White, Jared Wilkening, Jennifer Zieba. **Replication (PI: MacAlpine)**: Matthew L. Eaton, Heather K. MacAlpine, Jared T. Nordman, Sara K. Powell, Noa Sher. **coPI**: Terry L. Orr-Weaver, **Other**: Leyna C. DeNapoli, Queying Ding, Thomas Eng, Helena Kashevsky, Sharon Li, Joseph A. Prinz. **Small RNAs (PI: Lai)**: Berezikov, Dai, Okamura. **coPI**: Gregory J. Hannon, **other**: Martin Hirst, Marco Marra, Michelle Rooks, Yongjun Zhao. **Nucleosomes (PI: Henikoff)**: Henikoff, Sakai, **coPI**: Ahmad, **Other**: Terri D. Bryson. **Comparative transcription (PI: Oliver)**: Artieri, Li, Malone, Sturgill, **Other**: Lichun Jiang, Nicolas Mattiuzzo. **Data coordination center (PI: Stein)**: Feng, Perry, Washington. **coPIs**: Kent, Lewis, Micklem, **Other**: Galt Barber, Aurelien Chateigner, Hiram Clawson, Sergio Contrino, Francois Guillier, Angie S. Hinrichs, Ellen T. Kephart, Paul Lloyd, Rachel Lyne, Sheldon McKay, Richard A. Moore, Chris Mungall, Kim M Rutherford, Peter Ruzanov, Richard Smith, E. O. Stinson, Zheng Zha

Full list of author affiliations and author contributions is found in Table S8.

## **Supplementary Figure Legends**

**Figure S1: Conserved RNA secondary structures.** Examples of newly identified local RNA structure elements in the long noncoding RNAs roX2 and HSR-omega.

**Figure S2: Genome-wide correlation of chromatin marks with patterns of physical properties.** Correlation between nucleosome density, low salt soluble fraction of chromatin (80mM), CATCH-IT measurements (AHA 3h) and histone marks is shown. Each cell in the heat map is colored according to the value of Pearson correlation computed for the enrichment profiles of the corresponding chromatin mark (indicated as column name) and physical property (indicated as row name). Hierarchical clustering reveals distinct groups of marks associated with specific patterns of chromatin properties (indicated with green lines) with the marks of silent chromatin being on the left side of the heat map and marks of the active chromatin on the right side. The active chromatin marks positively correlate with chromatin solubility and nucleosome turnover rate, and the silent chromatin marks show opposite trends correlating positively with nucleosome density. Most marks in the middle group exhibit positive correlation with both high nucleosome turnover rate and high nucleosome density, while being anti-correlated with chromatin solubility. These patterns of physical properties are consistent with easier accessibility of active chromatin for regulatory proteins.

**Figure S3: Detecting transcribed regions based on enrichment of histone marks.**

The putative transcribed regions (TR) were detected based on a combination of two marks typically found in the expressed gene bodies: H2B ubiquitination and H3K36me3.

**A.** The plot illustrates the correspondence between the combined TR predictions and annotated expressed genes. The plot shows the fraction of each detected TR fragment overlapping with an expressed gene. Most of the fragments are found within the expressed genes. The mean fraction is indicated by the vertical red line. A fraction expected from randomized expression values is shown by the blue vertical line. The legend shows sensitivity and specificity of such fragments based on the current annotation. Note that because true negative regions are not available, the specificity is difficult to estimate. **B, C.** Analogous plots for the individual marks used in TR predictions. Out of all examined marks, H2B-ubiquitination appears to be the best predictor of TRs.

**Figure S4: Predictive model of active promoters in S2 cells.** **A.** Length distribution of promoters with evidence from 2 or more of the CAGE, RACE, RE-EST assays. **B.** Distance from promoter annotation midpoint to TSS (grey: annotation overlaps TSS, black: annotation does not overlap TSS). **C.** Precision-recall curve for the test set of valid promoters. The supervised predictions are compared to the chromatin state models described in the supporting text. **D.** Fold enrichment of the valid promoter sequences within promoter prediction sets. **E.** Fold enrichment of active first exons within promoter prediction sets. **F.** Fold enrichment of capped short RNA reads derived from PolIII within promoter prediction sets.

**Figure S5: Early origin activity, replication timing and cohesin binding in >10 kb H3K36me1 intergenic regions.** The large (>10 kb) intergenic H3K36me1 domains were analyzed in BG3 (133) and S2 (78) cells for their replication and cohesin activity. In the first panel, quantile normalized Nipped-B S2 and BG3 ChIP-chip probe values (S56) from these regions were graphed in a notched boxplot in S2 (blue) and BG3 (green). A background level of enrichment was determined by taking the median signal of 100 samples of sets of segments with identical width, number and chromosome distribution to the set of intergenic H3K36me1 domains and evaluating the Nipped-B enrichment within these segments. A representative sample set is shown as the boxplots of lighter shade, and the median background level across all of the random sample sets is shown as the dashed lines. In the second panel, quantile normalized early origin probes from these regions were treated analogously to the first panel. The third panel is analogous to the second using array probe values from quantile normalized replication timing tiling arrays.

**Figure S6. Data analysis used for annotation of chromatin states. A.** As in Fig. 4, the 30 state discrete model on the left, was grouped based on similarity to the 9 intensity-based states (color coded) based on correlations reported in Kharchenko *et al.* (S26). For each state, enrichments are shown for different gene elements (panel 1), expression levels (panel 2), non-histone proteins (panel 3), and chromosomes and regions (panel 9), as described in Kharchenko *et al.* (S26); numbers show fold-enrichment relative to random expectation. \_Het\_ refers to heterochromatin regions (S57). Intron enrichment is inferred from the presence of transcripts (panel 1, column 5) that do not display enrichment for other gene elements. Enrichments of states at other modENCODE mapped elements are shown in panels 4-8; ORC binding and early origins (S58), nucleosome solubility (S59) and turnover (S60), insulators proteins (S31), HOT regions, histone deacetylases (HDACs), and transcription factors (TFs) (S31). For insulators and HDACs, factors were subdivided into subgroups based on their differential genomic localization, and the TFs were subdivided into 5 large categories based on their main molecular or biological functions. **A.** shows the median enrichments for each insulator, HDAC and TF subgroup, based on the data shown in **B**. Data sets are named according to the following convention: Factor\_Sample\_Lab\_Comment. Factor being the protein for which the ChIP data has been generated, Sample being the material from which chromatin has been extracted, Lab reflecting which group generated the data sets (gk: Gary Karpen's group; kw: Kevin White's group) and Comments giving additional information to distinguish similar data sets (seq: generated by ChIP-seq instead of ChIP chip; 1-2-3: different antibodies for the same data set; GFP: use of a GFP-tagged epitope to perform the ChIP experiment)

**Figure S7: Maximum chromatin state enrichments with increasing numbers of states.** This figure shows the emission parameters of the state with the greatest enrichment for selected categories for models with every number of states between 2 and 60. The last two columns indicate the fold enrichment for this state for the category and the percent of the category covered by the state. The values in the last two columns

can also be viewed in a line graph form at the bottom of each page. In the line graph, the x-axis is the number of states and for the graph on the left the y-axis is the maximum fold enrichment and for the right the y-axis is the percent of the category covered. (a) Panel A shows the category is MB8 Transcription Start Sites (TSS). (b) A MB8 gene that overlaps heterochromatin where heterochromatin is defined by chr4, the ‘Het’ chromosomes and the heterochromatic portions of the regular chromosomes determined based on cytogenetics (S61, S62). (c) ORCs that were more than 2kb away from a MB8 or Flybase 5.12 TSS. These panels show that increasing number of states can increase the maximum enrichment of a state, in a step-like manner, while generally decreasing in the absolute percent of the category the maximum enriched state contains, enabling us to guide the selection of model complexity based on resolution and the available number of chromatin marks. For heterochromatin transcribed genes there was a relatively large increase in enrichment going from 29 to 30 states.

**Figure S8: Comparison of 30 chromatin states with 5 “colors” of Filion et al.** The columns on the right hand side show the fold enrichment of each state of the 30-state based on S2 cells for the 5 colors of chromatin in Kc cells as defined by (S63). This analysis demonstrates similarities in the chromatin types identified by the different methods, but importantly shows that individual colors often map to different states, and individual states can contain multiple colors. A similar comparison of the 9-state model is described at (S26).

**Figure S9: Chromatin properties for HOT regions at varying complexity levels.** By examining sequentially more complex subsets of the HOT regions throughout the genome with respect to nucleosome dynamics, we find that nucleosome density (A) and CATCH-IT (B) are both correlated with HOT complexity. In each panel, the x-axis denotes the lower bound of HOT complexity considered. For each value, a 2 kb region centered around HOT spots of sufficient complexity was examined for these four markers of nucleosome dynamics. The mean was taken for each HOT region, and then displayed in a box plot. The data was taken from the following GEO accessions: Nucleosome Density: GSM333835, GSM333840 and GSM333844; CATCH-IT (20 min): GSM494310 and GSM494311, H2Av: GSM333849; H3.3: GSM333860. These data were collected in (S59, S60). (D) Hotspots and chromatin states. Each row corresponds to a chromatin state from Figure 4. Each column indicates the fold enrichment or depletion of the chromatin state for locations which have a TF complexity equal to or above the threshold indicated in the column. The figure shows several chromatin states show strong enrichments or depletions for TF hotspots.

**Figure S10: Fraction of constrained and unconstrained bases in modENCODE data sets.** The 95% confidence intervals for random overlap of elements with constrained and unconstrained bases were calculated using 10000 iterations of the GSC algorithm (S64). The green columns indicate the overlap of the indicated element with constrained bases; when the top of the column is outside and above the confidence interval then there is statistical enrichment in conservation in the data set. Similarly, the orange columns indicate overlap of elements with unconstrained bases. Coding exons (p-value < 1.0e-299), Ago-bound short regulatory RNAs (p-value 1.3e-27), PolII (p-value

9.6e-19), HOT regions (p-value 3.9e-3), TFs (pvalue 3.7e-31), other binding proteins (p-value 4.7e-8), polycomb domains (p value 9.7e-8), ORC binding (p-value 7.0e-6), transcribed chromatin states (p-value 2.1e-5) and heterochromatic states (p-value 1.3e-4) show statistical enrichment for conserved regions of the genome, while UTRs (p-value 2.0e-5) and introns (p-value 1.7e-43) are statistically depleted for conserved bases. The conserved fraction within the noncoding RNAs, insulators and enhancer/promoter states sets are within the bounds expected by chance.

**Figure S11: Density correlations from broadly group modENCODE data sets.** **A.** shows a dot plot of protein binding element density vs. chromatin interacting element density using window size of 100bp. A very tight fitting positive linear relationship is observed between the data sets with slope of 0.077;  $R^2=0.849$ . 95% confidence interval for slope is (0.07617,0.0784). **B.** shows a dot plot of protein binding element density vs. transcribed element density using a window size of 20,000bp. A positive linear relationship with slope of 0.203;  $R^2=0.146$ . 95% confidence interval for slope is (0.19149, 0.21656) **C.** shows a dot plot of protein binding element density vs. Transcription Start Site (TSS) density using a window size of 20,000bp. A positive linear relationship is observed as in C, but tighter fitting and with a sharper slop of 0.782;  $R^2=0.406$ . 95% confidence interval for slope is (0.76257,0.80194).

**Figure S12: Chromatin state recovery with increasing subsets of marks.** The chart considers recovery of chromatin state assignments based on increasing sized subsets of the 18 profiled marks. The subset evaluation is based on the recovery of the BG3 maximum posterior state assignments for the 30-state chromatin model obtained with all the marks. In evaluating a subset of marks, those marks not in the subset would be omitted from the emission distribution of the HMM when computing state assignments. The y-axis indicates the percent of locations assigned to a state using all 18-marks that would receive the same state assignment based on the subset of marks averaged over all 30-states. The chart compares the recovery averaged over all subsets of a given size, applying the optimal subset from S2 to select marks for BG3, and using the optimal subset for BG3. The optimal subset for BG3 could only be known by actually mapping all the marks, but the chart indicates how information from a previously mapped cell type can suggest a near-optimal subset of informative marks, to map in a new cell type.

**Figure S13: Heatmap of TF-TF and miR-TF interactions.** The first panel of this figure shows the expression correlations among TFs projected onto the space of regulatory interactions in between TFs. TFs are ordered according to their locations in Fig.7A (left to right then top to bottom) and green and red colors indicate where the TF of the corresponding row regulates the TF in the corresponding column. Darker shades of green (red) imply positive (negative) correlation between the expression profiles. Gray lines are separating the levels and white color implies no regulatory interaction. The second panel shows the regulatory interplay between TFs and miRNAs. The color code is: green indicates a regulatory edge from TF to miRNA, red shows a regulatory edge from miRNA to TF and blue means that there is an edge from TF to miRNA as well as an edge from miRNA to TF. We identified 89 regulatory edges from TFs to miRNAs and 106 edges from miRNAs to TFs. The average out-degree (the number of outgoing

regulatory edges) is 2085.12 for a transcription factor and 75.39 for a miRNA (medians are 1229 and 53, respectively). A target gene is regulated by ~12 TFs and ~0.45 miRNAs on average, suggesting combinatorial regulation of TFs and additional tuning from miRNAs on the regulatory program of a gene.

**Figure S14: Comparison between motif, bound, known and functional network.** **A.** Definitions of the four networks compared. Predicted motif-based physical network (motif network, or mRN), experimentally-derived ChIP-based physical network (“binding network” or bRN), the known curated functional regulatory network (REDfly, known, or kRN) and the functional regulatory network predicted by integration of motifs, binding, and correlated activity patterns (fRN). **B.** Statistics on network size. Number of TFs, number of targets, and average number of regulatory connections per target shown for the four networks. **C.** Pairwise similarity between all four networks. For each pair of networks, the number of common TF and targets is shown, and the enrichment in common connections. **D.** Functional enrichments for the four networks. We evaluated each network based on four measures of functional enrichment for co-targeted genes. For each network, we compared the enrichment found for common gene ontology terms (Function), correlation of gene expression profiles across the developmental time-course (DevExpr), co-expression in similar tissues as observed by in situ hybridization experiments (Tissues) and protein-protein interactions (PPI). Reported true positive rate is calculated from the enrichment compared to a shuffled network preserving the number of TF, targets, and the degree distribution. For example, a 1.5-fold enrichment for each metric compared to a random network corresponds to a true positive rate of 33% (for every 150 predictions, 100 are expected to also occur by chance, leading to 50 true positives, or 1/3).

**Figure S15: Heatmap showing expression, motif-enrichment and GO-enrichment for each cluster.** Each row corresponds to a cluster. Left panel are average expression profiles of genes in the cluster. Middle panel shows GO process terms (top 5 most enriched), with colored block indicating that the cluster is enriched with the term in the column. Right panel shows the motifs (top 5 most enriched) that are enriched in the cluster.

**Figure S16: Precision-recall curves of feature specific and integrated functional regulatory network.** Precision recall curves were computed based on REDfly database as the true positive, and the complement of REDfly as negative. Different feature specific networks are shown in different colors, functional regulatory network in black. The curve for ChIP network is essentially flat because edge weights had only 3 values: 0.1 for all edges from a TF if no ChIP data is available for the TF, 0 if ChIP is available but no binding is measured, and 1 if there is binding.

**Figure S17: Number of GO process predictions versus FDR.** Number of predictions for genes with no previous annotation as well as genes with some annotations are shown as a function of different FDR cutoffs. The number of predictions is estimated using expression only (diamond), regulatory network only (square) and using both expression and regulatory network (triangle).

**Figure S18: Fold-enrichment of tissues among novel and known members of a GO process category.** Enrichment of Imago tissues was estimated for novel and known gene sets of a GO process. Each pair of rows shows the fold enrichment of a set of genes with a tissue annotation within the known and novel gene set associated with a GO process. The intensity of the blue is proportional to fold enrichment. Overall, 141 tissues showed enrichment across 115 GO processes with predicted new members (FDR<0.4). Of these, 95 tissue terms were also enriched and associated with 35 GO processes

**Figure S19: Information on developmental time course expression splits identified by DREM.** Each row corresponds to a path out of a split identified by DREM. Paths coming out of the same split are grouped. The first column indicates the sample of the split. The second column reports the mean expression of genes on that path and the third column the expression level of the child at the next time point. The fourth column reports the difference between these two values. The fifth columns reports the number of genes assigned to that path. The next two sets of four columns pertain to the physical and functional networks respectively. These four columns for each network reports the number of significant TFs at an uncorrected p-value of 0.001, the number of non-significant TFs at that p-value, the average absolute expression change of significant TFs at the time point immediately following a split compared to the time point of the split, and the same for non-significant TFs. The last set of four columns report the top significant Gene Ontology biological process for a paths out of splits conditioned on the set of genes going into the split, the uncorrected p-value, an estimated corrected p-value for multiple testing based on a randomization, and the fold enrichment for the category.

**Figure S20: Squared error, correlation coefficients, and regulators of 1487 predictable genes.** Each row corresponds to a gene, columns from left to right are: mean squared error between true and predicted expression profiles, correlation coefficient between true and predicted expression profile, correlation coefficients of the expression profile of the target and 10 strongest predictors (absolute value of coefficient) in decreasing order of the coefficient, the names of the 10 strongest predictors, colored based on the sign of correlation: green: negative correlation, red: positive correlation.

**Figure S21: Squared error, correlation coefficients, and regulators of 324 unpredictable genes.** Same as in figure S20 but for the 324 genes for which we could not predict expression better than random.

**Figure S22: Cumulative distribution of correlations of individual gene time courses obtained from the two developmental RNA-seq data sets.** Cumulative distribution of Pearson's correlation between time courses of predictable genes (blue) and not predictable genes (red) in the two developmental time courses (S5, S60).

**Figure S23: Expanded set of gene expression predictions from the expression of inferred regulators.** Expression levels for each gene were predicted as a linear

combination of the expression level of its regulators, with each 30-stage time-course data set split into 10 groups of 3 stages each, and weights learned on 9 and evaluated on the 10th. The predicted expression is shown (dotted blue line) compared to the true expression (top of shaded area) for target genes. The expression of the top 5 positive and negative regulators is shown for each target gene, ranked by their contribution to the expression prediction (the weights of the linear regression model). Top row shows examples from the top 25 scoring reconstructions, second row shows examples from 30-100, third row from 100-500 and fourth row from 500-1488, ordered by their prediction quality scores (rank in upper right corner), evaluated as the averaged squared error between predicted and true expression levels across the time-course. Color indicates the sign of correlation of a regulator and a target, red for positive (+) correlation, green for negative (-) correlation. Color scale is proportional to strength of correlation or anti-correlation.

### **Supplementary Table Legends:**

**Table S1: Data sets overview.** Cell lines, strains, and developmental stages assayed for each type of genomic element and function with accession numbers of all data generated. Each row indicates the annotation types sought for each class of functional elements and the corresponding assays carried out towards obtaining them. The actual number of unique experiments performed for each assay is indicated (# Expts), and the cell lines, developmental stages, and number of tissues in which they were carried out. The specific factors used in each type of assay are also listed, though not every factor has yet been tested in each experimental condition. All data sets, and extensive metadata of experimental designs, protocols and reagents are accessible through the Data Coordination Center (DCC) (<http://www.modencode.org>), modMine (<http://intermine.modencode.org>), and through a cloud computing environment. Results will also become available through the FlyBase community resource (<http://www.flybase.org>) (S2). modENCODE DCC, GEO, and SRA identifiers are listed and hyperlinked for each data set. Footnotes: (1) 22 other cell lines tested for RNA expression are 1182-4H, CME W2, CMEL1, DmBG1-c1, DmBG2-c2, DmD11, DmD16, DmD17, DmD20-c2, DMD20-c5, DmD21, DmD32, DmD4- c1, DmD8, DmD9, GM2, S1, S2R+, S3, Sg4, mbn2; 17 other cell lines tested for ncRNA expression are 1182-4H, GM2, Kc, Kc-Rubin, ML-DmBG1c1, ML-DmD21, ML-DmD20c5, OSS, S2, S2-NP, S2-Rubin, S2R+, Sg4, DmD9, DmD32, CME-L1, S2-NP. (2) For developmental stages, a single generic stage for embryo, larva, pupa, and adult are listed for brevity, though several sub-stages were assayed. Generally, the sequenced strain *y; cn bw sp* was used. In addition, Oregon-R was tested in RNA and small ncRNA expression, as well as in copy number variation assays. Additional strains tested for ncRNA expression include Canton-S, *r2d2[1]*, *Dcr-2[L811fsX]*, and *lqsf00791*. Some GFP-fusion strains were used to identify TFs, as described in 6. (3) Factors are listed as protein names. (4) The cDNA and RACE sequencing experiments were performed in replicate libraries derived from multiple targeted PCR primers; each library is counted as a separate experiment. (5) Short ncRNAs were also tested in isolated tissues: Testes, Ovary, Head, Body, Imaginal Disk. (6) Copy number variations were ascertained by CGH in the following isolated tissues: embryonic salivary glands and follicle cells. (7) Though listed under "factors", the following TFs were assayed in strains built with GFP or YFP fusions: CAD, DFD, ECR, EVE, EXD, EN, SENS, TLL.

**Table S2: *Drosophila* cell lines characterized by modENCODE.** This table describes the cell lines characterized in the course of the modENCODE project. The first four rows list the cell lines comprehensively characterized in the course of the transcription, replication, and transcriptional regulation aspects of the project. The remaining 29 cell lines were characterized for transcriptional expression and splicing patterns only.

**Table S3: Protein symbols used in the paper.** This table describes the chromatin modifying proteins, transcription factors, and other chromatin-associated proteins used in the manuscript.

**Table S4. cis-NAT-siRNA production in *Drosophila melanogaster*.** Shown are the MB8 genomic coordinates and gene names of sense-antisense overlapping transcript

models (cis-NATs) that exhibit enrichment for siRNAs. We considered those loci that satisfied the following criteria: (1) the overlap region generated at least fifty 21 nt reads, (2) 21 nt reads comprised at least 2/3 of all reads mapped to a given overlap, and (3) 21 nt reads derived from both strands of the overlap. The input for analysis was small RNA data aggregated from 187 libraries, including published data and additional unpublished data generated for the modENCODE project (S65).

**Table S5. ChIP data sources used to create the pRN.** This table lists the sources of the ChIP binding data sets for the transcriptional regulators of the pRN.

**Table S6: Number of gene predictions for GO process terms.** This table lists the set of GO process in which we can make predictions at a confidence of 40% or higher. The first column is the total number of calls. The second column is the number of hits for the term. The third column is the number of predictions of genes with some existing annotation, which are not in the negative set of the term of interest and not annotated with the term of interest. The fourth column is the number of predictions of genes with no previous annotations. FDR is calculated as the ratio of the number of false positives to the number of true positives + false positives. The number of false positives is equal to the number of negatives that are predicted as false. Rows in bold were selected for showing ROC curves in Fig. 8.

**Table S7: Relative performance of different networks for predicting expression of target genes.** Columns in gray specify the number of TFs, targets and average number of regulators per target for predicting expression for the four regulatory networks. Columns in green report the cross-validation performance of different networks. # predictions better than random: average number of genes on which the true network was better (lower prediction error) than a random network; # predictions worse than random: average number of genes on which the true network was worse than a random network Enrichment: ratio of the numbers in the previous columns. Columns in pink report the number of predictable and unpredictable genes in different networks. Columns in blue report performance on cell lines: first column is the average number of predictable genes that can predict expression in cell lines, and the second column is the proportion of the number of cell-line predictable genes to the total number of predictable genes.

**Table S8: Full author list and author contributions.** For each larger group (Integrative Analysis, Transcripts, Chromatin, Transcription Factors, Replication, Small RNAs, Nucleosomes, Data Coordination Center, Comparative Transcription, RNA families, and Program Management) are listed Co-First Authors (Authors 1-17, Red), Co-Second Authors (Authors 18-27, orange), Other Major Contributors (Authors 28-62, Yellow), Program Management Staff (Authors 63-66, light green), Principal Investigators (PI) (Authors 93-99, Purple), Co-Principal Investigators with primary role in analysis (Authors 87-92, light brown), Other co-PIs (Authors 67-86, darker green), and consortium co-investigators (Authors 101-173, white). Author 100 (dark blue) directed the integrative analysis. For each author are listed the larger group in which they worked, and the immediate PI that they worked with, their contribution, section to which they

contributed, and their affiliation.

## List of supplementary data sets.

All data is available at:

[http://www.modencode.org/publications/integrative\\_fly\\_2010/](http://www.modencode.org/publications/integrative_fly_2010/)

**Data S1: Listing of all modENCODE data sets and accession numbers.**

**Data S2: Novel transcripts with complete or incomplete conserved ORFs.**

**Data S3: Noncoding RNA-seq transcripts with conserved RNA secondary structure.**

**Data S4: Novel promoter predictions based on chromatin signatures.**

**Data S5: Genome annotated using the 9-state chromatin model.**

**Data S6: Genome annotated using the 30-state chromatin model.**

**Data S7: Predicted TFBS.**

**Data S8: HOT regions.**

**Data S9: Machine-parseable version of ChIP-based physical regulatory network (pRN) combined with miRNA regulations.**

**Data S10: Summary of connectivity properties of pRN.**

**Data S11: Lists of all motif instances for the 8 motifs shown in Figure 7B.**

**Data S12: Feature-specific data sets for functional network reconstruction.** This zip file contains six tab-separated files representing three types of data sources used for the reconstruction of the functional regulatory network: evolutionary conserved motifs (*motif\_net.txt*), chromatin profiles (*Chromatin\_profiles\_cell\_lines.txt*, *Chromatin\_profiles\_time\_course.txt*) and expression (*RNAseq\_profiles\_time\_course.txt*, *GSE6186\_profiles\_time\_course.txt*, *flyatlas\_profiles.txt*). Genes are indicated by their FBgn numbers. *motif\_net.txt* contains the regulatory network derived from evolutionarily conserved motif instances (mRN) in promoter regions of target genes. Each line in this file corresponds to an edge: the first column is a TF, second column is a target gene, and the third column is the confidence of a TF-target interaction, which is proportional to the extent of evolutionary conservation of a motif instance.

*Chromatin\_profiles\_cell\_lines.txt* and *Chromatin\_profiles\_time\_course.txt* contain the chromatin profiles of every gene. Each line corresponds to a gene and each column is a combination mark/gene region (See supporting text and methods for details).

*RNAseq\_profiles\_time\_course.txt*, *GSE6186\_profiles\_time course.txt* and *flyatlas\_profiles.txt* are the three expression data sets. Each line corresponds to a gene and the columns represent developmental time points or tissues where expression measurements were made.

**Data S13: Functional regulatory network.** This file has three columns, first column is the name of the regulator, second column is the name of the target gene, and third column is the score, which is directly proportional to the confidence of the regulator regulating the target.

**Data S14: Gene co-expression clusters.** This zipped directory has two files. The tab-delimited text file has the cluster assignments for the genes. The first column is the name of the gene and the second column is the cluster ID. The pdfs show the expression profiles, enriched motifs, and enriched GO process terms per cluster. The name of the file is associated with the cluster ID. Each line in these files corresponds to a gene. The first panel is the expression profile of every gene in the developmental time course, the middle panel has the GO process terms with the red blocks showing if a gene is a member of a particular process. Gray lines indicate that the gene has no annotated process. The last panel shows the motifs associated with each cluster, with blue blocks indicating which genes have an instance for a particular motif.

**Data S15: Function prediction of all genes.** This zipped directory has one 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), 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 are included.

**Data S16: Genes predicted as regulators of key developmental events.** The first column specifies whether the prediction is based on the physical or functional network. The second column gives the sample at which the expression split identified by DREM begins to occur. The third column indicates the ID of the associated split. The fourth column indicates the path out of the split ('H' for higher, 'M' for middle, 'L' for lower). The fifth column gives the predicted TF. The sixth column gives the uncorrected p-value for target gene enrichments of the TF. The seventh column gives the % of target genes on the path among those going into the split.

**Data S17: Gene expression model based on regulator binding.** This zipped file has the four input regulatory networks, and the models learned for every fold of cross validation. The included README has more details.

## **Supporting References**

- S1. D. L. Wheeler *et al.*, *Nucleic Acids Res.* **36**, D13 (2008).
- S2. R. J. Wilson, J. L. Goodman, V. B. Strelets, *Nucleic Acids Res.* **36**, D588 (2008).
- S3. M. D. Adams *et al.*, *Science* **287**, 2185 (2000).
- S4. L. Florea, G. Hartzell, Z. Zhang, G. M. Rubin, W. Miller, *Genome Res.* **8**, 967 (1998).
- S5. B. R. Graveley *et al.*, *Nature*, 10.1038/nature09715.
- S6. C. Trapnell, L. Pachter, S. L. Salzberg, *Bioinformatics* **25**, 1105 (2009).
- S7. C. Trapnell *et al.*, *Nat. Biotechnol.* **28**, 511 (2010).
- S8. M. Blanchette *et al.*, *Genome Res.* **14**, 708 (2004).
- S9. A. Stark *et al.*, *Nature* **450**, 219 (2007).
- S10. M. F. Lin, I. Jungreis, M. Kellis, *Nat. Precedings*, <http://hdl.handle.net/10101/npre.2010.4784.1> (2010).
- S11. S. Washietl, I. L. Hofacker, P. F. Stadler, *Proc. Natl. Acad. Sci. U.S.A.* **102**, 2454 (2005).
- S12. A. R. Gruber, S. Findeiss, S. Washietl, I. L. Hofacker, P. F. Stadler, *Pac. Symp. Biocomput.* **15**, 69 (2010).
- S13. S. Will, K. Reiche, I. L. Hofacker, P. F. Stadler, R. Backofen, *PLoS Comput. Biol.* **3**, e65 (2007).
- S14. S. W. Park *et al.*, *Genetics* **177**, 1429 (2007).
- S15. S. C. Lakhotia, A. Sharma, *Genetica* **97**, 339 (1996).
- S16. R. Valgardsdottir *et al.*, *Mol. Biol. Cell* **16**, 2597 (2005).
- S17. S. Washietl, I. L. Hofacker, *J. Mol. Biol.* **342**, 19 (2004).
- S18. E. Berezikov *et al.*, *Nat. Genet.* **42**, 6 (2010).
- S19. W.-J. Chung *et al.*, *Genome Res.*, 10.1101/gr.113050.110.
- S20. N. Robine *et al.*, *Curr. Biol.* **19**, 2066 (2009).
- S21. K. Okamura, E. C. Lai, *Nat. Rev. Mol. Cell Biol.* **9**, 673 (2008).
- S22. X. Wang, Z. Xuan, X. Zhao, Y. Li, M. Q. Zhang, *Genome Res.* **19**, 266 (2009).
- S23. D. P. Fergenson *et al.*, *Anal. Chem.* **76**, 373 (2004).
- S24. T. Abeel, T. Helleputte, Y. Van de Peer, P. Dupont, Y. Saeys, *Bioinformatics* **26**, 392 (2009).
- S25. S. Nechaev *et al.*, *Science* **327**, 335 (2009).
- S26. P. V. Kharchenko *et al.*, *Nature*, 10.1038/nature09725.
- S27. J. Ernst, M. Kellis, *Nat. Biotechnol.* **28**, 817 (2010).
- S28. X. Y. Li *et al.*, *PLoS Biol.* **6**, e27 (2008).
- S29. S. MacArthur *et al.*, *Genome Biol.* **10**, R80 (2009).
- S30. R. P. Zinzen, C. Girardot, J. Gagneur, M. Braun, E. E. Furlong, *Nature* **462**, 65 (2009).
- S31. TF binding, hotspots, TF motif instances, promoter and enhancer validations, 12-point expression, and chromatin time course are available at [www.cistrack.org](http://www.cistrack.org).
- S32. A. Siepel *et al.*, *Genome Res.* **15**, 1034 (2005).
- S33. R. A. Hoskins *et al.*, *Genome Res.*, 10.1101/gr.12466.110 (2010).
- S34. B. Z. Shilo, *Exp. Cell. Res.* **284**, 140 (2003).
- S35. N. C. Grieder, D. Nellen, R. Burke, K. Basler, M. Affolter, *Cell* **81**, 791 (1995).

- S36. U. Alon, *Nat. Rev. Genet.* **8**, 450 (2007).
- S37. S. Wernicke, F. Rasche, *Bioinformatics* **22**, 1152 (2006).
- S38. F. Biemar *et al.*, *Proc. Natl. Acad. Sci. U.S.A.* **102**, 15907 (2005).
- S39. N. S. Sokol, V. Ambros, *Genes Dev.* **19**, 2343 (2005).
- S40. P. Kheradpour, A. Stark, S. Roy, M. Kellis, *Genome Res.* **17**, 1919 (2007).
- S41. S. D. Hooper *et al.*, *Mol. Syst. Biol.* **3**, 72 (2007).
- S42. V. R. Chintapalli, J. Wang, J. A. Dow, *Nat. Genet.* **39**, 715 (2007).
- S43. S. M. Gallo, L. Li, Z. Hu, M. S. Halfon, *Bioinformatics* **22**, 381 (2006).
- S44. R. Polikar, A. Topalis, D. Green, J. Kounios, C. M. Clark, *Comput. Biol. Med.* **37**, 542 (2007).
- S45. P. Tomancak *et al.*, *Genome Biol.* **8**, R145 (2007).
- S46. J. Ernst, O. Vainas, C. T. Harbison, I. Simon, Z. Bar-Joseph, *Mol. Syst. Biol.* **3**, 74 (2007).
- S47. V. Matys *et al.*, *Nucleic Acids Res.* **31**, 374 (2003).
- S48. W. W. Wasserman, A. Sandelin, *Nat. Rev. Genet.* **5**, 276 (2004).
- S49. T. A. Down, C. M. Bergman, J. Su, T. J. Hubbard, *PLoS Comput. Biol.* **3**, e7 (2007).
- S50. A. Ivan, M. S. Halfon, S. Sinha, *Genome Biol.* **9**, R22 (2008).
- S51. M. B. Noyes *et al.*, *Cell* **133**, 1277 (2008).
- S52. M. B. Noyes *et al.*, *Nucleic Acids Res.* **36**, 2547 (2008).
- S53. D. E. Reed, X. M. Huang, J. A. Wohlschlegel, M. S. Levine, K. Senger, *Proc. Natl. Acad. Sci. U.S.A.* **105**, 8351 (2008).
- S54. A. Sen, B. G. Stultz, H. Lee, D. A. Hursh, *Dev. Biol.* **343**, 167 (2010).
- S55. H. Touzet, J. S. Varre, *Algorithms Mol. Biol.* **2**, 15 (2007).
- S56. Z. Misulovin *et al.*, *Chromosoma* **117**, 89 (2008).
- S57. N. C. Riddle *et al.*, *Genome Res.*, 10.1101/gr.110098.110.
- S58. H. K. MacAlpine, R. Gordan, S. K. Powell, A. J. Hartemink, D. M. MacAlpine, *Genome Res.* **20**, 201 (2010).
- S59. S. Henikoff, J. G. Henikoff, A. Sakai, G. B. Loeb, K. Ahmad, *Genome Res.* **19**, 460 (2009).
- S60. R. B. Deal, J. G. Henikoff, S. Henikoff, *Science* **328**, 1161 (2010).
- S61. R. A. Hoskins *et al.*, *Science* **316**, 1625 (2007).
- S62. C. D. Smith, S. Shu, C. J. Mungall, G. H. Karpen, *Science* **316**, 1586 (2007).
- S63. G. J. Filion *et al.*, *Cell* **143**, 212 (2010).
- S64. E. Birney *et al.*, *Nature* **447**, 799 (2007).
- S65. E. Berezikov *et al.*, *Genome Res.*, 10.1101/gr.116657.110.

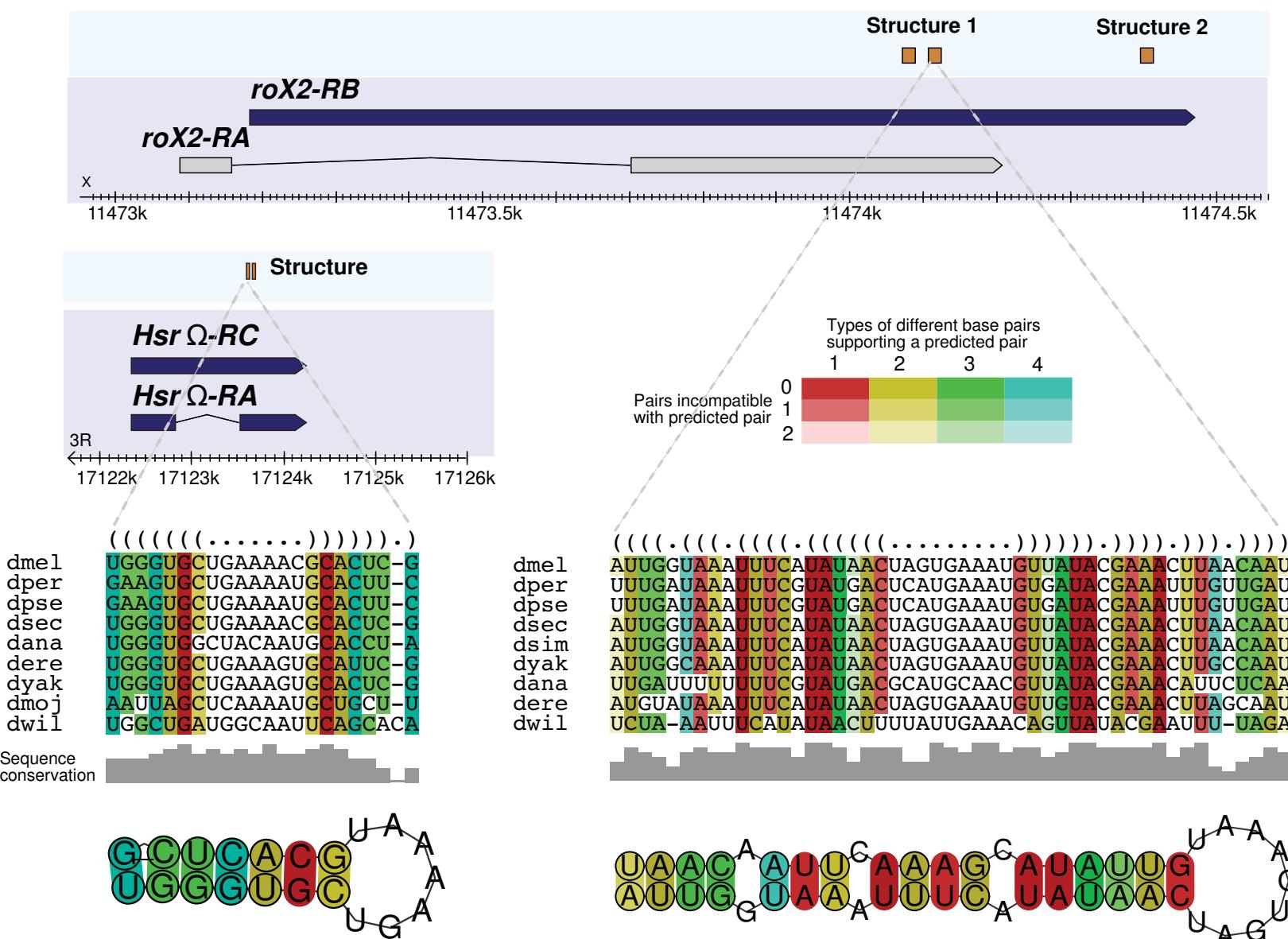


Figure S1: Conserved RNA secondary structures.

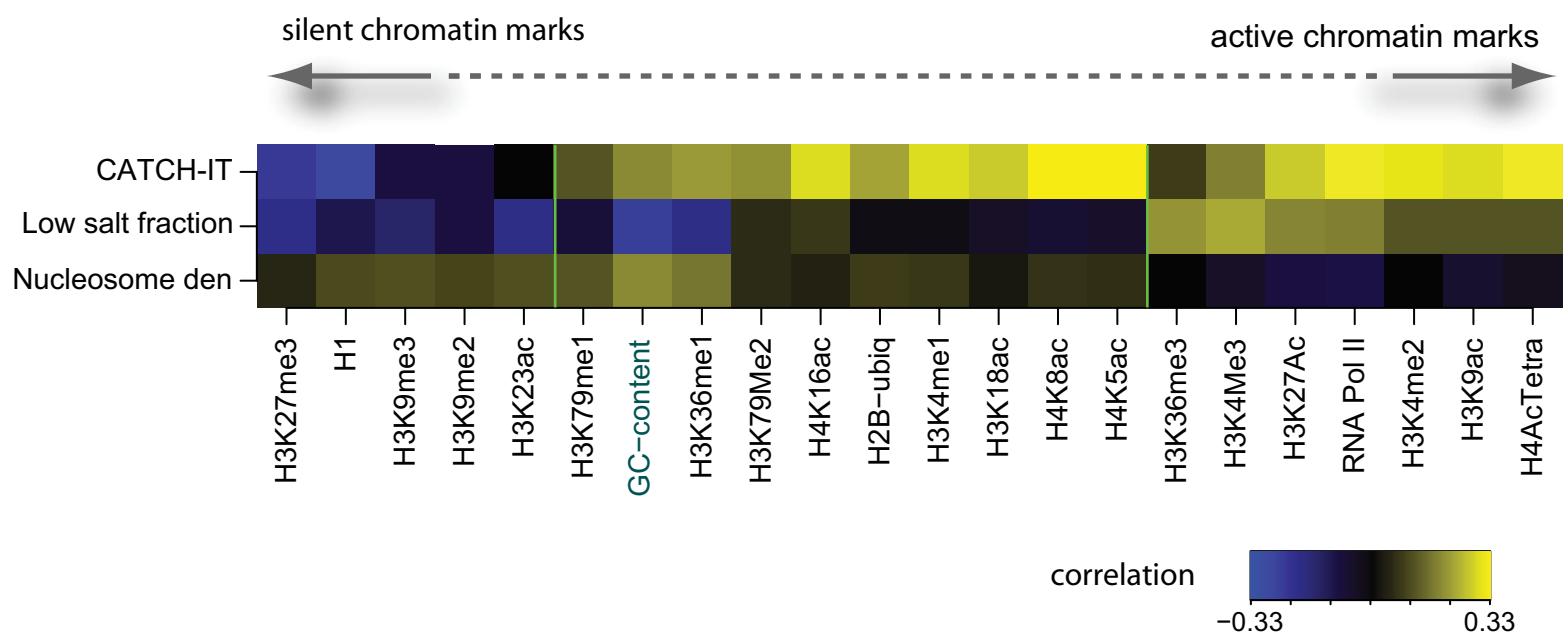


Figure S2 - Genome-wide correlation of chromatin marks with patterns of physical properties

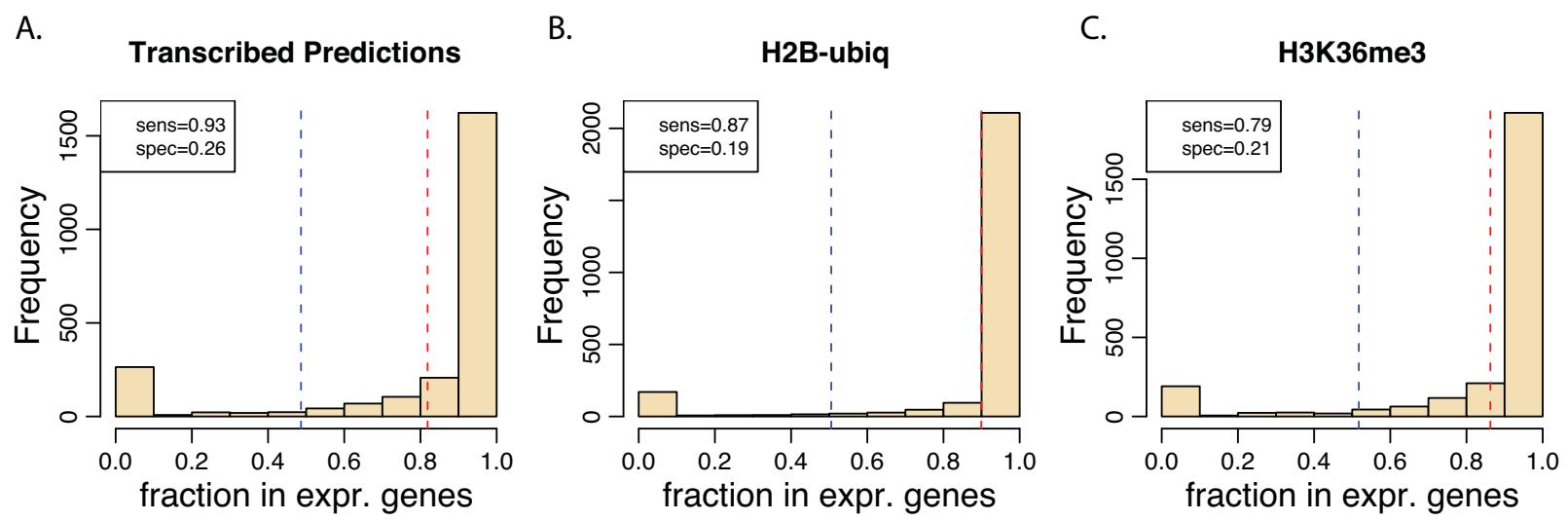


Figure S3 - Detecting transcribed regions based on enrichment of histone marks.

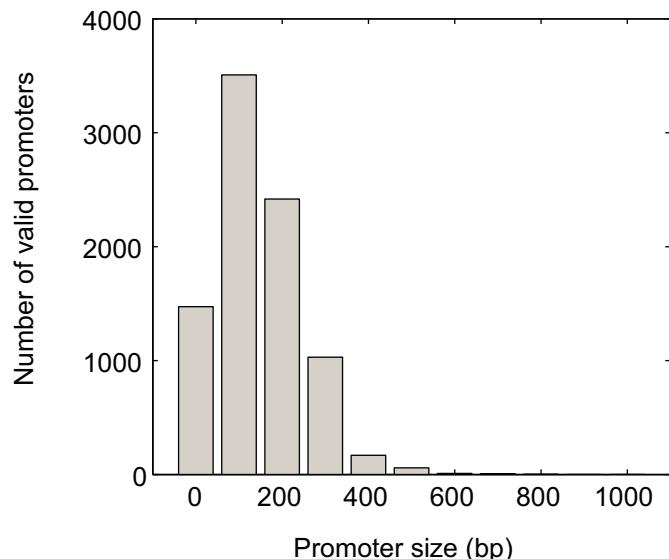
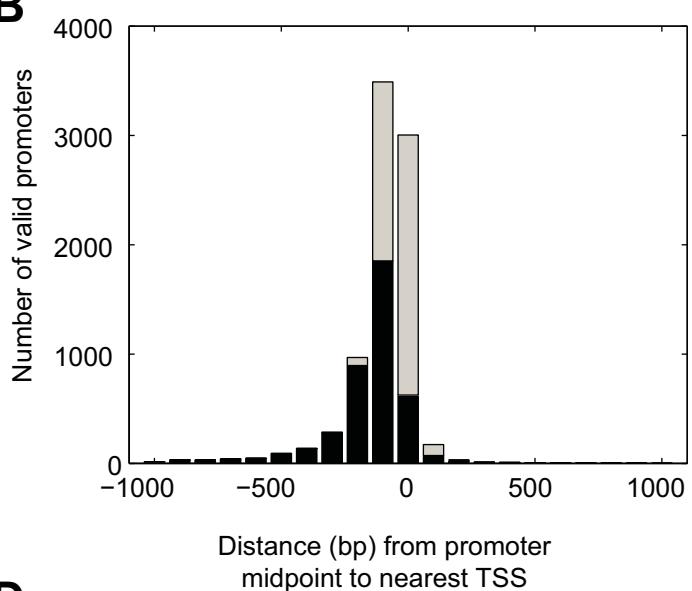
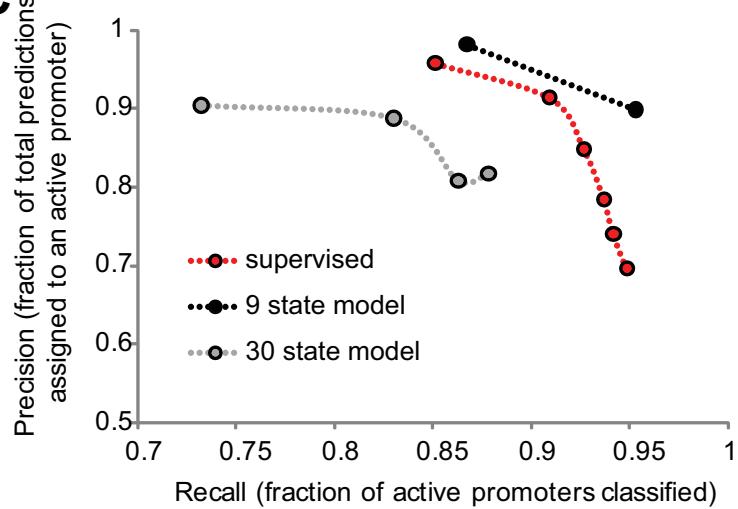
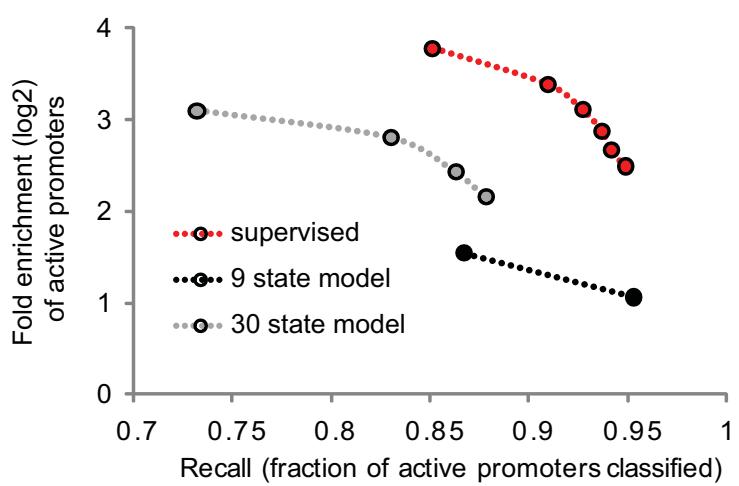
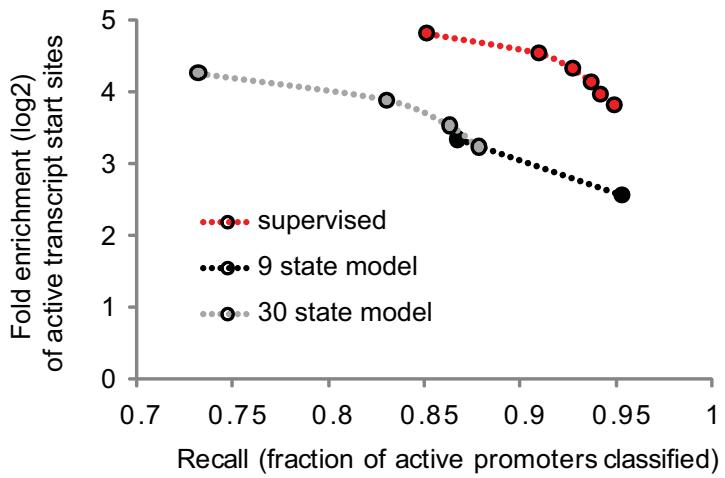
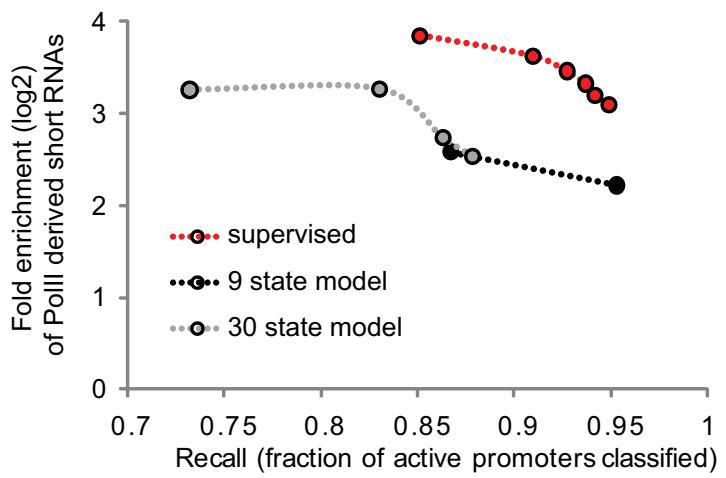
**A****B****C****D****E****F**

Figure S4 - Predictive model of active promoters in S2 cells.

### Replication and cohesin activity at >10kb intergenic H3K36me1 regions

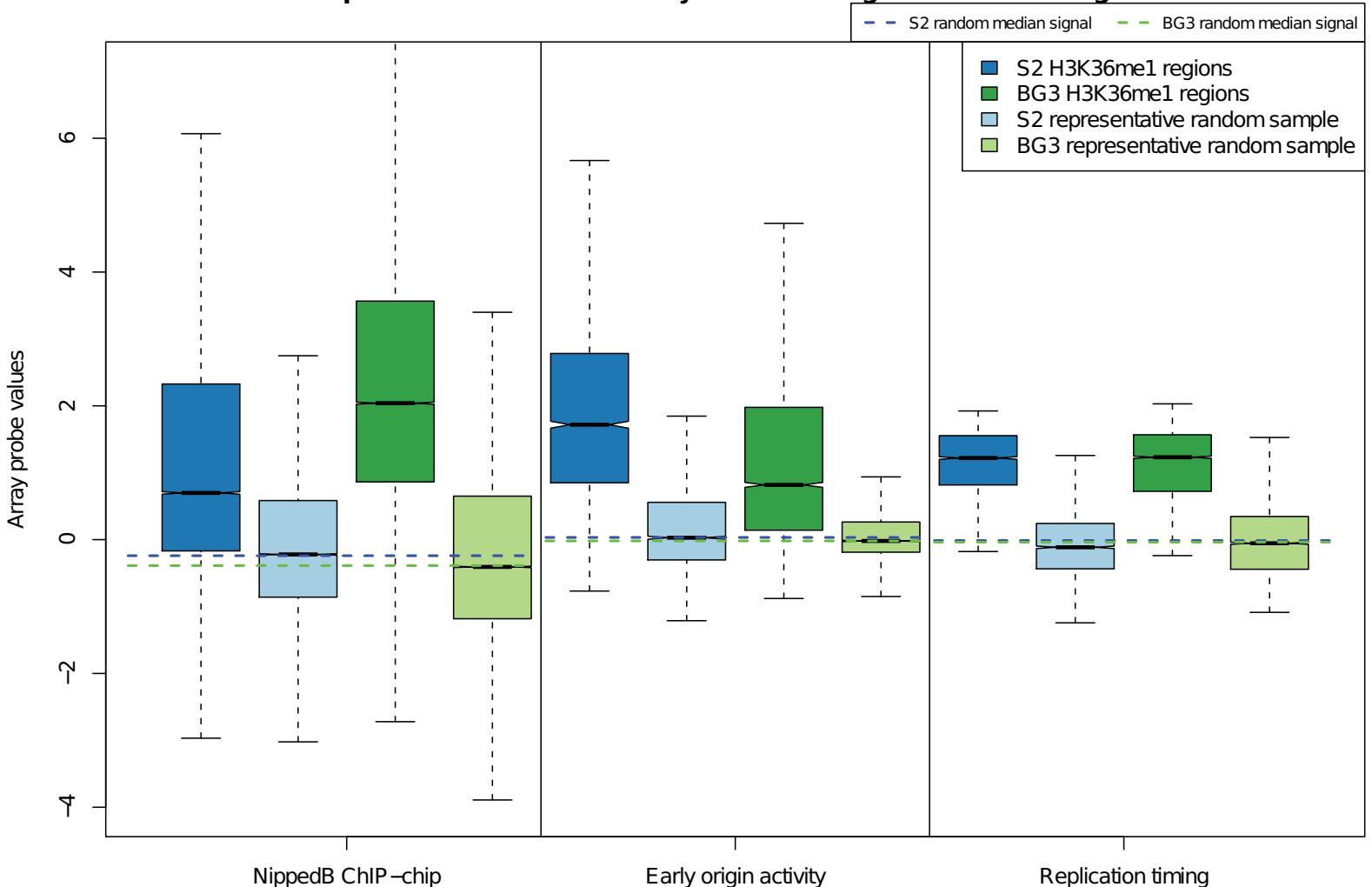


Figure S5 - Early origin activity, replication timing and cohesin binding in H3K36me1 intergenic regions.

A

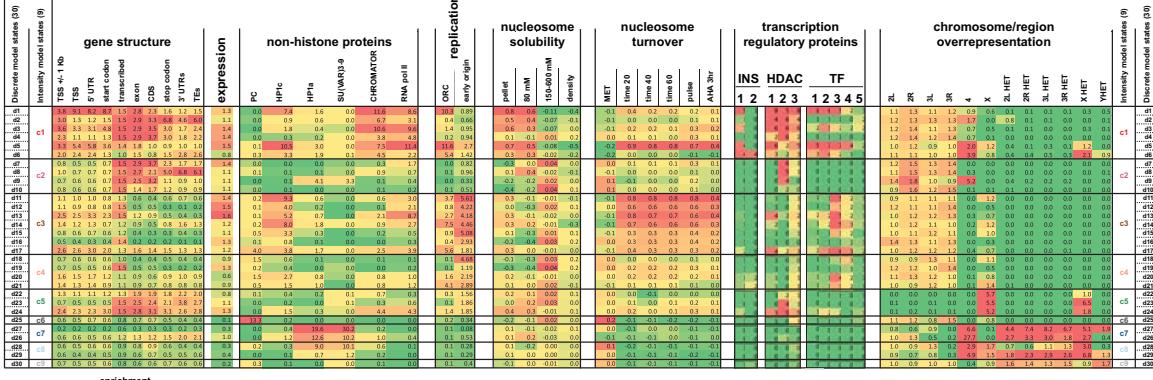
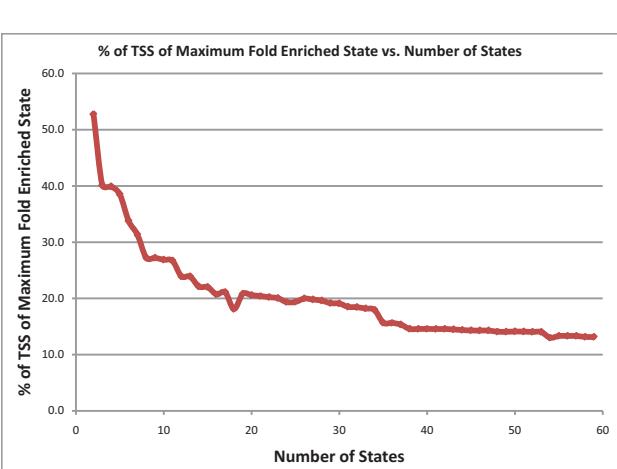
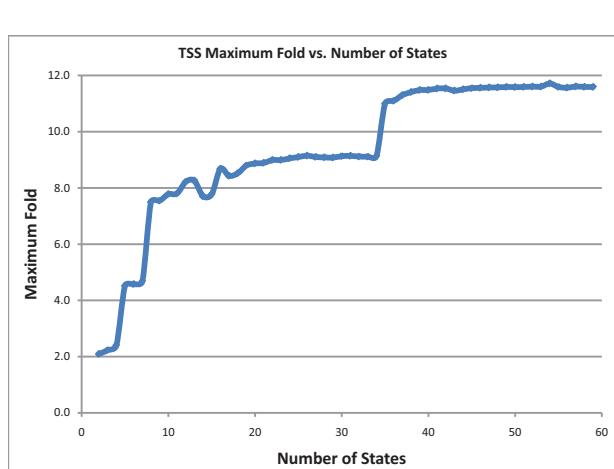


Figure S7 - Maximum chromatin state enrichments with increasing numbers of states

## a Maximum TSS Enriched State for Increasing Number of States

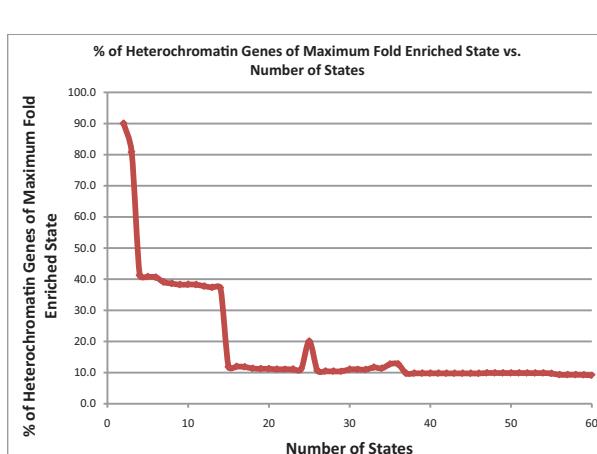
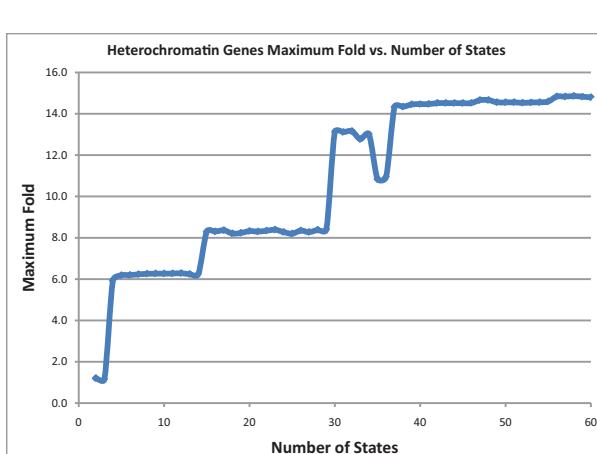
# States	H3K36ME3	H3K79ME1	H2B-ubiq	H3K79ME2	H3K4ME2	H3K4ME3	H3K9AC	H4K16AC	H3K4ME1	H3K36ME1	H3K18AC	H3K27AC	H1 depletion	H4 depletion	H3K23AC depletion	H3K9ME3	H3K9ME2	H3K27ME3	Max Fold	% of TSS
2	25	13	33	22	36	30	26	23	19	25	25	19	1	5	2	1	1	1	2.1	52.8
3	45	17	27	16	31	33	18	27	4	0	0	3	0	5	4	3	3	0	2.2	40.1
4	44	16	28	17	33	35	19	28	5	0	0	1	4	0	5	2	0	1	2.4	39.9
5	28	5	25	29	77	83	54	28	8	4	14	21	1	8	3	0	0	0	4.5	38.5
6	33	7	23	27	74	81	46	26	6	0	2	10	1	7	3	0	0	0	4.6	33.8
7	34	7	25	29	76	88	49	26	6	1	2	10	1	7	3	0	0	0	4.7	31.3
8	1	0	1	1	46	66	48	13	2	0	5	17	2	16	8	0	0	0	7.5	27.2
9	1	0	1	1	47	67	49	13	2	0	5	17	1	16	8	0	0	0	7.5	27.2
10	2	0	1	2	49	73	52	15	1	0	4	14	1	15	8	0	0	0	7.8	26.9
11	2	0	1	1	47	71	50	14	1	0	3	13	1	16	8	0	0	0	7.8	26.7
12	3	0	1	2	52	86	56	17	0	0	1	11	1	13	6	0	0	0	8.2	23.9
13	2	0	1	2	52	85	56	17	0	0	1	11	1	13	7	0	0	0	8.3	23.9
14	16	0	1	3	58	88	50	19	0	0	0	2	1	9	5	1	0	0	7.7	22.1
15	15	0	1	3	58	88	50	19	0	0	0	2	1	9	5	0	0	0	7.8	22.0
16	2	0	1	2	47	81	52	4	0	0	1	8	1	14	8	0	0	0	8.7	20.7
17	1	0	1	3	49	82	52	4	1	0	1	7	1	14	7	0	0	0	8.4	21.1
18	18	0	1	3	57	95	54	6	0	0	0	3	1	9	5	0	0	0	8.5	18.1
19	1	0	1	3	46	86	55	5	0	0	0	7	1	15	8	0	0	0	8.8	20.8
20	1	0	1	3	46	87	56	6	0	0	0	7	1	14	8	0	0	0	8.9	20.6
21	1	0	1	3	47	88	56	6	0	0	0	7	1	14	8	0	0	0	8.9	20.4
22	1	0	1	2	46	89	56	6	0	0	0	7	1	14	8	0	0	0	9.0	20.2
23	1	0	1	2	46	88	56	6	0	0	0	7	1	14	8	0	0	0	9.0	20.1
24	1	0	1	3	47	93	57	6	0	0	0	7	1	14	7	0	0	0	9.1	19.4
25	1	0	1	3	48	94	58	7	0	0	0	7	1	14	7	0	0	0	9.1	19.4
26	1	0	1	2	47	92	58	6	0	0	0	8	1	14	7	0	0	0	9.1	20.0
27	1	0	1	2	47	91	57	6	0	0	0	6	1	12	6	0	0	0	9.1	19.8
28	1	0	1	3	47	91	57	7	0	0	0	5	1	12	6	0	0	0	9.1	19.6
29	1	0	1	3	47	91	56	7	0	0	0	3	1	11	6	0	0	0	9.1	19.2
30	1	0	1	3	47	92	57	7	0	0	0	3	1	12	7	0	0	0	9.1	19.1
31	1	0	1	3	47	92	56	4	0	0	0	4	1	13	7	0	0	0	9.1	18.5
32	1	0	1	3	47	92	56	4	0	0	0	4	1	13	7	0	0	0	9.1	18.4
33	1	0	1	3	47	92	55	4	0	0	0	3	1	12	7	0	0	0	9.1	18.2
34	1	0	1	3	47	93	56	4	0	0	0	3	1	12	7	0	0	0	9.1	18.0
35	0	0	1	1	3	78	54	3	0	0	0	6	2	18	12	0	0	0	11.0	15.7
36	0	0	1	1	4	80	55	4	0	0	0	6	2	19	12	0	0	0	11.1	15.7
37	0	0	1	1	4	83	56	4	0	0	0	6	2	17	10	0	0	0	11.3	15.4
38	0	0	1	1	2	83	54	3	0	0	0	2	1	16	10	0	0	0	11.4	14.6
39	0	0	1	1	2	85	56	4	0	0	0	2	1	16	10	0	0	0	11.5	14.6
40	0	0	1	1	2	85	56	4	0	0	0	2	1	16	10	0	0	0	11.5	14.6
41	0	0	1	1	2	85	56	4	0	0	0	2	1	16	10	0	0	0	11.5	14.6
42	0	0	1	1	2	85	56	4	0	0	0	2	1	16	10	0	0	0	11.5	14.6
43	0	0	1	1	3	85	56	3	0	0	0	2	1	16	10	0	0	0	11.5	14.5
44	0	0	1	1	3	85	56	3	0	0	0	1	1	16	10	0	0	0	11.5	14.4
45	0	0	1	1	2	85	56	3	0	0	0	1	1	16	10	0	0	0	11.5	14.3
46	0	0	1	1	2	85	56	3	0	0	0	1	1	16	10	0	0	0	11.6	14.3
47	0	0	1	1	2	85	56	3	0	0	0	1	1	16	10	0	0	0	11.6	14.3
48	0	0	1	1	2	85	56	3	0	0	0	1	1	16	10	0	0	0	11.6	14.1
49	0	0	1	1	2	85	56	3	0	0	0	1	1	16	10	0	0	0	11.6	14.1
50	0	0	1	1	2	85	56	3	0	0	0	1	1	16	10	0	0	0	11.6	14.1
51	0	0	1	1	2	85	56	2	0	0	0	2	1	16	10	0	0	0	11.6	14.1
52	0	0	1	1	2	85	56	2	0	0	0	2	1	16	10	0	0	0	11.6	14.1
53	0	0	1	1	2	85	56	2	0	0	0	2	1	16	10	0	0	0	11.6	14.0
54	0	0	1	1	0	86	54	2	0	0	0	0	1	15	10	0	0	0	11.7	13.0
55	0	0	1	1	0	84	53	2	0	0	0	0	1	15	9	0	0	0	11.6	13.3
56	0	0	1	1	0	83	53	2	0	0	0	0	1	14	9	0	0	0	11.6	13.3
57	0	0	1	1	0	84	53	2	0	0	0	0	1	15	9	0	0	0	11.6	13.3
58	0	0	1	1	0	84	53	2	0	0	0	0	1	14	9	0	0	0	11.6	13.2
59	0	0	1	1	0	84	53	2	0	0	0	0	1	14	9	0	0	0	11.6	13.2
60	0	0	1	1	0	84	53	2	0	0	0	0	1	14	9	0	0	0	11.6	13.1



b

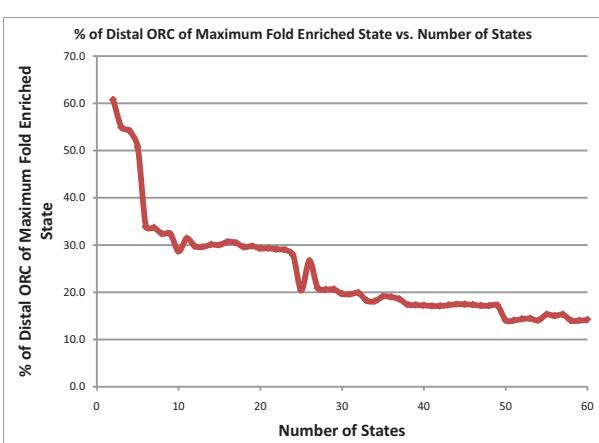
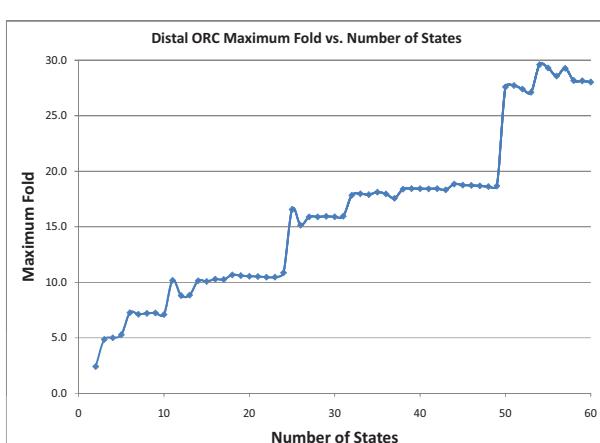
# Maximum Heterochromatin Gene Enriched State for Increasing Number of States

# of States	H3K36ME3	H3K79ME1	H2B-ubiq	H3K79ME2	H3K4ME2	H3K4ME3	H3K9AC	H4K16AC	H3K4ME1	H3K36ME1	H3K18AC	H3K27AC	H1 depletion	H4 depletion	H3K23AC depletion	H3K9ME3	H3K9ME2	H3K27ME3	Fold	% of all Het Genes	
2	3	1	0	0	0	0	0	2	1	1	0	0	0	1	4	5	7	6	1.2	90.0	
3	0	0	0	0	0	0	0	0	1	1	0	0	0	1	4	5	7	6	1.1	80.8	
4	9	4	3	0	0	0	0	0	2	0	0	0	0	4	43	53	74	0	5.9	41.3	
5	8	3	2	0	0	0	0	0	1	0	0	0	0	5	44	54	74	1	6.2	40.8	
6	8	3	2	0	0	0	0	0	1	0	0	0	0	5	44	54	74	1	6.2	40.6	
7	7	2	1	0	0	0	0	0	1	0	0	0	0	4	44	55	75	1	6.2	39.0	
8	7	2	1	0	0	0	0	0	1	0	0	0	0	4	44	55	75	1	6.3	38.6	
9	7	2	1	0	0	0	0	0	1	0	0	0	0	4	44	56	76	0	6.3	38.2	
10	7	2	1	0	0	0	0	0	1	0	0	0	0	4	44	56	76	0	6.3	38.3	
11	7	2	1	0	0	0	0	0	1	0	0	0	0	4	44	56	76	0	6.3	38.2	
12	7	2	1	0	0	0	0	0	1	0	0	0	0	4	44	56	77	0	6.3	37.7	
13	7	2	1	0	0	0	0	0	1	0	0	0	0	4	44	56	77	0	6.2	37.3	
14	7	2	1	0	0	0	0	0	1	0	0	0	0	4	44	56	77	0	6.2	37.1	
15	74	31	24	5	1	1	0	0	0	0	0	0	0	7	58	57	53	0	8.3	11.9	
16	74	30	23	5	1	1	0	0	0	0	0	0	0	7	58	56	52	0	8.3	11.9	
17	75	31	24	5	1	1	0	0	0	0	0	0	0	7	59	57	53	0	8.4	11.8	
18	76	31	24	5	1	1	0	0	0	0	0	0	0	6	58	58	54	0	8.2	11.3	
19	77	32	25	5	1	1	0	1	0	0	0	0	0	7	59	58	53	0	8.2	11.2	
20	78	31	24	5	1	1	0	1	0	0	0	0	0	7	59	58	53	0	8.3	11.2	
21	79	31	25	5	1	1	0	1	0	0	0	0	0	7	59	57	52	0	8.3	11.1	
22	79	31	24	5	1	1	0	1	0	0	0	0	0	7	59	58	52	0	8.3	11.1	
23	79	31	24	5	1	1	0	1	0	0	0	0	0	8	60	58	52	0	8.4	11.1	
24	78	32	24	5	1	1	0	1	0	0	0	0	0	7	59	57	52	0	8.3	11.1	
25	1	0	0	0	0	0	0	0	0	0	0	0	0	6	96	79	88	0	8.2	20.1	
26	80	31	25	5	1	1	0	1	0	0	0	0	0	7	60	59	53	0	8.3	10.7	
27	81	32	25	5	1	1	0	1	0	0	0	0	0	6	60	60	54	0	8.3	10.4	
28	81	31	25	5	1	1	0	1	0	0	0	0	0	6	60	60	54	0	8.4	10.4	
29	81	31	25	5	1	1	0	1	0	0	0	0	0	6	60	60	54	0	8.4	10.4	
30	61	5	4	0	1	1	0	1	0	0	0	0	0	19	82	55	35	0	13.1	11.0	
31	61	5	4	0	1	1	0	1	0	0	0	0	0	19	82	54	35	0	13.1	11.0	
32	62	5	4	0	1	1	0	1	0	0	0	0	0	18	82	55	35	0	13.2	11.0	
33	54	4	3	0	1	1	0	1	0	0	0	0	0	21	82	51	34	0	12.8	11.7	
34	57	4	4	0	1	1	0	1	0	0	0	0	0	20	82	53	34	0	13.0	11.4	
35	34	1	1	0	1	1	0	0	0	0	0	0	0	1	28	81	37	26	0	10.8	12.7
36	35	1	1	0	1	1	0	0	0	0	0	0	0	1	28	81	38	27	0	10.9	12.7
37	79	9	8	0	2	2	0	1	0	0	0	0	0	11	82	62	36	0	14.3	9.8	
38	79	9	8	0	2	2	0	1	0	0	0	0	0	11	82	62	36	0	14.3	9.8	
39	78	9	8	0	2	2	0	1	0	0	0	0	0	11	82	64	34	0	14.5	9.8	
40	78	9	8	0	2	2	0	1	0	0	0	0	0	11	82	64	34	0	14.5	9.8	
41	78	9	8	0	2	2	0	1	0	0	0	0	0	11	82	64	34	0	14.5	9.8	
42	78	8	8	0	2	2	0	1	0	0	0	0	0	12	82	64	34	0	14.5	9.8	
43	78	8	8	0	2	2	0	1	0	0	0	0	0	12	82	64	34	0	14.5	9.7	
44	78	8	8	0	2	2	0	1	0	0	0	0	0	12	82	64	34	0	14.5	9.8	
45	78	8	8	0	2	2	0	1	0	0	0	0	0	12	82	64	34	0	14.5	9.7	
46	78	8	8	0	2	2	0	1	0	0	0	0	0	11	82	64	34	0	14.5	9.7	
47	78	9	8	1	2	2	0	1	0	0	0	0	0	12	82	63	34	0	14.7	9.9	
48	78	9	8	1	2	2	0	1	0	0	0	0	0	12	82	63	34	0	14.7	9.9	
49	78	9	8	1	2	2	0	1	0	0	0	0	0	12	82	63	34	0	14.6	9.9	
50	78	9	8	1	2	2	0	1	0	0	0	0	0	12	82	63	34	0	14.6	9.9	
51	78	9	8	1	2	2	0	1	0	0	0	0	0	12	82	63	34	0	14.6	9.9	
52	77	8	8	0	2	2	0	1	0	0	0	0	0	12	82	63	34	0	14.5	9.8	
53	77	8	8	0	2	2	0	1	0	0	0	0	0	12	82	63	34	0	14.5	9.8	
54	78	9	8	1	2	2	0	1	0	0	0	0	0	12	82	63	34	0	14.6	9.9	
55	78	9	8	0	2	2	0	1	0	0	0	0	0	11	82	63	34	0	14.6	9.7	
56	87	10	9	1	2	1	0	1	0	0	0	0	0	11	84	62	37	0	14.8	9.3	
57	87	10	9	1	2	1	0	1	0	0	0	0	0	11	84	62	37	0	14.8	9.3	
58	87	10	9	1	2	1	0	1	0	0	0	0	0	11	84	62	37	0	14.9	9.3	
59	87	10	9	1	2	1	0	1	0	0	0	0	0	11	84	62	37	0	14.8	9.3	
60	87	10	9	0	2	1	0	1	0	0	0	0	0	11	84	62	37	0	14.8	9.1	



**C Maximum Distal ORC Enriched State for Increasing Number of States**

# of States	H3K36ME3	H3K79ME1	H2B-ubiq	H3K79ME2	H3K4ME2	H3K4ME3	H3K9AC	H4K16AC	H3K4ME1	H3K36ME1	H3K18AC	H3K27AC	H1 depletion	H4 depletion	H3K23AC depletion	H3K9ME3	H3K9ME2	H3K27ME3	Max Fold	% of Distal ORC
2	25	13	33	22	36	30	26	23	19	25	25	19	1	5	2	1	1	1	2.4	60.7
3	1	6	29	22	30	13	28	18	35	55	52	36	1	5	0	0	0	0	4.9	55.0
4	1	6	30	23	30	13	29	18	36	56	55	36	1	6	0	0	0	0	5.0	54.2
5	0	6	27	20	19	1	19	14	39	61	49	27	1	4	0	0	0	0	5.3	50.7
6	2	2	31	31	64	28	58	42	49	45	87	81	2	10	1	0	0	0	7.3	33.9
7	2	2	29	29	62	27	57	40	51	45	88	82	2	10	0	0	0	0	7.1	33.7
8	2	2	33	32	63	25	57	43	53	48	90	81	2	9	0	0	0	0	7.2	32.4
9	2	2	33	32	64	26	57	43	53	48	89	81	2	9	0	0	0	0	7.2	32.3
10	2	2	30	31	69	31	63	49	53	45	91	87	3	10	0	0	0	0	7.1	28.7
11	2	2	4	8	29	0	28	36	67	39	80	79	2	10	1	0	0	0	10.2	31.4
12	2	2	17	19	43	1	38	42	72	49	90	82	2	9	0	0	0	0	8.8	29.7
13	2	2	16	18	42	1	37	43	73	48	89	83	2	9	0	0	0	0	8.8	29.6
14	2	2	2	7	32	1	31	39	69	36	83	79	2	10	0	0	0	0	10.1	30.1
15	2	2	2	7	32	1	31	39	69	37	84	79	2	9	0	0	0	0	10.1	30.0
16	2	2	2	6	32	1	31	38	68	35	83	79	2	10	0	0	0	0	10.3	30.7
17	2	2	2	6	32	1	31	38	68	35	83	79	2	10	0	0	0	0	10.3	30.5
18	1	2	2	7	31	1	33	41	70	38	88	83	2	10	0	0	0	0	10.7	29.6
19	1	2	2	7	29	1	33	40	69	37	87	81	2	10	0	0	0	0	10.6	29.8
20	1	2	2	7	30	1	33	41	70	38	88	81	2	10	0	0	0	0	10.5	29.3
21	1	2	2	7	29	1	33	40	70	38	87	81	2	10	0	0	0	0	10.5	29.3
22	1	2	3	7	29	1	32	40	70	38	87	81	2	10	0	0	0	0	10.5	29.1
23	1	2	3	7	29	1	32	40	70	38	88	81	2	10	0	0	0	0	10.5	29.0
24	1	2	3	8	33	1	35	44	70	39	90	87	2	10	0	0	0	0	10.9	27.8
25	2	0	1	1	7	0	16	6	33	2	54	56	0	8	1	0	0	0	16.6	20.4
26	3	1	2	27	1	31	33	53	1	72	87	2	13	1	0	0	0	15.1	26.7	
27	2	0	1	1	11	0	17	9	39	1	57	66	1	10	1	0	0	0	15.9	20.9
28	2	0	1	1	10	0	17	9	39	2	57	65	1	10	1	0	0	0	15.9	20.6
29	2	1	1	1	11	0	16	9	40	1	57	66	1	10	1	0	0	0	15.9	20.6
30	3	1	1	1	13	0	15	11	42	2	56	75	1	11	1	0	0	0	15.9	19.6
31	3	1	1	1	12	0	15	11	42	2	57	75	1	11	1	0	0	0	16.0	19.6
32	0	0	1	1	12	0	16	8	39	2	59	72	1	10	1	0	0	0	17.8	19.9
33	0	0	1	1	12	0	11	6	41	3	56	80	1	11	1	0	0	0	18.0	18.2
34	0	0	1	1	13	0	12	6	41	3	57	81	1	10	1	0	0	0	17.9	18.1
35	0	0	1	1	12	0	12	5	40	3	56	75	1	10	2	0	0	0	18.1	19.1
36	0	0	1	1	10	0	12	6	41	3	57	74	1	10	1	0	0	0	18.0	19.0
37	0	0	1	1	10	0	13	7	42	3	58	76	1	10	1	0	0	0	17.6	18.5
38	0	0	1	1	10	0	10	4	41	4	55	76	1	10	1	0	0	0	18.4	17.4
39	0	0	1	1	10	0	10	4	41	4	56	77	1	10	1	0	0	0	18.4	17.3
40	0	0	1	1	10	0	10	3	41	4	56	77	1	10	1	0	0	0	18.4	17.2
41	0	0	1	1	11	0	10	3	41	4	56	76	1	10	1	0	0	0	18.4	17.1
42	0	0	1	1	11	0	10	3	41	4	56	77	1	10	1	0	0	0	18.4	17.1
43	0	0	1	1	11	0	11	3	41	4	56	76	1	10	1	0	0	0	18.3	17.3
44	0	0	1	0	2	0	13	3	41	2	55	72	1	10	1	0	0	0	18.8	17.5
45	0	0	1	1	1	0	13	2	41	2	54	70	1	9	1	0	0	0	18.8	17.4
46	0	0	1	1	1	0	13	2	41	2	55	71	1	9	1	0	0	0	18.7	17.4
47	0	0	1	1	1	0	13	2	42	2	55	72	1	9	1	0	0	0	18.7	17.1
48	0	0	1	1	1	0	13	2	42	2	55	72	1	9	1	0	0	0	18.6	17.1
49	0	0	1	1	1	0	13	2	42	2	55	72	1	9	1	0	0	0	18.7	17.2
50	0	0	1	0	1	0	9	2	7	2	37	63	1	10	2	0	0	0	27.6	14.1
51	0	0	1	0	1	0	9	2	7	2	37	64	1	10	2	0	0	0	27.7	14.0
52	0	0	1	0	1	0	10	2	6	2	37	61	1	10	2	0	0	0	27.4	14.3
53	0	0	0	0	1	0	11	2	5	2	39	59	1	10	2	0	0	0	27.1	14.4
54	0	0	0	0	1	0	4	1	10	2	34	64	1	11	2	0	0	0	29.6	14.0
55	0	0	0	0	1	0	3	1	6	2	30	56	1	9	2	0	0	0	29.3	15.3
56	0	0	0	0	1	0	4	1	6	2	31	57	0	9	2	0	0	0	28.6	15.0
57	0	0	0	0	1	0	3	1	6	2	30	56	1	9	2	0	0	0	29.3	15.3
58	0	0	0	0	1	0	6	1	7	2	36	64	0	10	2	0	0	0	28.2	14.0
59	0	0	0	0	1	0	6	1	7	2	36	64	0	10	2	0	0	0	28.1	14.0
60	0	0	1	0	2	0	6	2	8	2	36	64	1	10	2	0	0	0	28.0	14.1



		H3K36ME3	H3K79ME1	H2B-Ubiq	H3K79ME2	H3K4ME2	H3K4ME3	H3K9AC	H4K16AC	H3K4ME1	H3K36ME1	H3K18AC	H3K27AC	H1 depletion	H4 depletion	H3K23AC depletion	H3K9ME3	H3K9ME2	H3K27ME3	Active Yellow	Active Red	Polycomb Blue	HP1 Green	Repressive Black
State																								
1	1	1	0	1	3	47	92	57	7	0	0	0	3	1	12	7	0	0	0	5.0	1.2	0.1	0.5	0.0
2	95	20	10	10	79	93	24	9	1	0	0	0	0	0	0	0	1	0	0	5.2	0.2	0.1	1.3	0.1
3	52	3	55	79	99	100	92	45	7	0	1	13	1	2	1	0	0	0	0	5.5	0.5	0.0	0.6	0.0
4	57	22	73	77	93	64	5	7	23	1	0	4	0	1	0	1	1	1	0	5.3	0.5	0.1	1.4	0.0
5	2	0	8	11	78	87	92	39	4	1	59	89	4	22	3	0	0	0	0	2.3	6.1	0.3	0.5	0.0
6	0	0	0	0	1	0	1	0	0	0	1	1	1	24	11	0	0	0	0	2.1	1.3	0.9	1.3	0.7
7	73	44	88	37	0	1	0	1	4	0	1	1	0	0	0	0	0	0	0	4.6	0.3	0.1	3.3	0.1
8	82	14	2	1	0	0	0	1	2	0	0	0	0	3	0	0	0	0	0	4.3	0.4	0.4	1.7	0.2
9	73	67	54	14	1	0	0	0	1	1	0	0	0	1	26	50	77	0	0	1.5	0.2	0.2	11.4	0.4
10	1	37	34	2	0	0	0	0	0	2	0	0	0	0	0	0	0	1	1	1.6	0.4	1.2	1.0	1.0
11	2	2	7	14	63	7	64	82	84	46	98	98	4	12	0	0	0	0	0	0.5	8.9	0.7	0.0	0.1
12	2	2	88	69	79	2	47	23	55	87	84	59	0	3	0	0	0	0	0	0.7	8.3	0.9	0.0	0.0
13	4	1	79	73	100	94	87	32	24	66	83	73	1	5	0	0	0	0	0	1.7	7.3	0.5	0.1	0.0
14	3	1	1	1	13	0	15	11	42	2	56	75	1	11	1	0	0	0	0	0.8	7.3	1.0	0.1	0.1
15	0	8	3	17	8	0	13	15	63	93	81	26	0	3	0	0	0	1	0	0.2	6.3	2.1	0.1	0.1
16	0	5	88	64	3	0	30	3	15	95	84	3	0	1	0	0	0	0	0	0.2	4.6	2.7	0.0	0.1
17	3	2	4	4	92	12	13	2	33	12	9	11	1	4	1	1	0	13	0	2.1	4.4	1.2	0.3	0.1
18	0	10	2	2	0	0	3	1	6	88	21	0	0	1	0	0	0	0	7	0.1	2.3	3.6	0.1	0.3
19	0	15	84	36	2	0	3	3	7	72	4	0	0	0	0	0	0	0	0	0.5	2.2	3.0	0.0	0.4
20	1	4	0	1	1	0	3	1	89	11	10	2	0	5	0	1	0	3	0	0.6	3.7	2.1	0.1	0.5
21	0	0	0	0	0	0	5	0	1	3	7	0	0	0	0	0	0	0	0	0.7	2.8	2.4	0.1	0.5
22	1	1	0	0	0	1	0	78	3	2	1	0	0	2	0	1	3	1	0	2.1	0.7	1.0	1.2	0.8
23	65	21	29	5	3	1	0	99	5	1	1	1	0	4	1	1	2	0	0	3.1	0.5	0.6	3.8	0.4
24	28	5	21	11	88	58	26	98	11	1	2	7	0	3	0	0	0	0	0	5.1	0.8	0.1	0.9	0.1
25	0	0	0	0	0	0	1	0	2	0	0	0	0	0	0	0	1	0	81	0.1	0.2	3.8	0.0	0.7
26	0	0	0	0	0	0	0	0	0	0	0	0	0	4	96	82	91	0	0	0.0	0.2	0.4	3.1	0.3
27	61	5	4	0	1	1	0	1	0	0	0	0	0	19	82	55	35	0	0	0.4	0.1	0.2	10.9	0.3
28	0	0	0	0	0	0	0	3	0	0	0	0	0	0	3	91	85	2	0	0.3	0.1	0.8	2.1	1.4
29	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4	1	35	0	0	0.3	0.1	0.4	2.5	1.2
30	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.2	0.2	0.7	0.6	1.5

Figure S8 - Comparison of 30 chromatin states with 5 'colors' of chromatin

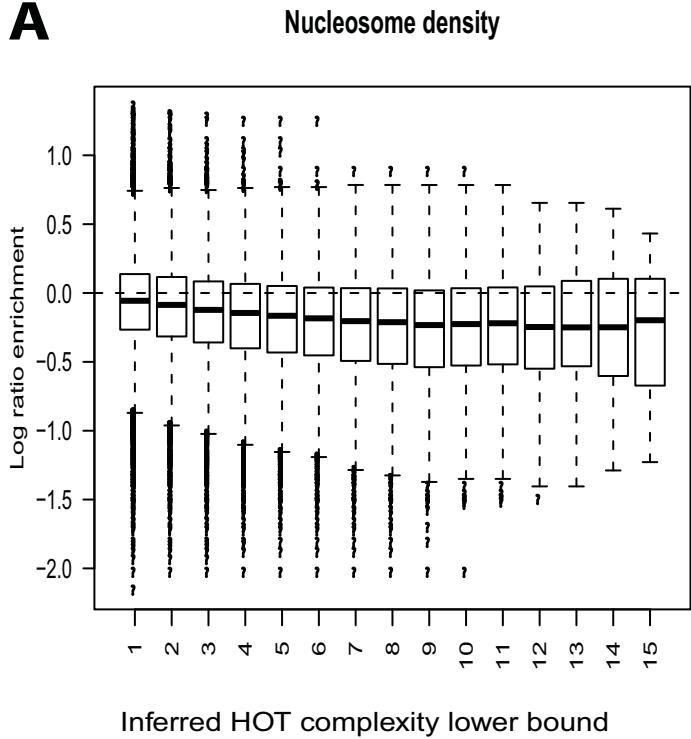
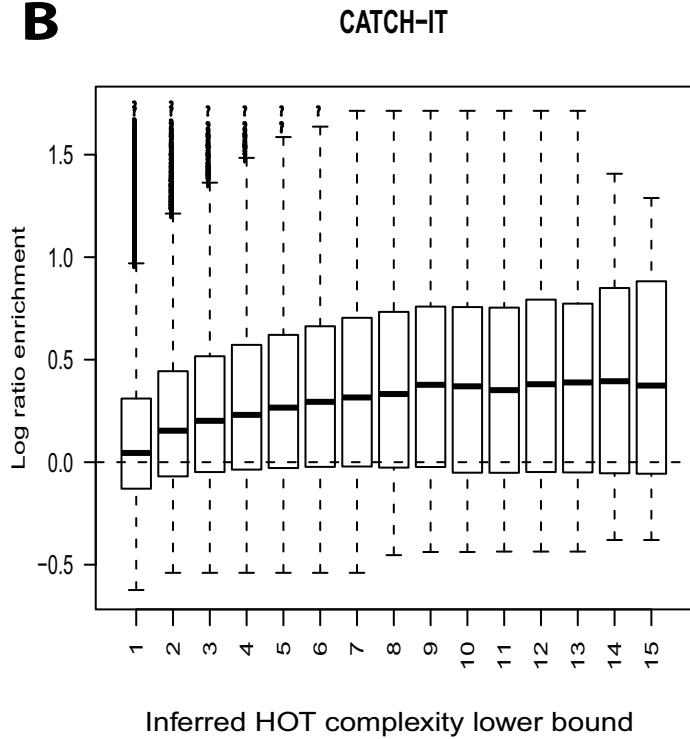
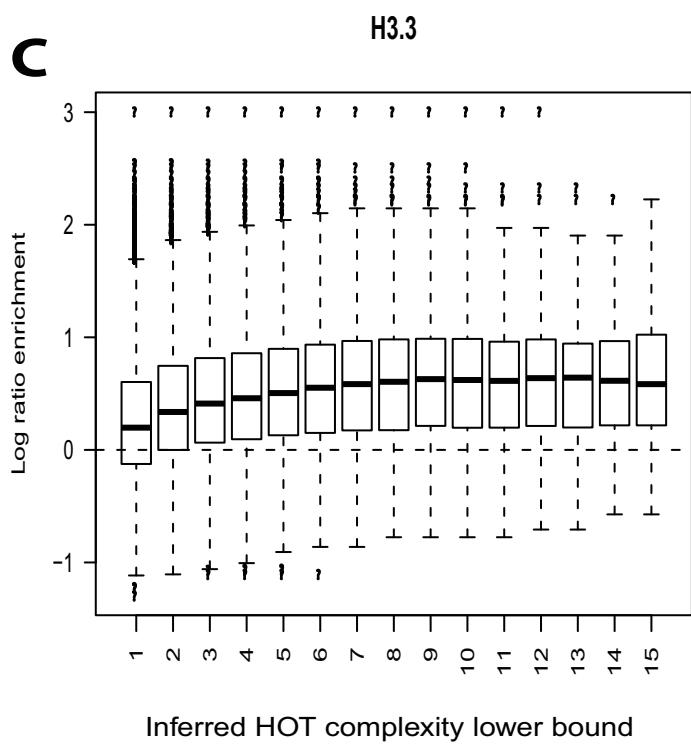
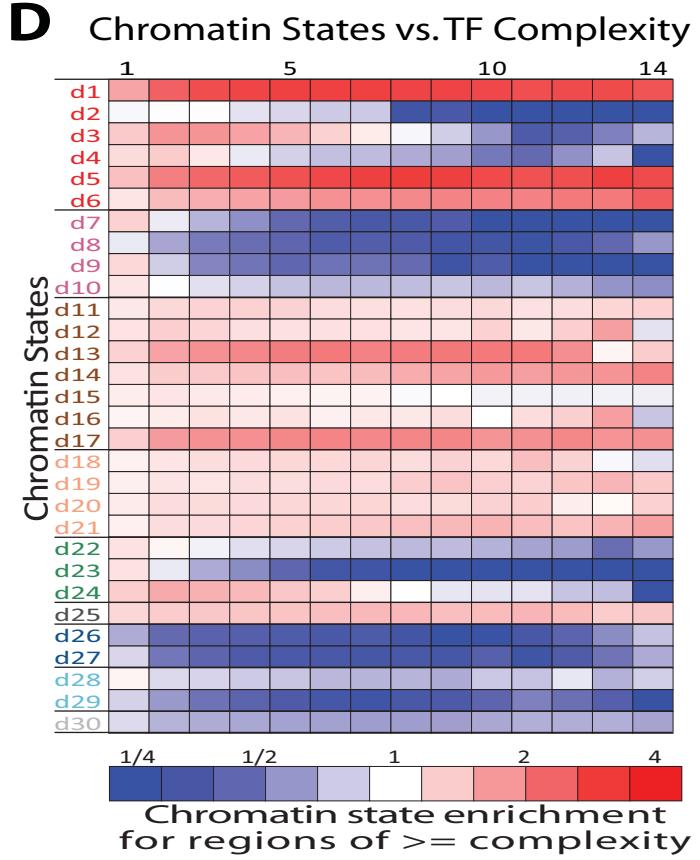
**A****B****C****D**

Figure S9 - Chromatin properties for HOT regions at varying complexity levels

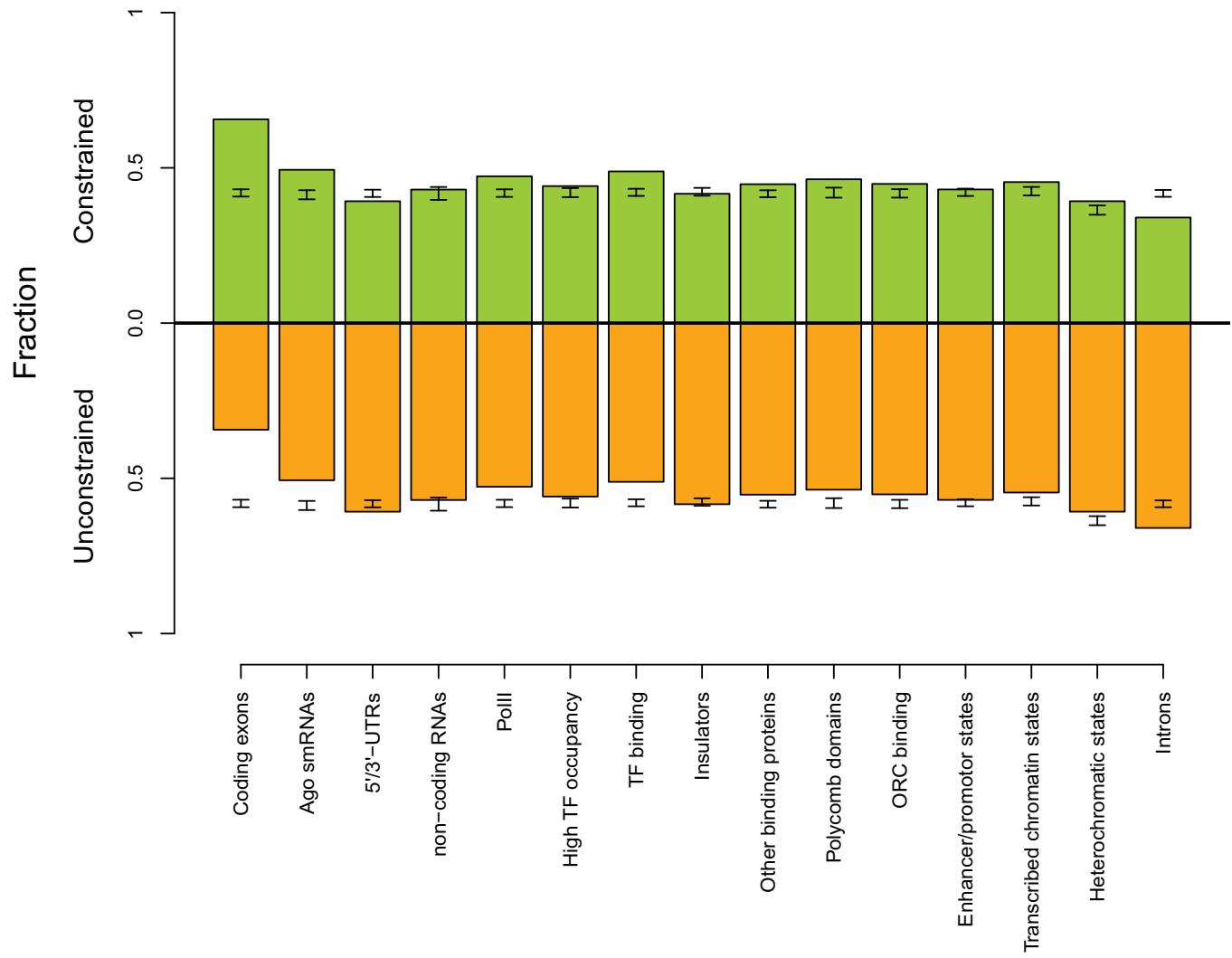
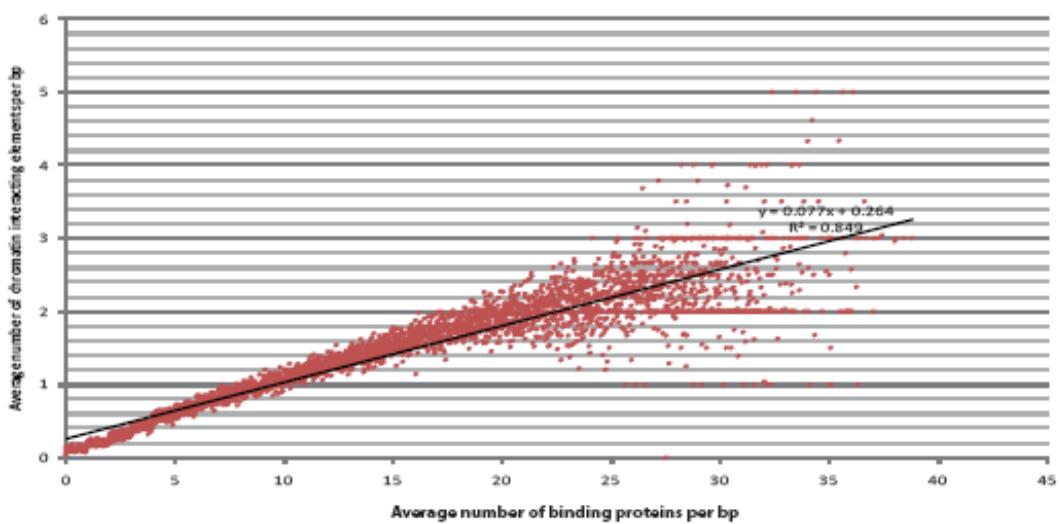
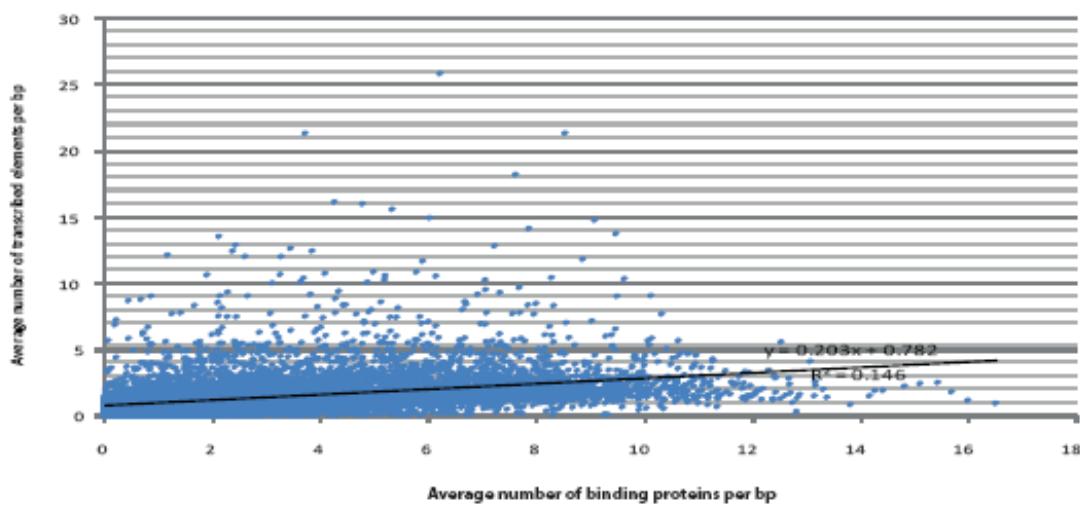


Figure S10 - Fraction of constrained and unconstrained bases in modENCODE datasets

A.



B.



C.

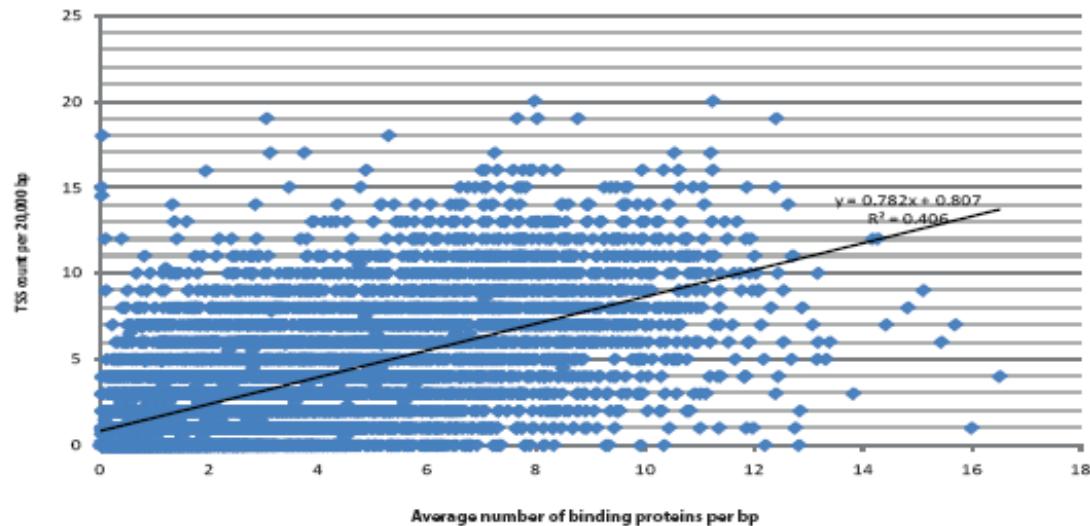


Figure S11 - Density correlation amongst broadly grouped modENCODE datasets

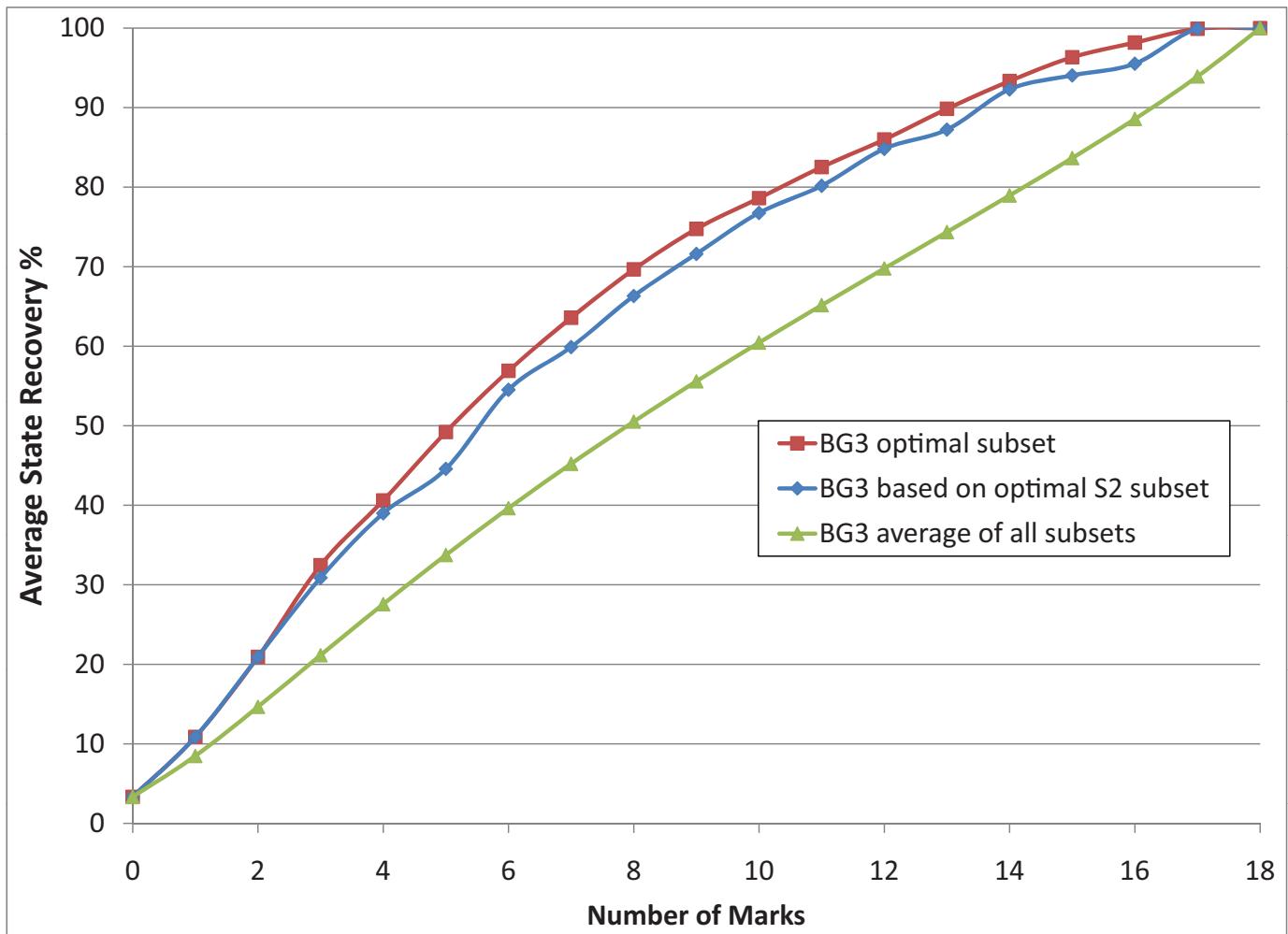


Figure S12 - Chromatin State Recovery with Increasing Subsets of Marks

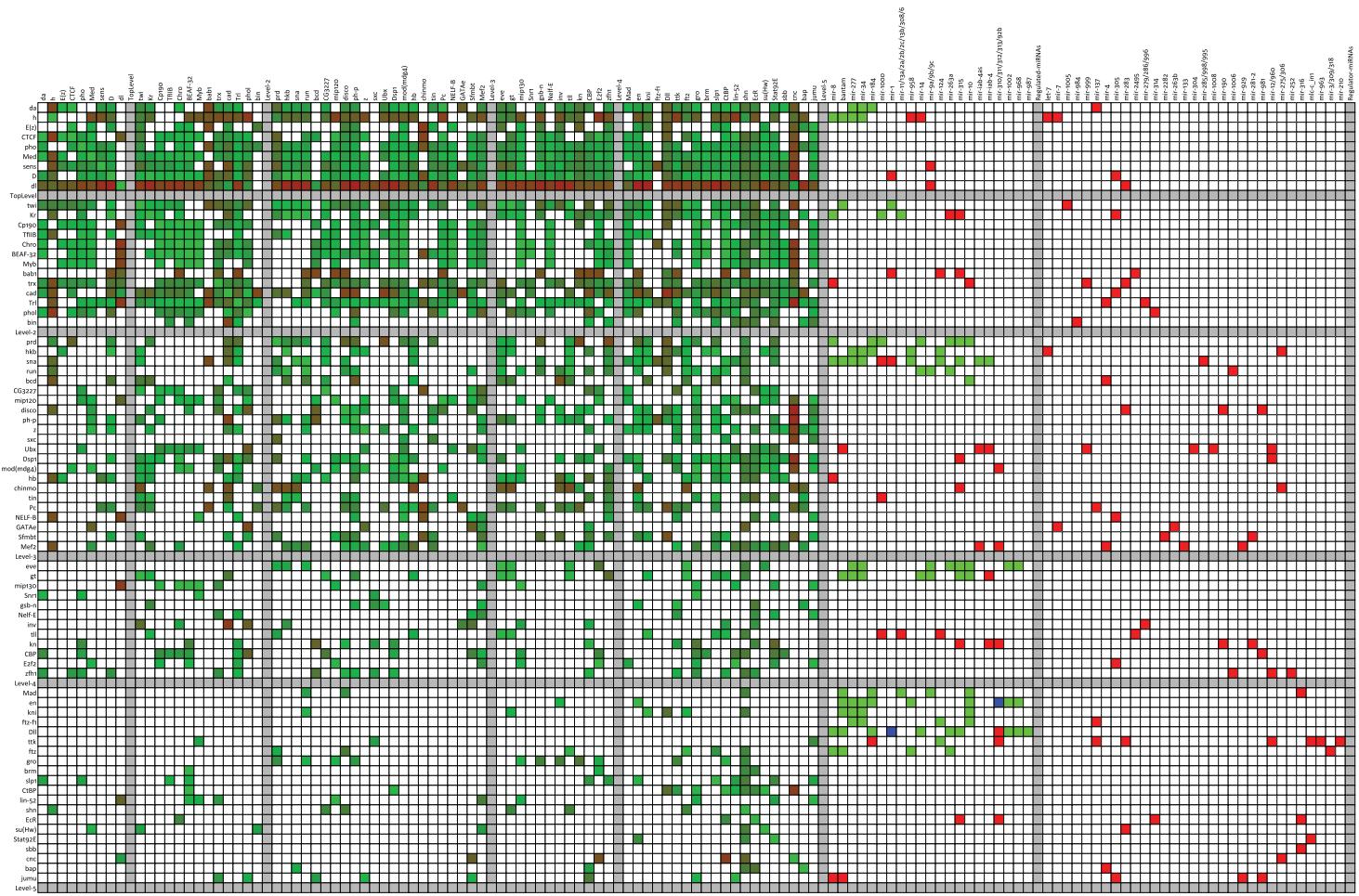


Figure S13 - Heatmap of TF-TF and miR-TF interactions

## A. Definitions and abbreviations for the four regulatory networks compared

Short	Network name	Network description	
mRN	motifNet	Conserved regulatory motifs in promoter regions	"Physical" regulatory network (motif-predicted, and ChIP-inferred)
bRN	boundNet	ChIP-inferred experimental TF binding sites	
kRN	REDfly	Literature-based known regulatory network	"Functional" regulatory networks (known/curated, and integrative/predicted)
fRN	FuncNet	Functional network with chrom/expr activity	

## B. Number of TFs, Targets, and Edges for each network

Short	Network size		
	NumTFs	NTargets	Edges/TG
mRN	104	11,090	7
bRN	76	12,482	13
kRN	82	88	3
fRN	576	9,436	24

## C. Pairwise comparison of networks for common TFs, common targets, and common edges

Short	mRN	Number of targets in common			Enrichment in common edges		
		bRN	kRN	fRN	bRN	kRN	fRN
mRN		10,135	83	10,135	1.56	2.43	4.78
bRN	31		84	9,251	1.90		4.35
kRN	46	35		79		2.79	
fRN	101	57	70				
	Number of TFs in common						

## D. Enrichments for each network for common function, developmental expression, tissues, and protein-protein interactions

Short	True positive rate for predictions			
	Function	DevExpr	Tissues	PPI
mRN	-11%	-4%	-18%	35%
bRN	7%	12%	16%	16%
kRN	37%	64%	43%	61%
fRN	19%	46%	19%	60%
	mult. TFs	mult. TFs	mult. TFs	single TF

Figure S14 - Comparison between motif, bound, known, and functional network

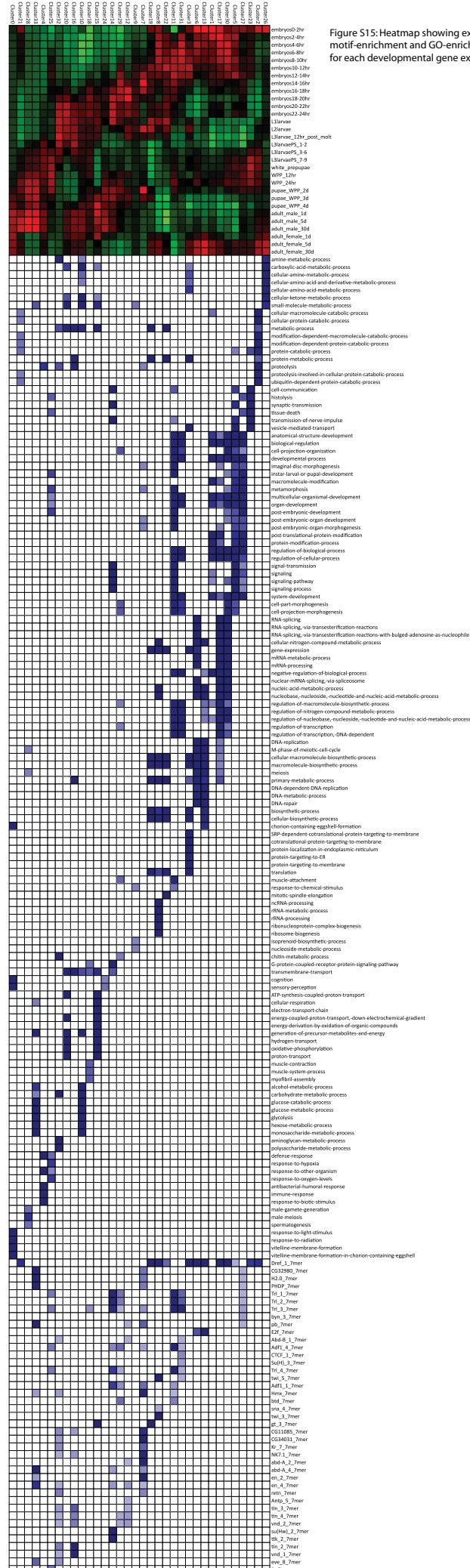


Figure S15: Heatmap showing expression, motif-enrichment and GO-enrichment for each developmental gene expression cluster

## Precision-Recall Curves

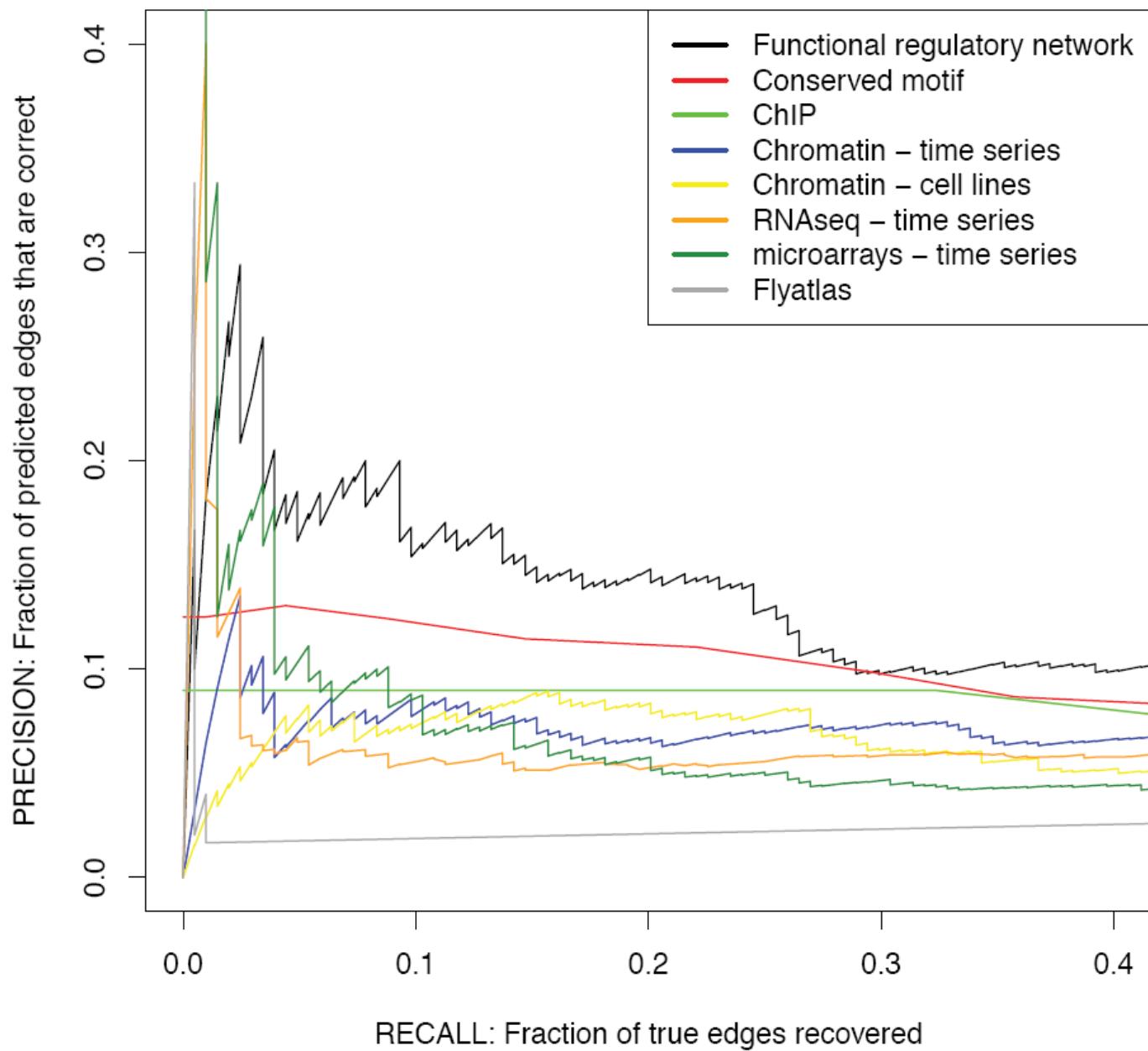


Figure S16 - Precision Recall curves of feature-specific and functional regulatory network

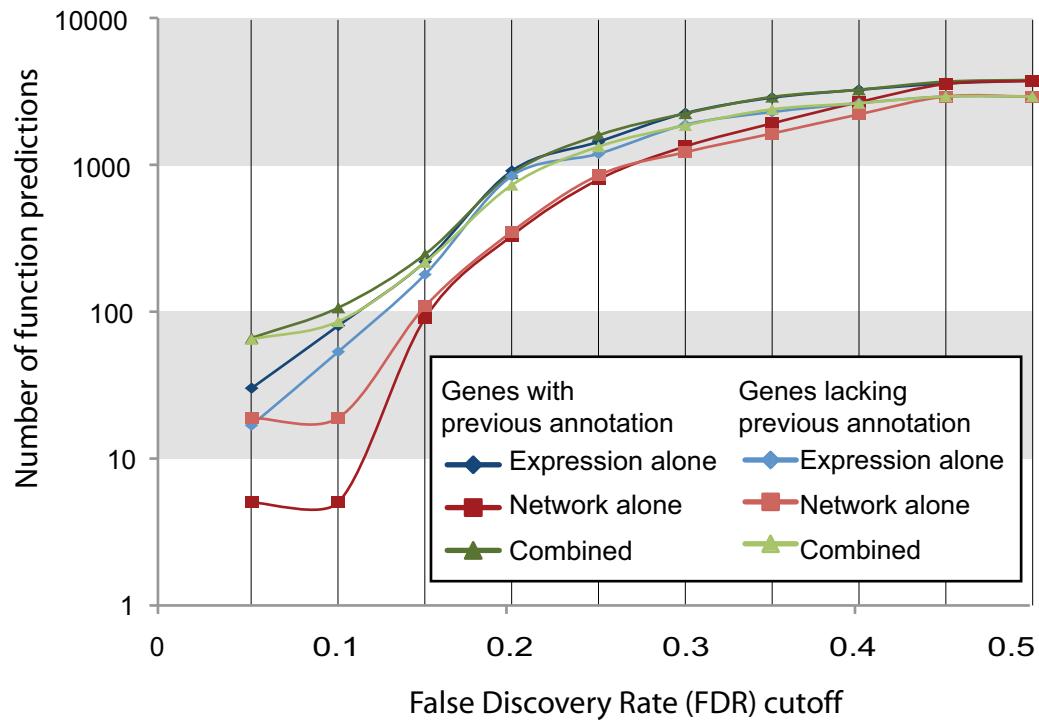


Figure S17 - Number of genes and process terms at different FDRs.

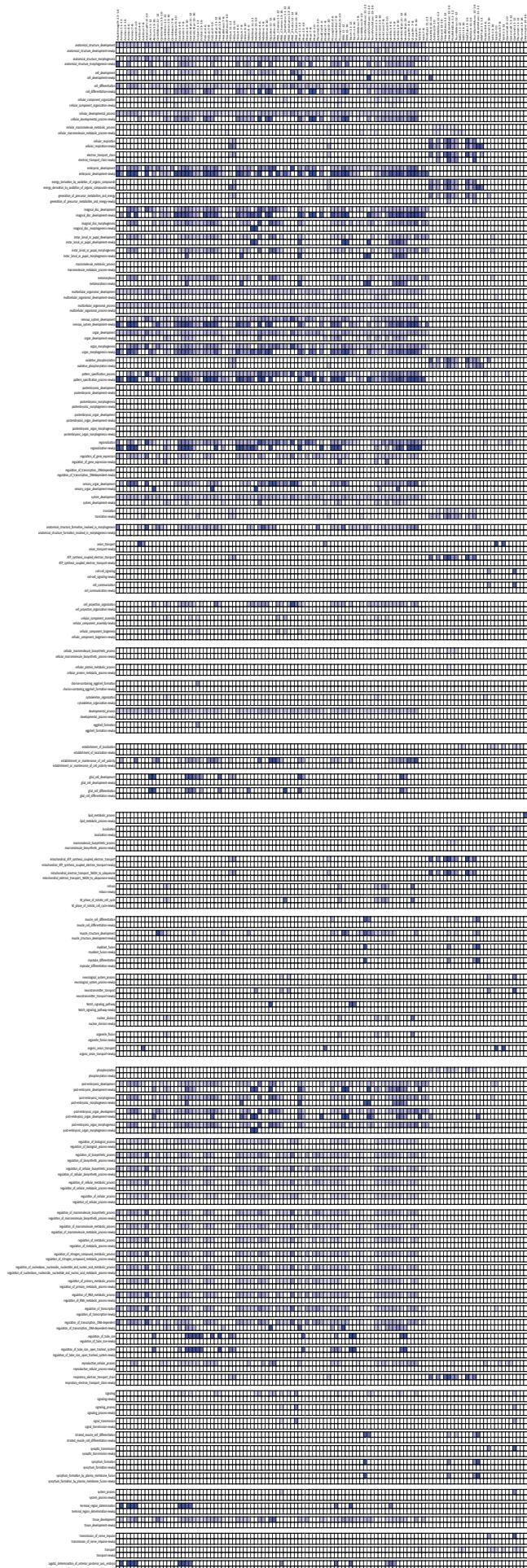


Figure S18 - Fold-enrichment of tissues among novel and known members of a GO process category

Split Number	Sample	Path	Physical Network			Functional Network			Top GO Biological Process on path conditional on genes into split	old p-value	corrected p-value	
			giltmean	child mean	diffmean	burn Genes	num TF Sig [p<0.01]	num TF Not Sig [p>0.01]	Avg. Abs. Express Change Sig TF	Wt. Abs. Express Change Non Sig TF		
1	embryos0-2hr	H	0.0	0.8	0.8	1361	82	34	1.54	1.01	112	464
	M	0.0	0.0	0.0	3353	1	75	66	1.31	3	573	0.99
	L	-1.3	-0.9	-0.9	1395	65	11	1.36	0.97	293	283	0.83
2	embryo4-8hr	H	-1.3	-1.0	-0.4	757	25	51	0.32	0.47	207	369
	L	-2.1	-2.1	-0.8	638	0	76	0.00	0.42	20	556	0.45
3	embryos10-12hr	H	0.2	0.5	0.3	1559	4	72	0.28	0.37	34	542
	L	0.2	0.1	-0.1	1794	2	74	0.29	0.37	5	571	0.10
4	embryos10-12hr	H	-2.0	-1.3	-0.7	400	6	70	0.33	0.37	7	569
	L	-2.0	-3.4	-1.4	238	1	75	0.50	0.36	172	404	0.33
5	embryos14-16hr	H	3.3	4.5	1.2	385	0	76	0.00	0.39	38	538
	L	3.3	2.7	-0.6	996	25	51	0.41	0.38	30	546	0.39
6	embryos14-16hr	H	1.3	2.2	0.9	908	0	76	0.00	0.39	0	576
	L	1.3	2.0	0.8	651	21	55	0.34	0.41	5	571	0.38
7	embryos16-18hr	H	2.7	2.8	0.1	671	13	63	0.39	0.31	0	576
	L	2.7	1.4	-1.3	325	1	75	0.09	0.32	33	543	0.26
8	1LarvaeP5_1-2	H	1.0	1.9	0.9	1482	1	75	0.41	0.52	4	572
	L	1.0	0.1	-0.9	312	0	76	0.00	0.52	0	576	0.00
9	1LarvaeP5_3-6	H	1.9	2.6	0.7	682	0	76	0.00	0.25	0	576
	L	1.9	1.9	0.0	804	0	76	0.00	0.25	0	576	0.00
10	1LarvaeP5_3-6	H	2.6	2.9	0.3	276	0	76	0.00	0.16	0	576
	L	2.6	2.6	0.0	455	0	76	0.00	0.16	0	576	0.00
11	white_prepupae	H	3.0	4.1	2.3	402	1	75	0.07	0.35	2	574
	L	3.0	0.4	-1.6	562	0	76	0.00	0.34	0	576	0.00
12	WPP_12hr	H	2.5	3.3	0.8	433	0	76	0.00	0.35	0	576
	L	2.5	2.6	0.1	367	0	76	0.00	0.35	0	576	0.00
13	WPP_24hr	H	1.7	1.7	0.0	378	0	76	0.00	0.51	1	576
	L	-1.7	-2.5	-0.8	481	0	76	0.00	0.51	154	422	0.57
14	pupae_WPP_2d	H	4.0	6.1	2.3	118	0	76	0.00	0.53	0	576
	L	4.0	2.3	-1.7	287	0	76	0.00	0.53	0	576	0.00
15	pupae_WPP_2d	H	2.3	2.7	0.3	349	1	75	0.40	0.31	2	574
	L	2.3	1.5	-0.9	322	14	62	0.40	0.29	15	561	0.26
16	pupae_WPP_3d	H	0.2	1.7	1.5	186	2	74	0.48	0.31	1	575
	L	0.2	0.1	-0.1	317	0	76	0.00	0.31	1	575	0.07
17	pupae_WPP_4d	H	4.0	5.2	1.2	210	0	76	0.00	0.28	4	572
	L	4.0	1.3	-2.7	192	0	76	0.00	0.28	1	575	0.50
18	pupae_WPP_4d	H	2.5	3.5	1.0	275	0	76	0.00	0.28	1	575
	L	2.5	2.0	-0.5	376	0	76	0.00	0.28	0	576	0.00
19	pupae_WPP_4d	H	1.8	2.7	0.9	121	0	76	0.00	0.28	0	576
	L	1.8	0.4	-1.5	166	0	76	0.00	0.28	0	576	0.00

Figure S19 - Information on developmental time course expression splits identified by DREM

Figure S20 - Squared error, correlation coefficients and Regulators of 1487 predictable genes

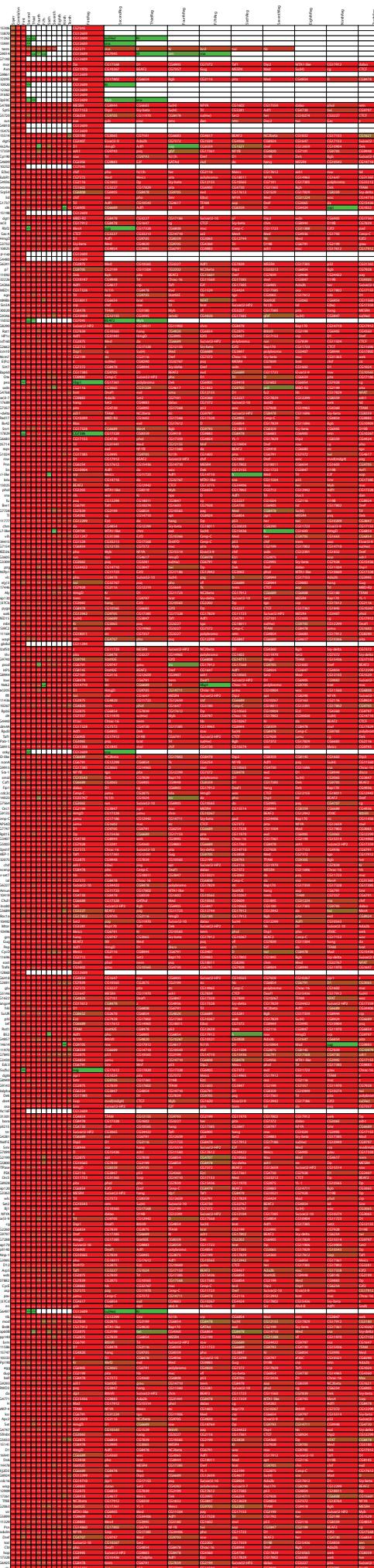


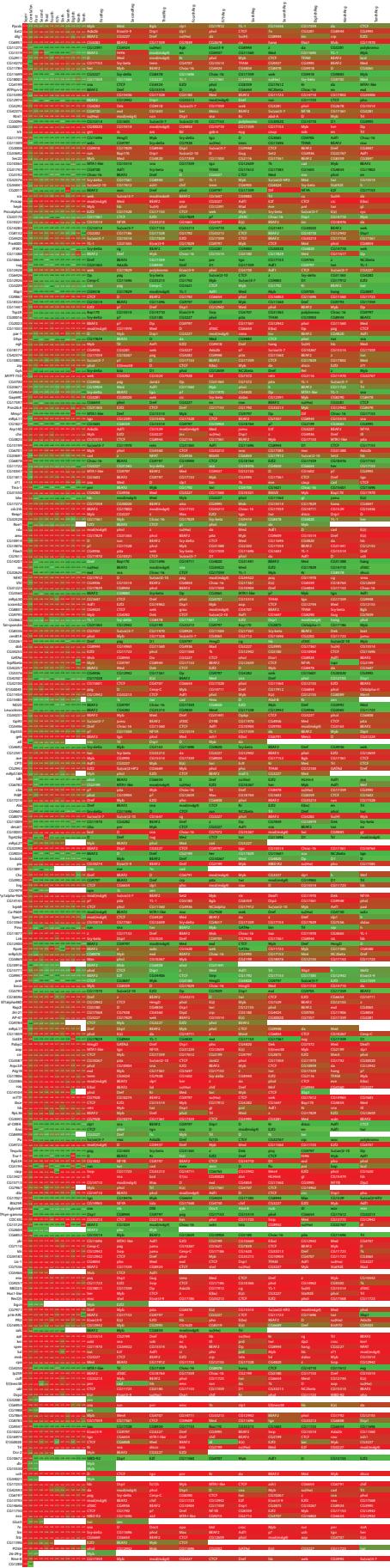








Figure S21 - Squared error, correlation coefficients and Regulators of 324 unpredictable genes



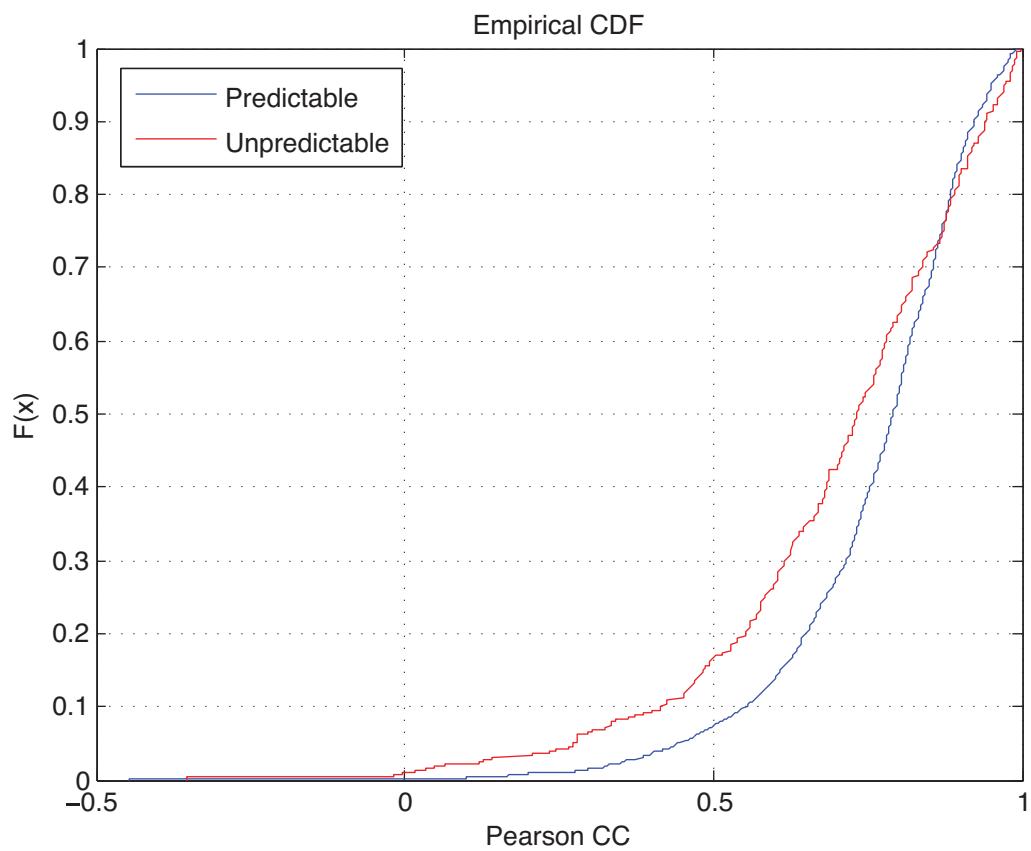


Figure S22 - Cumulative distribution of correlations of individual gene time courses obtained from the two developmental RNAseq datasets.

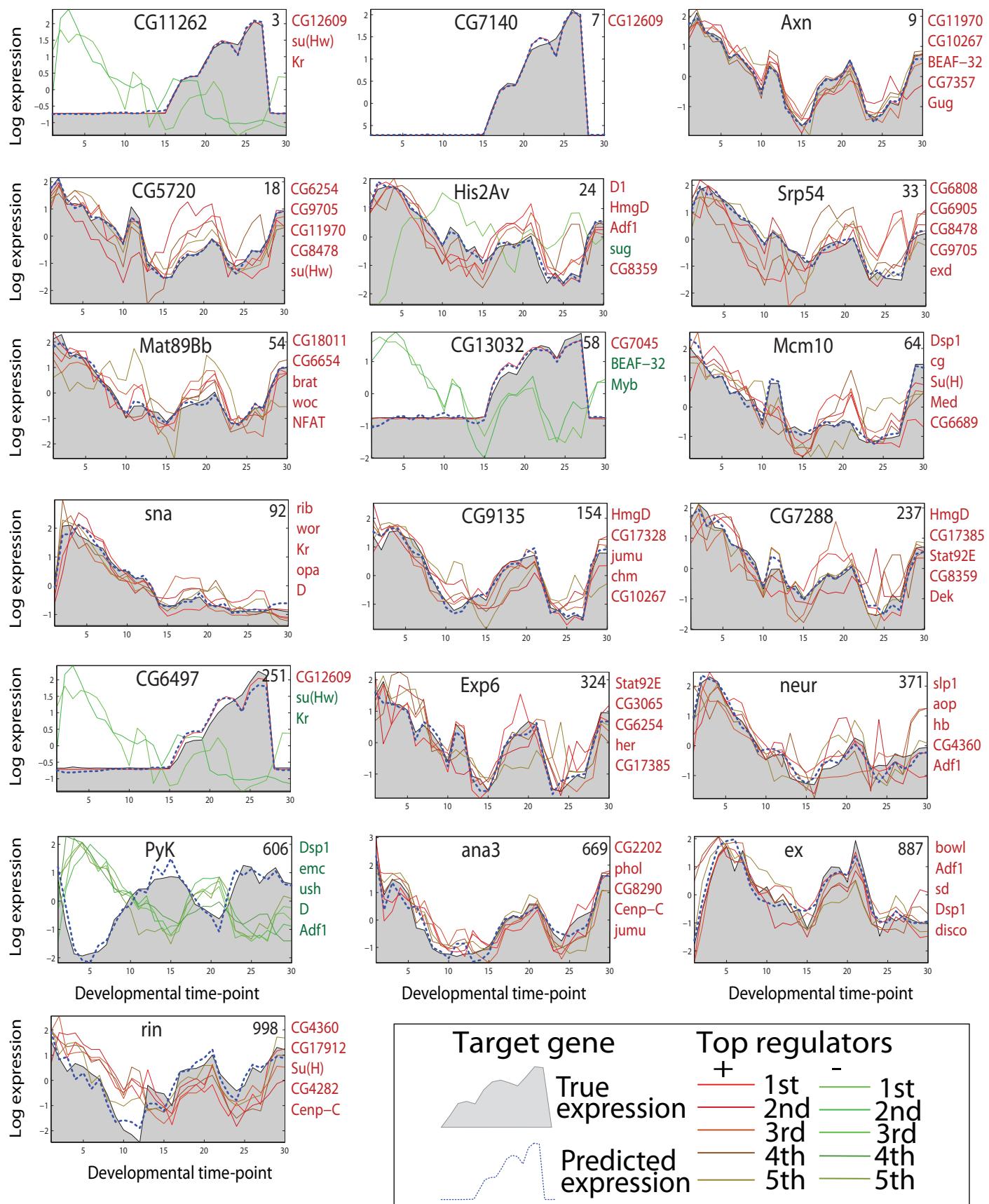


Figure S23 - Expanded set of gene expression prediction examples

Category	Annotation types	Assay	# Experiments	Cell Lines			Developmental Stage <sup>2</sup>			Factors <sup>3,7</sup>			modENCODE ids	GEO/SRA ids		
				BG3	Cl.8+	Kc167	S2-DRSC	Others	Embryo	Larva	Pupa	Adult/Male	Adult/Female			
Transcription/Splicing																
Gene Expression	gene, transcripts, cell-type & tissue-specific expression	RNA tiling array (polyA, total RNA)	80	+ + + +	22	27	6	6	3	3	3			GSM409470, GSM409471, GSM409472, GSM409473, GSM409474, GSM409475, GSM409476, GSM409477, GSM409478, GSM409479, GSM409480, GSM409481, GSM409482, GSM409483, GSM409484, GSM409485, GSM409486, GSM409487, GSM409488, GSM409489, GSM409490, GSM409491, GSM409492, GSM409493, GSM409494, GSM409495, GSM409496, GSM409497, GSM409498, GSM409499, GSM409500, GSM409501, GSM409502, GSM409503, GSM409504, GSM409505, GSM409506, GSM409507, GSM409508, GSM409509, GSM409510, GSM409511, GSM409512, GSM409513, GSM409514, GSM409515, GSM409516, GSM409517, GSM409518, GSM409519, GSM409520, GSM409521, GSM409522, GSM409523, GSM409524, GSM409525, GSM409526, GSM409527, GSM409528, GSM409529, GSM409530, GSM409531, GSM409532, GSM409533, GSM409534, GSM409535, GSM409536, GSM409537, GSM409538, GSM409539, GSM409540, GSM409541, GSM409542, GSM409543, GSM409544, GSM409545, GSM409546, GSM409547, GSM409548, GSM409549, GSM409550, GSM409551, GSM409552, GSM409553, GSM409554, GSM409555, GSM409556, GSM409557, GSM409558, GSM409559, GSM409560, GSM409561, GSM409562, GSM409563, GSM409564, GSM409565, GSM409566, GSM409567, GSM409568, GSM409569, GSM409570, GSM409571, GSM409572, GSM409573, GSM409574, GSM409575, GSM409576, GSM409577, GSM409578, GSM409579, GSM409580, GSM409581, GSM409582, GSM409583, GSM409584, GSM409585, GSM409586, GSM409587, GSM409588, GSM409589, GSM409590, GSM409591, GSM409592, GSM409593, GSM409594, GSM409595, GSM409596, GSM409597, GSM409598, GSM409599, GSM409600, GSM409601, GSM409602, GSM409603, GSM409604, GSM409605, GSM409606, GSM409607, GSM409608, GSM409609, GSM409610, GSM409611, GSM409612, GSM409613, GSM409614, GSM409615, GSM409616, GSM409617, GSM409618, GSM571093, GSM571094, GSM571095, GSM571096, GSM571097, GSM571098, GSM571099, GSM571100, GSM571101, GSM571102, GSM571103, GSM571104, GSM571105, GSM571106, GSM571107, GSM409643, GSM409644, GSM409645, GSM409646, GSM409647, GSM409648, GSM571108, GSM571109, GSM571110, GSM571111, GSM571112, GSM571113, GSM571114, GSM571115, GSM571116, GSM571117, GSM571118, GSM571119, GSM571120, GSM571121, GSM571122, GSM571123, GSM571124, GSM571125, GSM571126, GSM571127, GSM571128, GSM571129, GSM571130, GSM571131, GSM571132, GSM571133, GSM571134, GSM571135, GSM571136, GSM571137, GSM571138, GSM571139, GSM624330, GSM624331, GSM624332, GSM624333, GSM624334, GSM624335, GSM624336, GSM624337, GSM624338, GSM624339, GSM624340, GSM624341, GSM624342, GSM624343, GSM624344, GSM624345, GSM624346, GSM624347, GSM624348, GSM624349, GSM624350, GSM624351, GSM624352, GSM624353, GSM624354, GSM624355, GSM624356, GSM624357, GSM624358, GSM624359		
coding and non-coding genes, transcripts, splice junctions	Total RNA-seq	16	+ + + +		12									2676, 2677, 2678, 2679, 2680, 2681, 2682, 2683, 2684, 2685, 2686, 2687, 360, 362, 364, 366, SRX015641, SRX015642, SRX015643, SRX015644, SRX015645, SRX015646, SRX015647, SRX015648, SRX015649, SRX015650, SRX015651, SRX015652, GSM451044, GSM451127, GSM451031, GSM450972		
	RNA-seq (polyA+)	38	+ + + +		13	6	6	3	3					342, 343, 633, 650, 651, 652, 712, 2866, 2867, 2868, 2869, 2870, 2871, 2872, 2873, 2874, 2875, 2876, 2877, 2878, 2879, 2880, 2881, 2882, 2883, 2884, 2885, 2886, 2887, 2888, 2889, 2890, 2891, 2892, 2893, 2894, 2895, SRR023595, SRR023662, SRR023680, SRR023702, SRR023737, SRR023828, SRR024220, SRR035393, SRR023548, SRR023598, SRR023674, SRR023681, SRR023712, SRR023757, SRR029232, SRR029236, SRR023547, SRR023607, SRR023645, SRR023651, SRR023717, SRR023730, SRR029230, SRR029234, SRR023543, SRR023648, SRR023664, SRR023678, SRR023741, SRR023825, SRR035394, SRR023601, SRR023650, SRR023698, SRR023714, SRR023752, SRR023823, SRR024218, SRR035395, SRR023550, SRR023605, SRR023606, SRR023642, SRR023658, SRR023672, SRR023679, SRR023713, SRR029176, SRR029231, SRR029233, SRR029235, SRR023597, SRR023646, SRR023661, SRR023666, SRR023706, SRR023835, SRR035410, SRR023542, SRR023670, SRR023719, SRR023761, SRR023824, SRR024016, SRR035223, SRR035411, SRR035412, SRR023507, SRR023649, SRR023677, SRR023731, SRR023760, SRR027111, SRR023546, SRR023608, SRR023638, SRR023640, SRR023697, SRR023726, SRR023758, SRR035413, SRR023505, SRR023676, SRR023683, SRR023690, SRR023692, SRR023742, SRR027108, SRR023729, SRR023739, SRR023782, SRR023865, SRR026433, SRR035219, SRR035392, SRR023687, SRR023695, SRR023704, SRR023784, SRR023837, SRR026432, SRR035217, SRR023544, SRR023639, SRR023647, SRR023689, SRR023716, SRR023833, SRR035414, SRR023541, SRR023656, SRR023708, SRR023756, SRR023830, SRR035415, SRR023667, SRR023721, SRR023743, SRR023785, SRR023829, SRR026431, SRR023738, SRR023748, SRR023783, SRR024216, SRR026430, SRR035218, SRR035396, SRR023681, SRR023688, SRR023691, SRR023732, SRR024009, SRR024217, SRR027113, SRR035407, SRR035408, SRR023549, SRR023596, SRR023603, SRR023657, SRR023701, SRR023749, SRR023750, SRR023754, SRR023759, SRR024219, SRR035409, SRR023604, SRR023644, SRR023694, SRR023727, SRR023831, SRR024010, SRR024221, SRR035396, SRR023746, SRR023836, SRR024014, SRR035420, SRR035405, SRR035406, SRR023504, SRR023654, SRR023668, SRR023691, SRR023732, SRR024009, SRR024217, SRR027113, SRR035407, SRR035408, SRR023549, SRR023596, SRR023603, SRR023657, SRR023701, SRR023749, SRR023750, SRR023754, SRR023759, SRR024219, SRR035409, SRR023604, SRR023644, SRR023694, SRR023727, SRR023831, SRR024010, SRR024221, SRR035396, SRR023746, SRR023602, SRR023643, SRR023703, SRR023724, SRR023832, SRR024012, SRR035221, SRR035398, SRR035399, SRR023540, SRR023641, SRR023725, SRR023734, SRR023736, SRR023822, SRR024011, SRR035400, SRR035401, SRR023600, SRR023707, SRR023715, SRR023720, SRR023751, SRR023826, SRR035402, SRR023503, SRR023599, SRR023673, SRR023682, SRR023686, SRR023699, SRR023700, SRR024215, SRR027114, SRR035422, SRR035403, SRR023506, SRR023665, SRR023684, SRR023718, SRR023728, SRR027109, SRR023548, SRR023685, SRR023711, SRR023733, SRR023834, SRR024015, SRR035404, GSM451119, GSM439966, GSM390059, GSM390063, GSM390065, GSM390061		
	cDNA sequencing	3204 <sup>4</sup>												3275		
	splicing	RNAi/RNA-seq	27			1								Splicing	B52	1013











**Table S2 - Drosophila Cell Lines Characterized by modENCODE**

Short Name	Official Name	Genotype	Sex	Source Tissue	Source Stage
<b>BG3</b>	ML-DmBG3-c2	y v f mal	-	CNS	third instar larval
<b>Cl.8+</b>	CME W1 Cl.8+	-	M	dorsal mesothoracic disc	third instar larval
<b>Kc167</b>	Kc167	se/e	*	embryo	late embryonic
<b>S2-DRSC</b>	S2-DRSC	-	M	embryo	late embryonic
<b>1182-4H</b>	1182-4H	mh	-	embryo	late embryonic
<b>CME W2</b>	CME W2	-	-	dorsal mesothoracic disc	third instar larval
<b>CME-L1</b>	CME L1	-	-	ventral prothoracic disc	third instar larval
<b>DmBG1-c1</b>	ML-DmBG1-c1	y v f mal	-	CNS	third instar larval
<b>DmBG2-c2</b>	ML-DmBG1-c2	y v f mal	-	CNS	third instar larval
<b>DmD11</b>	ML-DmD11	y v f mal	-	eye-antenna disc	third instar larval
<b>DmD16</b>	ML-DmD16-c3	y v f mal	*	dorsal mesothoracic disc	third instar larval
<b>DmD17</b>	ML-DmD17-c3	y v f mal	-	dorsal mesothoracic disc	third instar larval
<b>DmD20-c2</b>	ML-DmD20-c2	y v f mal	M	antenna disc	third instar larval
<b>DmD20-c5</b>	ML-DmD20-c5	y v f mal	M	antenna disc	third instar larval
<b>DmD21</b>	ML-DmD21	y v f mal	M	dorsal mesothoracic disc	third instar larval
<b>DMD32</b>	ML-DmD32	y v f mal	-	dorsal mesothoracic disc	third instar larval
<b>DmD4-c1</b>	ML-DmD4-c1	y v f mal	-	imaginal disc	third instar larval
<b>DmD8</b>	ML-DmD8	y v f mal	-	dorsal mesothoracic disc	third instar larval
<b>DmD9</b>	ML-DmD9	y v f mal	-	dorsal mesothoracic disc	third instar larval
<b>GM2</b>	GM2	-	-	embryo	late embryonic
<b>Kc-Rubin</b>	Kc-GMR	se/e	*	embryo	late embryonic
		I(2) malignant			
<b>mbn2</b>	mbn2	blood	-	embryo/larval hemocyte	late embryonic
		neoplasm			
<b>ML-DmBG1c1</b>	ML-DmBG1-c1	y v f mal	-	CNS	third instar larval
<b>ML-DmD20c5</b>	ML-DmD20-c5	y v f mal	M	antenna disc	third instar larval
<b>ML-DmD21</b>	ML-DmD21	y v f mal	M	dorsal mesothoracic disc	third instar larval
<b>OSS</b>	Ovary Somatic Sheet	-	*	follicle cell, ovary	adult
<b>S1</b>	S1	Oregon R	-	embryo	late embryonic
<b>S2</b>	S2	Oregon R	-	embryo	late embryonic
<b>S2-NP</b>	S2-NP	Oregon R	M	embryo	late embryonic
<b>S2-Rubin</b>	S2-GMR	Oregon R	M	embryo	late embryonic
<b>S2R+</b>	S2R+	Oregon R	M	embryo	late embryonic
<b>S3</b>	S3	Oregon R	-	embryo	late embryonic
<b>Sg4</b>	Sg4	Oregon R	M	embryo	late embryonic

Table S3 Protein symbols used in text

Gene Symbol	Name	FlyBase ID	Annotation ID
<b>AGO1</b>	Argonaute-1	FBgn0026611	CG6671
<b>AGO2</b>	Argonaute 2	FBgn0087035	CG7439
<b>B52</b>	B52	FBgn0004587	CG10851
<b>BAB1</b>	bric a brac 1	FBgn0004870	CG9097
<b>BEAF-32</b>	Boundary element-associated factor of 32kD	FBgn0015602	CG10159
<b>BL</b>	bancal	FBgn0015907	CG13425
<b>BRE1</b>	Bre1	FBgn0086694	CG10542
<b>BRM</b>	brahma	FBgn0000212	CG5942
<b>CAD</b>	caudal	FBgn0000251	CG1759
<b>CG17838</b>	CG17838	FBgn0038826	CG17838
<b>CG30122</b>	CG30122	FBgn0050122	CG30122
<b>CHINMO</b>	Chronologically inappropriate morphogenesis	FBgn0086758	CG31666
<b>CNC</b>	cap-n-collar	FBgn0000338	CG17894
<b>CP190</b>	Centrosomal protein 190kD	FBgn0000283	CG6384
<b>CTCF</b>	CTCF	FBgn0035769	CG8591
<b>DISCO</b>	disconnected	FBgn0000459	CG9908
<b>E(BX)</b>	Enhancer of bithorax	FBgn0000541	CG32346
<b>E(Z)</b>	Enhancer of zeste	FBgn0000629	CG6502
<b>ECR</b>	Ecdysone receptor	FBgn0000546	CG1765
<b>EN</b>	engrailed	FBgn0000577	CG9015
<b>EVE</b>	even skipped	FBgn0000606	CG2328
<b>FTZ-F1</b>	ftz transcription factor 1	FBgn0001078	CG4059
<b>GATAE</b>	GATAe	FBgn0038391	CG10278
<b>GLO</b>	glorund	FBgn0259139	CG6946
<b>GRO</b>	groucho	FBgn0001139	CG8384
<b>GSB-N</b>	gooseberry-neuro	FBgn0001147	CG2692
<b>H</b>	hairy	FBgn0001168	CG6494
<b>HDAC3</b>	Histone deacetylase 3	FBgn0025825	CG2128
<b>HDAC4</b>	HDAC4	FBgn0041210	CG1770
<b>HDAC6</b>	HDAC6	FBgn0026428	CG6170
<b>HDACX</b>	Histone deacetylase X	FBgn0051119	CG31119
<b>HEPH</b>	hephaestus	FBgn0011224	CG31000
<b>HKB</b>	huckebein	FBgn0261434	CG9768
<b>HP1C</b>	HP1c	FBgn0039019	CG6990
<b>HRB87F</b>	Heterogeneous nuclear ribonucleoprotein at 87F	FBgn0004237	CG12749
<b>HRB98DE</b>	Heterogeneous nuclear ribonucleoprotein at 98DE	FBgn0001215	CG9983
<b>INV</b>	invected	FBgn0001269	CG17835
<b>JIL-1</b>	JIL-1	FBgn0020412	CG6297
<b>MBD-R2</b>	MBD-R2	FBgn0038016	CG10042
<b>Mi-2</b>	Mi-2	FBgn0013591	CG8103
<b>MOD(MDG4)</b>	modifier of mdg4	FBgn0002781	CG32491
<b>MSI</b>	musashi	FBgn0011666	CG5099
<b>MUB</b>	mushroom-body expressed	FBgn0014362	CG7437
<b>NEJ</b>	Nejire, homologue of CREB-binding protein (CBP)	FBgn0261617	CG15319
<b>ORC1</b>	Origin recognition complex subunit 1	FBgn0022772	CG10667
<b>ORC2</b>	Origin recognition complex subunit 2	FBgn0015270	CG3041
<b>ORC4</b>	Origin recognition complex subunit 4	FBgn0023181	CG2917
<b>ORC5</b>	Origin recognition complex subunit 5	FBgn0015271	CG7833

### Protein symbols used in text

<b>ORC6</b>	Origin recognition complex subunit 6	FBgn0023180	CG1584
<b>PC</b>	Polycomb	FBgn0003042	CG32443
<b>PCL</b>	Polycomblike	FBgn0003044	CG5109
<b>PIWI</b>	piwi	FBgn0004872	CG6122
<b>QKR58E-2</b>	quaking related 58E-2	FBgn0022985	CG5821
<b>RBP1</b>	RNA-binding protein 1	FBgn0260944	CG17136
<b>RBP1-LIKE</b>	Rbp1-like	FBgn0030479	CG1987
<b>RPD3</b>	Rpd3	FBgn0015805	CG7471
<b>RSF1</b>	Repressor splicing factor 1	FBgn0011305	CG5655
<b>RUMP</b>	rumpelstiltskin	FBgn0260010	CG9373
<b>RUN</b>	runt	FBgn0003300	CG1849
<b>SBB</b>	scribbler	FBgn0010575	CG5580
<b>SC35</b>	SC35	FBgn0040286	CG5442
<b>SCE</b>	Sex combs extra	FBgn0003330	CG5595
<b>SENS</b>	senseless	FBgn0002573	CG32120
<b>SF1</b>	Splicing factor 1	FBgn0025571	CG5836
<b>SNR1</b>	Snf5-related 1	FBgn0011715	CG1064
<b>SQD</b>	squid	FBgn0086897	CG16901
<b>SRP54</b>	Srp54	FBgn0024285	CG4602
<b>STAT92E</b>	Signal-transducer and activator of transcription protein at 92E	FBgn0016917	CG4257
<b>SU(VAR)205</b>	Suppressor of variegation 205	FBgn0003607	CG8409
<b>SU(VAR)3-9</b>	Suppressor of variegation 3-9	FBgn0003600	CG6476
<b>TRL</b>	Trithorax-like	FBgn0013263	CG33261
<b>TRX</b>	trithorax	FBgn0003862	CG8651
<b>TTK</b>	tramtrack	FBgn0003870	CG1856
<b>UBX</b>	Ultrabithorax	FBgn0003944	CG10388
<b>UPF1</b>	Upf1	FBgn0030354	CG1559
<b>WDS</b>	will die slowly	FBgn0040066	CG17437
<b>XL6</b>	xl6	FBgn0028554	CG10203
<b>YTR</b>	yantar	FBgn0021895	CG18426
<b>YU</b>	yu	FBgn0029764	CG3249
<b>ZFH1</b>	Zn finger homeodomain 1	FBgn0004606	CG1322

Table S3 - Protein Symbols used in the text

**Table S4. cis-NAT-siRNA loci in *D. melanogaster***

Loci were detected at a cutoff of 50 21nt reads (gray column) and >2/3 representation of 21mers within a sense/antisense transcript overlap.

CisNAT	gene1	gene2	position	Overlap	overlap size (bp)	total read number in overlap	total 21mers in overlap	totalN_plus	totalN_minus	totalN_plus_sum	totalN_minus_sum	totalN_plus_21	totalN_minus_21
CG7739-RA	AGO2-RA	3L:115533779...11554211	catt3	CG7739-RA	AGO2-RA	3L:115533779...11554211	432	32940	27374	19647	13293	16547	10827
MED12-RA	AGO2-RA	3L:116000000...116000000	catt3	CG11884-RA	AGO2-RA	3L:116000000...116000000	426	10643	5657	3000	4648	4168	4168
cher-RR	call-RA	3R:12917257...12917580	catt3	CG11884-RA	f9K1-RA	3R:12917257...12917580	1018	7373	5583	3371	4002	2635	2340
catatin-60-RA	CG1115-RA	3R:1053046...1053248	catt3	catatin-60-RA	CG1115-RA	3R:1053046...1053248	202	3504	2760	1321	2183	1785	1132
Tango-70-RA	O-fut1-RA	2R:10050358...1005054	catt3	Tango-70-RA	O-fut1-RA	2R:10050358...1005054	196	3296	2698	1423	1872	1198	1500
CG6218-RA	Sra-1-RA	3R:11163628...11164103	catt3	CG6218-RA	Sra-1-RA	3R:11163628...11164103	475	3111	2404	1165	1945	883	1521
alpha4G11-RA	CG1271-RA	3R:102130991...302130996	catt3	alpha4G11-RA	CG1271-RA	3R:102130991...302130996	197	2399	1944	1003	1368	1032	939
CG11884-RA	AGO2-RA	3L:12917257...12917580	catt3	CG11884-RA	AGO2-RA	3L:12917257...12917580	217	2538	2741	1323	1432	1432	1432
CG5728-RA	BWD3-RA	2R:102154151...20154344	catt3	CG5728-RA	BWD3-RA	2R:102154151...20154344	193	2268	1758	1110	1158	858	900
CG10214-RA	Lsd-1-RC	3R:1959241...19592453	catt3	CG10214-RA	Lsd-1-RC	3R:1959241...19592453	212	2132	1698	819	1313	635	1063
CG33969-RA	fbl-RA	3L:20383986...20385014	catt3	CG33969-RA	fbl-RA	3L:20383986...20385014	1028	2307	1662	1112	1195	741	921
CG1365-RA	cerB1A-RA	3R:18949676...18949941	catt3	CG1365-RA	cerB1A-RA	3R:18949676...18949941	265	2106	1597	1039	1067	751	846
CG9062-RA	Vhl-RA	2R:7167919...7168538	catt3	CG9062-RA	Vhl-RA	2R:7167919...7168538	619	2103	1551	1203	901	856	695
AdoP-RA	DNase-1-RA	3R:1663500...1663687	catt3	AdoP-RA	DNase-1-RA	3R:1663500...1663687	315	2044	1546	938	1101	734	813
Rip4-RA	WIF1-RA	X:1635050...1635050	catt3	Rip4-RA	WIF1-RA	X:1635050...1635050	101	1897	1435	869	864	535	535
CG9779-RA	MPI-RA	3R:135319...135559	catt3	CG9779-RA	MPI-RA	3R:135319...135559	317	1715	1419	849	866	696	733
Magr-Mago-RA	2R:17016965...17017272	catt3	Magr-Mago-RA	2R:17016965...17017272	catt3	307	1833	1402	791	1042	593	809	
CG1553-RA	mus202-RA	2R:3688250...3688549	catt3	CG1553-RA	mus202-RA	2R:3688250...3688549	299	1856	1363	1289	567	1005	358
Ets97D-RA	T48-RA	3R:22736497...22736722	catt3	Ets97D-RA	T48-RA	3R:22736497...22736722	225	1717	1288	774	943	561	727
Mes-4-a	CG5508-RA	3R:2375842...2375884	catt3	Mes-4-a	CG5508-RA	3R:2375842...2375884	418	1789	1276	747	1042	554	721
Rip4-RA	Sap30-RA	X:1663438...1663687	catt3	Rip4-RA	Sap30-RA	X:1663438...1663687	249	1552	1259	655	897	531	690
CG6218-RA	cerB1A-RA	3R:16630768...16631932	catt3	CG6218-RA	cerB1A-RA	3R:16630768...16631932	192	1738	1381	831	860	542	542
PPT1-RA	C17177-RA	2R:6356934...6357150	catt3	PPT1-RA	C17177-RA	2R:6356934...6357150	216	1500	1205	692	808	567	630
THIIP-beta-H-Cap-H2-RA	3R:6592191...6593248	catt3	THIIP-beta-H-Cap-H2-RA	3R:6592191...6593248	catt3	287	1540	1187	798	1065	475	380	
CG11523-RA	Cap-H2-RA	3L:22258906...22259124	catt3	CG11523-RA	Cap-H2-RA	3L:22258906...22259124	218	1565	1165	850	714	661	505
CG16903-RA	CC4025-RA	X:1972324...1972488	catt3	CG16903-RA	CC4025-RA	X:1972324...1972488	164	1428	1141	716	712	542	599
My5a45a-tsu-RA	2R:499906...4999122	catt3	My5a45a-tsu-RA	2R:499906...4999122	catt3	116	1531	1057	463	1068	314	743	
CG33080-RA	CG4096-RA	X:5602545...5602666	catt3	CG33080-RA	CG4096-RA	X:5602545...5602666	121	1349	1055	623	726	533	522
CG4257-RA	cerB1A-RA	3R:16630768...16631932	catt3	CG4257-RA	cerB1A-RA	3R:16630768...16631932	122	1370	1056	597	449	349	349
CG3397-RA	(f2)P-RA	2L:1954376...1955259	catt3	CG3397-RA	(f2)P-RA	2L:1954376...1955259	151	1398	1031	652	546	405	426
CG7816-RA	Ros8-RA	3R:2568935...2568970	catt3	CG7816-RA	Ros8-RA	3R:2568935...2568970	347	1378	1029	523	855	378	651
CG3004-RA	Gga-RA	X:998140...998128	catt3	CG3004-RA	Gga-RA	X:998140...998128	188	1377	1023	801	576	482	440
curl-4	CG1120-RA	TA:3989765...3990023	catt3	curl-4	CG1120-RA	TA:3989765...3990023	141	1326	1008	872	453	679	329
CG1599-RA	CG1665-RA	2R:5728572...5728667	catt3	CG1599-RA	CG1665-RA	2R:5728572...5728667	95	1265	1001	876	706	295	295
wmd3	CG3018-RA	2R:19432260...19432517	catt3	wmd3	CG3018-RA	2R:19432260...19432517	257	1271	984	923	348	721	263
bcr9-RA	Pgd-RA	2R:2042496...2042500	catt3	bcr9-RA	Pgd-RA	2R:2042496...2042500	119	1238	939	301	703	231	231
Alb-Alpha-RA	2R:169810...169810	catt3	Alb-Alpha-RA	2R:169810...169810	catt3	432	1342	919	823	519	560	560	
SyXb-RA	mtd-RA	3R:20449436...20449621	catt3	SyXb-RA	mtd-RA	3R:20449436...20449621	185	1268	889	687	282	436	436
EndoG-RA	Cr5-RA	2R:8064950...8065042	catt3	EndoG-RA	Cr5-RA	2R:8064950...8065042	192	1317	885	574	388	497	497
sec15-RA	CG5745-RA	3R:16930768...16931030	catt3	sec15-RA	CG5745-RA	3R:16930768...16931030	262	1076	845	699	285	560	560
Ppat-Dock-RA	Txi-RA	3L:5760074...5760278	catt3	Ppat-Dock-RA	Txi-RA	3L:5760074...5760278	204	1135	796	355	780	265	531
CG10214-RA	CG7376-RA	3L:6171777...6172041	catt3	CG10214-RA	CG7376-RA	3L:6171777...6172041	264	1013	789	338	675	260	530
vari-RA	Pomp-RA	2L:20788642...20788642	catt3	vari-RA	Pomp-RA	2L:20788642...20788642	458	1059	771	562	497	414	357
CG6502-RA	CG6502-RA	2R:1908416...1908416	catt3	CG6502-RA	CG6502-RA	2R:1908416...1908416	275	927	736	438	329	295	295
CG14149-RA	CG5357-RA	3R:12404766...12404980	catt3	CG14149-RA	CG5357-RA	3R:12404766...12404980	234	911	727	401	374	403	403
CG1514-RA	CG1514-RA	2L:17413248...17413503	catt3	CG1514-RA	CG1514-RA	2L:17413248...17413503	255	1204	911	704	522	389	303
CG6194-RA	CG6194-RA	3R:11179366...11179468	catt3	CG6194-RA	CG6194-RA	3R:11179366...11179468	282	1220	972	702	524	297	405
CG3313-RA	CG3352-RA	3R:7840500...7840824	catt3	CG3313-RA	CG3352-RA	3R:7840500...7840824	233	1200	984	698	485	389	309
mad2-a	CG5537-RA	3L:5805232...5805330	catt3	mad2-a	CG5537-RA	3L:5805232...5805330	98	986	696	277	708	218	477
act-RD	fzy-RA	TL:16312940...16313094	catt3	act-RD	fzy-RA	TL:16312940...16313094	159	857	676	734	123	595	81
CG6194-RA	CG6194-RA	3R:11179366...11179468	catt3	CG6194-RA	CG6194-RA	3R:11179366...11179468	180	945	694	495	541	311	363
MinP-RA	CG6194-RA	2L:20245881...20245881	catt3	MinP-RA	CG6194-RA	2L:20245881...20245881	181	945	694	495	533	307	311
bigmax-RA	CG368-RA	3L:21130244...21130244	catt3	bigmax-RA	CG368-RA	3L:21130244...21130244	230	1204	914	694	494	309	346
CG7739-RA	CG7737-RA	3L:13621955...13622044	catt3	CG7739-RA	CG7737-RA	3L:13621955...13622044	89	715	548	226	490	173	375
jet-RA	beta-RA	2L:4948202...4948377	catt3	jet-RA	beta-RA	2L:4948202...4948377	175	656	530	398	258	334	196
hpo-b	Rep-RA	2R:1581300...1581300	catt3	hpo-b	Rep-RA	2R:1581300...1581300	130	689	529	277	466	230	99
CG12412-RA	CG12412-RA	3L:8192349...8192349	catt3	CG12412-RA	CG12412-RA	3L:8192349...8192349	272	1207	960	613	476	295	295
Hdax-X	NC	3L:2048045...2048045	catt3	Hdax-X	NC	3L:2048045...2048045	163	545	396	237	188	223	223
CG12795-RA	CG12795-RA	2L:2812949...2812949	catt3	CG12795-RA	CG12795-RA	2L:2812949...2812949	231	1255	933	300	285	209	223
CG16909-RA	CG16909-RA	2L:13782337...13782315	catt3	CG16909-RA	CG16909-RA	2L:13782337...13782315	278	1268	906	427	388	275	152
DphD-RA	aux-RA	2R:3405003...3405234	catt3	DphD-RA	aux-RA	2R:3405003...3405234	239	1254	985	485	370	257	196
CG12341-RA	mm5s4-RA	2R:6699585...6699729	catt3	CG12341-RA	mm5s4-RA	2R:6699585...6699729	144	574	420	214	361	190	245
CG6429-RA	ATg4-RA	2L:1038081...1038081	catt3	CG6429-RA	ATg4-RA	2L:1038081...1038081	299	1203	845	415	301	210	225
CG11180-RA	CG9945-RA	2R:16508401...16508604	catt3	CG11180-RA	CG9945-RA	2R:16508401...16508604	205	520	412	237	283	177	235
Golginin-RA	CG6454-RA	2L:20181741...20181741	catt3	Golginin-RA	CG6454-RA	2L:20181741...20181741	499	434	377	211	311	164	213
CG3332-RA	CG3332-RA	3R:10741716...10741747	catt3	CG3332-RA	CG3332-RA	3R:10741716...10741747	132	533	377	198	335	238	238
CG2811-RA	CG3736-RA	3L:20245873...20245873	catt3	CG2811-RA	CG3736-RA	3L:20245873...20245873	174	535	373	231	337	137	137
CG3397-RA	CG3397-RA	2L:20871334...20871565	catt3	CG3397-RA	CG3397-RA	2L:20871334...20871565	231	409	304				



Data sources used to compile the physical regulatory network

Data source	Name	modENCODE	External	PMID	doi	Regulatory proteins
1	Gary Karpen group	X	-	-	-	BEAF-32, CTCF, Chro, Cp190, E(z), Pc, Trl, mod(mdg4), su(Hw), trx, Trl
2	Kevin White group, Nègre et al.	X	20084099	10.1371/journal.pgen.1000814	cnc, bab1, brm, chinmo, CtBP, D, Dll, en, ftz-f1, gro, gsb-n, inv, kn, run, sbb, Snr1, Stat92E, ttk, Ubx, zfh1, BEAF-32, Cp190, CTCF, Trl, mod(mdg4), su(Hw), sens, Kr, CBP, cad, GATAe, disco, h, hkb, jumu, eve, EcR, E(z),	
3	Kwong et al.	X	18773083	10.1371/journal.pgen.1000178	Pc, pho	
4	Georlette et al.	X	17978103	10.1101/qad.1600107	Myb, mip120, mip130, E2f2, lin-52	
5	Schuettengruber et al.	X	19143474	10.1371/journal.pbio.1000013	Pc, ph-p, trx, pho, phol, Dsp1, Trl	
6	Moses et al.	X	17040121	10.1371/journal.pcbi.0020130	z	
7	Li et al.	X	18271625	10.1371/journal.pbio.0060027	bcd, cad, gt, hb, kni, Kr	
8	Lee et al.	X	18332113	10.1128/MCB.02224-07	Trl, Nelf-E, NELF-B	
9	MacArthur et al.	X	19627575	10.1186/gb-2009-10-7-r80	hkb, tll, D, ftz, h, prd, run, slp1, da, dl, Mad, Med, shn, sna, twi, TFIIB	
10	Bushey et al.	X	19443682	10.1101/qad.1798209	su(Hw), CTCF, BEAF-32	
11	Oktaba et al.	X	18993116	10.1016/j.devcel.2008.10.005	pho, Sfmbt	
12	Gambetta et al.	X	19478141	10.1126/science.1169727	ph-p, O-GlcNAc	
13	Sandmann et al.	X	16740481	10.1016/j.devcel.2006.04.009	Mef2	
14	Sandmann et al.	X	17322403	10.1101/qad.1509007	twi	
15	Jakobsen et al.	X	17908931	10.1101/qad.437607	bin, bap	
16	Zinzen et al.	X	19890324	10.1038/nature08531	tin, twi, bap, bin, mef2	

Table S5 - ChIP data sources used to create the physical (binding) Regulatory Network

<b>Term</b>	<b>Total predictions</b>	<b>KnownHits</b>	<b>Predictions WithAnnot</b>	<b>Predictions NoAnnot</b>
multicellular_organismal_process	7476	1578	2399	2447
regulation_of_biological_process	5339	1341	1366	1741
developmental_process	5056	1176	1455	1641
cellular_component_organization	5006	1143	1501	1600
regulation_of_cellular_process	4606	1114	1240	1510
macromolecule_metabolic_process	3751	1110	678	1230
cellular_macromolecule_metabolic_process	3020	917	539	959
multicellular_organismal_development	2972	775	764	917
anatomical_structure_development	1083	344	230	293
system_development	438	154	89	93
<b>organ_development</b>	<b>378</b>	<b>133</b>	<b>75</b>	<b>82</b>
translation	312	132	21	71
anatomical_structure_morphogenesis	266	81	76	55
<b>organ_morphogenesis</b>	<b>263</b>	<b>85</b>	<b>69</b>	<b>54</b>
cellular_developmental_process	233	62	81	51
<b>cell_differentiation</b>	<b>216</b>	<b>57</b>	<b>75</b>	<b>47</b>
generation_of_precursor_metabolites_and_energy	178	70	16	46
<b>nervous_system_development</b>	<b>193</b>	<b>63</b>	<b>56</b>	<b>35</b>
oxidative_phosphorylation	123	50	13	28
regulation_of_transcription_DNA-dependent	116	54	6	22
regulation_of_nitrogen_compound_metabolic_process	71	28	8	21
regulation_of_nucleobase_nucleoside_nucleotide_and_nucleic_acid_metabolic_process	71	28	8	21
regulation_of_biosynthetic_process	88	39	4	20
regulation_of_cellular_biosynthetic_process	88	39	4	20
energy_derivation_by_oxidation_of_organic_compounds	98	45	7	19
<b>imaginal_disc_development</b>	<b>140</b>	<b>50</b>	<b>40</b>	<b>18</b>
<b>pattern_specification_process</b>	<b>137</b>	<b>55</b>	<b>36</b>	<b>18</b>
<b>regulation_of_macromolecule_biosynthetic_process</b>	<b>85</b>	<b>39</b>	<b>3</b>	<b>18</b>
<b>regulation_of_transcription</b>	<b>86</b>	<b>39</b>	<b>4</b>	<b>18</b>
embryonic_development	139	49	45	17
regulation_of_cellular_metabolic_process	65	27	7	16
regulation_of_RNA_metabolic_process	68	27	9	15
regulation_of_gene_expression	46	17	7	14
<b>cellular_respiration</b>	<b>84</b>	<b>42</b>	<b>6</b>	<b>11</b>
mitotic_spindle_elongation	71	28	14	11
regulation_of_primary_metabolic_process	41	17	4	11
<b>spindle_elongation</b>	<b>71</b>	<b>28</b>	<b>14</b>	<b>11</b>
phosphorylation	43	20	2	10
post-embryonic_development	66	23	20	10
cellular_macromolecule_biosynthetic_process	34	16	2	8
instar_larval_or_pupal_development	47	16	16	8
macromolecule_biosynthetic_process	32	16	2	8
post-embryonic_organ_development	55	20	16	8
regulation_of_metabolic_process	14	3	2	8
fertilization	8	1	0	7
organic_anion_transport	8	1	0	7
regulation_of_macromolecule_metabolic_process	15	4	2	7
single_fertilization	8	1	0	7
cell_development	18	3	8	6
imaginal_disc_morphogenesis	37	13	12	6
instar_larval_or_pupal_morphogenesis	39	14	12	6
metamorphosis	39	14	12	6
post-embryonic_morphogenesis	39	14	12	6
post-embryonic_organ_morphogenesis	37	13	12	6
signaling	14	4	3	6
<b>muscle_structure_development</b>	<b>9</b>	<b>3</b>	<b>0</b>	<b>5</b>
regionalization	54	21	15	5
regulation_of_tube_size_open_tracheal_system	6	1	0	5
electron_transport_chain	32	21	1	4
sensory_organ_development	25	8	10	4
anion_transport	5	1	1	3
cellular_protein_metabolic_process	10	4	1	3
reproductive_cellular_process	10	4	1	3
actin_cytoskeleton_organization	4	1	1	2
actin_filament-based_process	4	1	1	2
glial_cell_differentiation	4	1	1	2
tissue_development	9	3	4	2
ATP_synthesis_coupled_electron_transport	24	18	1	1
anatomical_structureFormation_involved_in_morphogenesis	2	1	0	1
cell_communication	3	1	1	1
chorion-containing_eggshellFormation	2	1	0	1
eggshellFormation	2	1	0	1

establishment_of_localization	4	2	0	1
glial_cell_development	2	1	0	1
lipid_metabolic_process	2	1	0	1
mitochondrial_ATP_synthesis_coupled_electron_transport	27	21	2	1
muscle_cell_differentiation	2	1	0	1
protein_amino_acid_phosphorylation	2	1	0	1
regulation_of_immune_system_process	2	1	0	1
respiratory_electron_transport_chain	20	14	1	1
striated_muscle_cell_differentiation	2	1	0	1
transport	4	2	0	1
M_phase_of_mitotic_cell_cycle	1	1	0	0
Notch_signaling_pathway	1	1	0	0
cell-cell_signaling	4	2	1	0
cell_projection_organization	1	1	0	0
cellular_component_assembly	1	1	0	0
cellular_component_biogenesis	1	1	0	0
cytoskeleton_organization	1	1	0	0
establishment_or_maintenance_of_cell_polarity	1	1	0	0
localization	3	2	0	0
macromolecule_catabolic_process	2	1	1	0
mannose_metabolic_process	2	2	0	0
mitochondrial_electron_transport_NADH_to_ubiquinone	1	1	0	0
mitosis	1	1	0	0
myoblast_fusion	1	1	0	0
myotube_differentiation	1	1	0	0
nerve_development	1	1	0	0
neurological_system_process	2	1	1	0
neurotransmitter_transport	1	1	0	0
nuclear_division	1	1	0	0
organelle_fission	1	1	0	0
protein_catabolic_process	1	1	0	0
regulation_of_cell_shape	1	1	0	0
regulation_of_tube_size	1	1	0	0
response_to_DNA_damage_stimulus	1	1	0	0
signal_transmission	5	2	3	0
signaling_process	5	2	3	0
synaptic_transmission	4	2	1	0
syncytiumFormation	1	1	0	0
syncytiumFormation_by_plasma_membrane_fusion	1	1	0	0
system_process	2	1	1	0
terminal_region_determination	2	1	1	0
transmision_of_nerve_impulse	3	2	1	0
zygotic_determination_of_anterior/posterior_axis,_embryo	3	1	2	0

**Table S7: Relative performance of different networks for predicting expression of target genes.**

	10-fold cross-validation in time course						Predictable genes		Unpredictable genes		Cell line predictable	
	# of TFs with expression	# targets w/ expr	# predictors per target	# predictions better than random	# predictions worse than random	Enrichment over random network						
bRN	72	10248	13.6	1707	1642	1.04	656	6%	1,030	10%	296	45%
mRN	101	9065	7.3	1055	977	1.08	411	5%	214	2%	215	52%
kRN	74	82	2.3	13	3	3.91	0	-	0	-	0	-
fRN	545	8439	25.8	1992	805	2.47	1,487	18%	324	4%	932	63%

Order	Group	Immediate PI	First Name	I	Last Name	Contribution	Sections/Figures	Affiliation, Additional affiliation
<b>Integrative analysis</b>								
1	Kellis (DAC)	Kellis	Sushmita	Roy	Development of methods for network network inference. Application of networks for expression and function prediction.	7: Network inference and topology, 8: Gene function prediction and modules, 9: Predictive models of gene expression	3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 USA, . 2. Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02140, USA	
2	Kellis (DAC)	Kellis	Jason	Ernst	Prediction of stage-specific regulators. Chromatin state analysis. Chromatin mark ordering for state recovery analysis.	4: Chromatin states, 9: Predictive models of gene expression, 6: Genome coverage/conservation	3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 USA, . 2. Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02140, USA	
4	Kellis (DAC)	Kellis	Pouya	Kheradpour	Motif analysis in Hot spots, prediction of cell-type specific chromatin regulators	3: Chromatin marks, 5: TFs and Hotspots, Supplement/Methods	3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 USA, . 2. Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02140, USA	
8	Kellis (DAC)	Kellis	Christopher	A Bristow	Developed integrative model for promoter prediction in S2 cells, contributed to network analysis data integration, selection of validation experiments	3: Chromatin marks	3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 USA, . 2. Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02140, USA	
10	Kellis (DAC)	Kellis	Michael	F Lin	Analyzed coding potential in novel transcripts	2: Coding and non-coding genes	3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 USA, . 2. Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02140, USA	
11	Kellis (DAC)	Kellis	Stefan	Washietl	Analyzed structural properties of noncoding RNAs, co-authored a section in the maintext, a section in the supplement and a figure panel	2: Coding and non-coding genes	3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 USA, . 2. Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02140, USA	
13	Kellis (DAC)	Kellis	Ferhat	Ay	Topological and network motif analysis, contributed to network function prediction.	7: Network inference and topology, 8: Gene function prediction and modules, Supplement/Methods	3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 USA, . 2. Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02140, USA	
14	Kellis (DAC)	Kellis	Patrick	E Meyer	Functional network inference.	7: Network inference and topology	3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 USA, . 2. Machine Learning Group, Université Libre de Bruxelles, Boulevard du Triomphe CP212, 1050 Brussels, Belgium	
17	Kellis (DAC)	Kellis	Luisa	Di Stefano	Manuscript production, assembly, and coordination	Leadership/Organization	3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 USA, . 81. Massachusetts General Hospital Cancer Center, Harvard Medical School, Charlestown, MA, 02129, USA	
20	Kellis (DAC)	Kellis	Rogerio	Candeias	Construction of the dataset used to build the functional regulatory network, analysis of network properties, network validation methods	7: Network inference and topology, Supplement/Methods	3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 USA, . 2. Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02140, USA	
23	Kellis (DAC)	Kellis	Irwin	Jungreis	Used RNA-seq data to investigate Drosophila stop codon readthrough.	2: Coding and non-coding genes	3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 USA, . 2. Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02140, USA	
24	Kellis (DAC)	Kellis	Daniel	Marbach	Physical/functional regulatory network inference and analysis, gene function prediction, predictive models of expression, method discussion, analysis of results	7: Network inference and topology, 8: Gene function prediction and modules, 9: Predictive models of gene expression	3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 USA, . 2. Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02140, USA	
25	Kellis (DAC)	Kellis	Rachel	Sealfon	Predicting TF binding from motif combinations, enhancer prediction	5: TFs and Hotspots	3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 USA, . 2. Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02140, USA	
100	Kellis (DAC)	Kellis	Manolis	Kellis	Envisioned and directed the integrative analysis and co-wrote the paper	1: Dataset generation, validation, organization, 2: Coding and non-coding genes, 3: Chromatin marks, 4: Chromatin states, 5: TFs and Hotspots, 6: Genome coverage/conservation, 7: Network inference and topology, 8: Gene function prediction and modules, 9: Predictive models of gene expression, Intro/Discussion, Supplement/Methods, Leadership/Organization	3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 USA, . 2. Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02140, USA	
<b>Transcripts</b>								
7	Celniker	Celniker	Jane	M Landolin	Analyze primary data to generate novel gene and novel promoter lists for integrative analysis		68. Genomics Division, Lawrence Berkeley National Laboratory, 1 Cyclotron Road, Mailstop 64-121, Berkeley, California 94720 USA	
21	Celniker	Celniker	Joseph	W Carlson	cDNA verification experiment design and analysis	1: Dataset generation, validation, organization, 2: Coding and non-coding genes	7. Department of Genome Dynamics, Lawrence Berkeley National Laboratory, and Department of Molecular and Cell Biology, University of California, Berkeley, One Cyclotron Road, Berkeley, California 94720, USA	
31	Celniker	Celniker	Benjamin	W Booth	informatics for production	1: Dataset generation, validation, organization		
32	Celniker	Brenner	Angela	N Brooks	Junction alignments to RNA-Seq reads from development samples, alternative splicing analyses	2: Coding and non-coding genes, 6: Genome coverage/conservation	32. Department of Molecular and Cell Biology, University of California, Berkeley, California 94720, USA	
34	Celniker	Gingeras	Carrie	A Davis	Arrays and RNA-Seq was performed in the gingeras lab	1: Dataset generation, validation, organization, Leadership/Organization	40. Cold Spring Harbor Laboratory, Cold Spring Harbor, New York 11724, USA	
35	Celniker	Graveley	Michael	O Duff	Developmental time-course transcriptome RNA-seq data quantitation and analysis; alternative-splicing analysis	1: Dataset generation, validation, organization	41. Department of Genetics and Developmental Biology, University of Connecticut Stem Cell Institute, 263 Farmington, Connecticut 06030-6403, USA	
40	Celniker	Gingeras	Philipp	Kapranov	Supervised microarray experiments	2: Coding and non-coding genes	Affymetrix, Inc., 3380 Central Expyw, Santa Clara CA 95051	
52	Celniker	Perrimon	Anastasia	A Samsonova	RNAseq data analysis and functional annotation of ncRNAs	2: Coding and non-coding genes	42. Department of Genetics and Drosophila RNAi Screening Center, Harvard Medical School, 77 Avenue Louis Pasteur, Boston, Massachusetts 02115, USA	
53	Celniker	Celniker	Jeremy	E Sandler	Operated the Illumina genome sequencers, analyzed preliminary data as it came off the sequencers, and diagnosed sequencer problems. Performed reverse transcription, PCR amplification, pooling, and sequencing of 5' RACE products for transcription start site project.	1: Dataset generation, validation, organization	10. Department of Genome Dynamics, Lawrence Berkeley National Laboratory, Berkeley, California 94720, USA	
58	Celniker	Brent	Marijke	J van Baren	Created genome annotation	2: Coding and non-coding genes	38. Center for Genome Sciences and Systems Biology, Washington University, 4444 Forest Park Boulevard, Saint Louis, Missouri 63108, USA	
60	Celniker	Celniker	Kenneth	H Wan	Managed sequence data production, performed sequence data analysis	1: Dataset generation, validation, organization, 2: Coding and non-coding genes	10. Department of Genome Dynamics, Lawrence Berkeley National Laboratory, Berkeley, California 94720, USA	
61	Celniker	Graveley	Li	Yang	Making all the solexa libraries and doing sequencing analysis of drosophila developmental samples and cell line samples	1: Dataset generation, validation, organization	41. Department of Genetics and Developmental Biology, University of Connecticut Stem Cell Institute, 263 Farmington, Connecticut 06030-6403, USA	
62	Celniker	Celniker	Charles	Yu	Produced RACE products used to generate a subset of the 1.2 million RACE reads.	1: Dataset generation, validation, organization, Supplement/Methods	10. Department of Genome Dynamics, Lawrence Berkeley National Laboratory, Berkeley, California 94720, USA	
68	Celniker	Cherbas	Justen	Andrews	directed RNA production	1: Dataset generation, validation, organization	36. Department of Biology, Indiana University, 1001 E. 3rd Street, Bloomington, Indiana 47405-7005, USA	
71	Celniker	Graveley	Steven	E Brenner	Supervised development and engaged in analysis of splicing data	2: Coding and non-coding genes	32. Department of Molecular and Cell Biology, University of California, Berkeley, California 94720, USA, 35. Department of Plant & Microbial Biology, University of California, Berkeley, California 94720, USA	
72	Celniker	Brent	Michael	R Brent	Gene prediction, PCR, sequencing, and integrative annotation	2: Coding and non-coding genes	Center for Genome Sciences and Systems Biology, Washington University School of Medicine, Saint Louis, MO 63108	
73	Celniker	Cherbas	Lucy	Cherbas	RNA production, data analysis	1: Dataset generation, validation, organization	37. Center for Genomics and Bioinformatics, Indiana University, 1001 E. 3rd Street, Bloomington, Indiana 47405-7005, USA, 36. Department of Biology, Indiana University, 1001 E. 3rd Street, Bloomington, Indiana 47405-7005, USA	
75	Celniker	Gingeras	Thomas	R Gingeras	Carried out array and next-gen	1: Dataset generation, validation, organization, 2: Coding and non-coding genes, Supplement/Methods, Leadership/Organization	40. Cold Spring Harbor Laboratory, Cold Spring Harbor, New York, 11724 39. Affymetrix Inc, Santa Clara, California 95051, USA	
77	Celniker	Hoskins	Roger	A Hoskins	I conceived and managed the RACE/CAGE/EST fly embryo promoter analysis. I have played an important role in most aspects of the Celniker project.	1: Dataset generation, validation, organization, 2: Coding and non-coding genes	7. Department of Genome Dynamics, Lawrence Berkeley National Laboratory, and Department of Molecular and Cell Biology, University of California, Berkeley, One Cyclotron Road, Berkeley, California 94720, USA	
78	Celniker	Cherbas	Thomas	C Kaufman	Provided biological materials for the Drosophila project, directed collection of staged animals	1: Dataset generation, validation, organization	36. Department of Biology, Indiana University, 1001 E. 3rd Street, Bloomington, Indiana 47405-7005, USA	
82	Celniker	Perrimon	Norbert	Perrimon	co-directed the work	Leadership/Organization	42. Department of Genetics and Drosophila RNAi Screening Center, Harvard Medical School, 77 Avenue Louis Pasteur, Boston, Massachusetts 02115, USA, Howard Hughes Medical Institute	
87	Celniker	Cherbas	Peter	Cherbas	coordinated RNA production, contributed to data analysis	1: Dataset generation, validation, organization	37. Center for Genomics and Bioinformatics, Indiana University, 1001 E. 3rd Street, Bloomington, Indiana 47405-7005, USA, 36. Department of Biology, Indiana University, 1001 E. 3rd Street, Bloomington, Indiana 47405-7005, USA	
88	Celniker	Graveley	Brenton	R Graveley	Oversaw RNA-Seq data generation and analysis, aligned/analyzed RNA-Seq data, determined transcript models, identified new small RNAs/miRNA precursors	1: Dataset generation, validation, organization, 2: Coding and non-coding genes, 7: Network inference and topology	41. Department of Genetics and Developmental Biology, University of Connecticut Stem Cell Institute, 263 Farmington, Connecticut 06030-6403, USA	
93	Celniker	Celniker	Susan	E Celniker	Analyzed data, Made figures, wrote sections	1: Dataset generation, validation, organization, 2: Coding and non-coding genes, 5: TFs and Hotspots, 6: Genome coverage/conservation, 8: Gene function prediction and modules, 9: Predictive models of gene expression, Intro/Discussion, Supplement/Methods, Leadership/Organization	10. Department of Genome Dynamics, Lawrence Berkeley National Laboratory, Berkeley, California 94720, USA	
110	Celniker	Celniker	Charles	L Comstock	Developed and ran software for RT-PCR/sequencing	1: Dataset generation, validation, organization	59. Department of Biology CB-1137, Washington University, Saint Louis, Missouri 63130, USA	
115	Celniker	Gingeras	Alex	Dobin	contributed to junction identification	transcripts	40. Cold Spring Harbor Laboratory, Cold Spring Harbor, New York 11724, USA	
117	Celniker	Celniker	Jorg	Drenkow	RNA production	1: Dataset generation, validation, organization	40. Cold Spring Harbor Laboratory, Cold Spring Harbor, New York 11724, USA	
118	Celniker	Brenner	Sandrine	Dudoit	Development of statistical methods and software	Supplement/Methods	33. Division of Biostatistics, School of Public Health, University of California, Berkeley, California 94720, USA	
119	Celniker	Celniker	Jackie	Dumais	RNA production	1: Dataset generation, validation, organization	Affymetrix Inc., Santa Clara Ca.	
121	Celniker	Gingeras	Delphine	Fagegalter	production	small RNAs	40. Cold Spring Harbor Laboratory, Cold Spring Harbor, New York 11724, USA	
123	Celniker	Gingeras	Srinka	Ghosh	contributed in data generation and initial analysis as part of the Affymetrix group under Tom Gingeras	1: Dataset generation, validation, organization, 2: Coding and non-coding genes	Genentech Inc., 1 DNA Way,South San Francisco, CA 94080	

127	Celniker	Brenner	Kasper	D Hansen	splicing	alternative transcripts	33. Division of Biostatistics, School of Public Health, University of California, Berkeley, California 94720, USA.
131	Celniker	Celniker	Sonal	Jha	RNA production	1: Dataset generation, validation, organization	40. Cold Spring Harbor Laboratory, Cold Spring Harbor, New York 11724, USA.
137	Celniker	Brent	Laura	Langton	new genes	validation	59. Department of Biology CB-1137, Washington University, Saint Louis, Missouri 63130, USA.
141	Celniker	Gingeras	Wei	Lin	Tiling array transfrag Summary. RNA expression profiling and clustering by tiling array	1: Dataset generation, validation, organization, 2: Coding and non-coding genes	40. Cold Spring Harbor Laboratory, Cold Spring Harbor, New York 11724, USA.
150	Celniker	Celniker	David	Miller	RNA production	1: Dataset generation, validation, organization	36. Department of Biology, Indiana University, 1001 E. 3rd Street, Bloomington, Indiana 47400
166	Celniker	Brent	Aaron	E Tenney	Modified CONTRAST gene predictor to make use of RNAseq data when making gene predictions	2: Coding and non-coding genes	38. Center for Genome Sciences and Systems Biology, Washington University, 4444 Forest Park Boulevard, Saint Louis, Missouri 63108, USA.
169	Celniker	Celniker	Hualien	Wang	RNA production	1: Dataset generation, validation, organization	40. Cold Spring Harbor Laboratory, Cold Spring Harbor, New York 11724, USA.
172	Celniker	Gingeras	Aaron	T Willingham	tiling array data for RNA transcript maps	1: Dataset generation, validation, organization	39. Affymetrix Inc, Santa Clara, California 95051, USA.
173	Celniker	Gingeras	Chris	Zaleski	RNA expression	RNA expression	40. Cold Spring Harbor Laboratory, Cold Spring Harbor, New York 11724, USA.
175	Celniker	Cherbas	Dayu	Zhang	RNA production	1: Dataset generation, validation, organization	37. Center for Genomics and Bioinformatics, Indiana University, 1001 E. 3rd Street, Bloomington, Indiana 47405-7005, USA.
<b>Chromatin</b>							
3	Karpen	Park	Peter	V Kharchenko	Managed and analyzed fly chromatin group data, analyzed chromatin enrichments for different classes of elements	3: Chromatin marks, 4: Chromatin states, Supplement/Methods	61. Center for Biomedical Informatics, Harvard Medical School, 10 Shattuck St Boston, Massachusetts 02115, USA., 83. Informatics Program at Children's Hospital, 300 Longwood Avenue, Boston, MA 02115, USA
26	Karpen	Park	Michael	Y Tolstorukov	Analyzed positional correlation between chromatin properties and histone marks, handled chromatin physical properties data	3: Chromatin marks, 4: Chromatin states	61. Center for Biomedical Informatics, Harvard Medical School, 10 Shattuck St Boston, Massachusetts 02115, USA
28	Karpen	Kuroda	Artyom	A Alekseyenko	Antibody validation, chromatin preparation, ChIP-chip, QC	1: Dataset generation, validation, organization, 3: Chromatin marks, 4: Chromatin states, Supplement/Methods	60. Department of Medicine, Brigham & Women's Hospital, Department of Genetics, Harvard Medical School, 77 Avenue Louis Pasteur, Boston, Massachusetts 02115, USA
37	Karpen	Kuroda	Andrey	A Gorchakov	antigene design, antibody validation, chromatin preparation, ChIP-chip, QC	1: Dataset generation, validation, organization, 3: Chromatin marks, 4: Chromatin states	60. Department of Medicine, Brigham & Women's Hospital, Department of Genetics, Harvard Medical School, 77 Avenue Louis Pasteur, Boston, Massachusetts 02115, USA
38	Karpen	Elgin	Tingting	Gu	data production, data analysis	1: Dataset generation, validation, organization, 3: Chromatin marks, 4: Chromatin states	59. Department of Biology CB-1137, Washington University, Saint Louis, Missouri 63130, USA
45	Karpen	Karpen	Aki	Minoda	Management and setting up of ChIP-chip data production	1: Dataset generation, validation, organization, 3: Chromatin marks	7. Department of Genome Dynamics, Lawrence Berkeley National Laboratory, and Department of Molecular and Cell Biology, University of California, Berkeley, One Cyclotron Road, Berkeley, California 94720, USA.
50	Karpen	Elgin	Nicole	C Riddle	data production, data analysis, supervision	1: Dataset generation, validation, organization, 3: Chromatin marks, 4: Chromatin states	59. Department of Biology CB-1137, Washington University, Saint Louis, Missouri 63130, USA
54	Karpen	Pirrota	Yuri	B Schwartz	antibody generation and characterization, chip-chip, chip-seq, primary data analysis and quality control	1: Dataset generation, validation, organization, 3: Chromatin marks	64. Department of Molecular Biology and Biochemistry, Rutgers University, Piscataway, New Jersey 08854, USA.
74	Karpen	Elgin	Sarah	C Elgin R	Co-PI Karpen project	1: Dataset generation, validation, organization, 3: Chromatin marks, 4: Chromatin states	59. Department of Biology CB-1137, Washington University, Saint Louis, Missouri 63130, USA.
80	Karpen	Kuroda	Mitzi	I Kuroda	Mapping of chromatin proteins and histone modifications	1: Dataset generation, validation, organization, 3: Chromatin marks, 4: Chromatin states	60. Department of Medicine, Brigham & Women's Hospital, Department of Genetics, Harvard Medical School, 77 Avenue Louis Pasteur, Boston, Massachusetts 02115, USA.
83	Karpen	Pirrota	Vincenzo	Pirrota	generated and evaluated data	1: Dataset generation, validation, organization, 3: Chromatin marks, 4: Chromatin states	64. Department of Molecular Biology and Biochemistry, Rutgers University, Piscataway, New Jersey 08854, USA.
92	Karpen	Karpen	Peter	J Park	Co-directed chromatin-related sections, overall contribution to analytical aspects	3: Chromatin marks, 4: Chromatin states, Intro/Discussion	61. Center for Biomedical Informatics, Harvard Medical School, 10 Shattuck St Boston, Massachusetts 02115, USA.
95	Karpen	Karpen	Gary	H Karpen	PI for the chromatin production group, substantial writing and editing, Figure 4 and Supp 4, overall leadership	1: Dataset generation, validation, organization, 3: Chromatin marks, 4: Chromatin states, 5: TFS and Hotspots, 6: Genome coverage/conservation, Intro/Discussion, Leadership/Organization	7. Department of Genome Dynamics, Lawrence Berkeley National Laboratory, and Department of Molecular and Cell Biology, University of California, Berkeley, One Cyclotron Road, Berkeley, California 94720, USA.
101	Karpen	Karpen	David	Acevedo	ChIP	1: Dataset generation, validation, organization	7. Department of Genome Dynamics, Lawrence Berkeley National Laboratory, and Department of Molecular and Cell Biology, University of California, Berkeley, One Cyclotron Road, Berkeley, California 94720, USA.
105	Karpen	Park	Eric	P Bishop	Analysis of GRO-seq data	4: Chromatin states	61. Center for Biomedical Informatics, Harvard Medical School, 10 Shattuck St Boston, Massachusetts 02115, USA.
122	Karpen	Elgin	Sarah	E Gadel	Antibody Characterization	1: Dataset generation, validation, organization, 3: Chromatin marks, 4: Chromatin states	59. Department of Biology CB-1137, Washington University, Saint Louis, Missouri 63130, USA
133	Karpen	Park	Youngsook	L Jung	analysis of 4th chromosome and other chromatin related topics	3: Chromatin marks	61. Center for Biomedical Informatics, Harvard Medical School, 10 Shattuck St Boston, Massachusetts 02115, USA.
135	Karpen	Karpen	Cameron	D Kennedy	ChIP	1: Dataset generation, validation, organization	7. Department of Genome Dynamics, Lawrence Berkeley National Laboratory, and Department of Molecular and Cell Biology, University of California, Berkeley, One Cyclotron Road, Berkeley, California 94720, USA.
138	Karpen	Karpen	Ok-Kyung	Lee	production	1: Dataset generation, validation, organization	60. Department of Medicine, Brigham & Women's Hospital, Department of Genetics, Harvard Medical School, 77 Avenue Louis Pasteur, Boston, Massachusetts 02115, USA.
142	Karpen	Karpen	Daniela	Linder-Basso	production	1: Dataset generation, validation, organization	64. Department of Molecular Biology and Biochemistry, Rutgers University, Piscataway, New Jersey 08854, USA.
145	Karpen	Elgin	Sarah	E Marchetti	ChIP-chip, chromatin preps	1: Dataset generation, validation, organization, 3: Chromatin marks, 4: Chromatin states	59. Department of Biology CB-1137, Washington University, Saint Louis, Missouri 63130, USA.
162	Karpen	Karpen	Gregory	Shanower	production	1: Dataset generation, validation, organization	64. Department of Molecular Biology and Biochemistry, Rutgers University, Piscataway, New Jersey 08854, USA.
<b>Transcription Factors</b>							
5	White	White	Nicolas	Nègre	Generated TF ChIP-chip data, co-directed the analysis of TF overlaps and HOT regions	4: Chromatin states, 5: TFS and Hotspots	5. Institute for Genomics & Systems Biology, Department of Human Genetics, The University of Chicago, 900 East 57th Street, Chicago, IL 60637, USA.
9	White	White	Ljija	Ma	Analyzed the co-localization of TFS	5: TFS and Hotspots	5. Institute for Genomics & Systems Biology, Department of Human Genetics, The University of Chicago, 900 East 57th Street, Chicago, IL 60637, USA.
19	White	White	Christopher	D Brown	Data analysis on Transcription Factor data, HDAC binding data analysis, transcript assembly/quantification, TF binding site analysis	5: TFS and Hotspots, 8: Gene function prediction and modules, 9: Predictive models of gene expression	5. Institute for Genomics and Systems Biology, Department of Human Genetics, The University of Chicago, 900 E. 57th ST., Chicago, IL 60637, USA
56	White	White	Rebecca	F Spokony	produced TF binding datasets using tagging strategy		
76	White	Grossman	Robert	L Grossman	led the development of some of the informatics infrastructure used by the project; analyzed some of the data with my students	informatics infrastructure	84. Institute for Genomics & Systems Biology, University of Chicago, 900 East 57th Street, Chicago, IL 60637 , 49. National Center for Data Mining, University of Illinois at Chicago, 700 Science and Engineering Offices, MC 249, 851 South Morgan Street, Chicago, Illinois 60607, USA.
84	White	Posakony	James	W Posakony	Supervised biological validation effort; analyzed data; helped write CisReg paper		55. Division of Biological Sciences, Section of Cell & Developmental Biology, University of California San Diego, 9500 Gilman Drive, La Jolla, California 92093 USA.
85	White	Ren	Bing	Ren	coPI, author		
86	White	Russell	Steven	Russell	validation of insulators and silencers, analysis of TF binding and chromatin mark data	3: Chromatin marks, 5: TFS and Hotspots, Leadership/Organization	47. Department of Genetics and Cambridge Systems Biology Centre, University of Cambridge, Downing Street, Cambridge, CB2 3EH, UK.
99	White	White	Kevin	P White	PI, production, 5ab, 7a		5. Institute for Genomics & Systems Biology, Department of Human Genetics, The University of Chicago, 900 East 57th Street, Chicago, IL 60637, USA.
102	White	Russell	Richard	Auburn	analysis of primary ChIP data	1: Dataset generation, validation, organization	5. Institute for Genomics & Systems Biology, Department of Human Genetics, The University of Chicago, 900 East 57th Street, Chicago, IL 60637, USA.
104	White	Bellen	Hugo	J. Bellen	contributed Pacman clones for TF tagging	TF datasets	48. Department of Molecular and Human Genetics, Howard Hughes Medical Institute, Baylor College of Medicine, Houston, TX, USA.
108	White	Grossman	Jia	Chen	TF primary data analysis/management	1: Dataset generation, validation, organization	47. Department of Genetics and Cambridge Systems Biology Centre, University of Cambridge, Downing Street, Cambridge, CB2 3EH, UK.
116	White	White	Marc	H Domann	performed sequencing	TF datasets	6. Institute for Genomics & Systems Biology, Argonne National Laboratory, 9700 South Cass Ave, Lemont, IL, USA.
125	White	Grossman	David	Harley	TF primary data analysis/management	1: Dataset generation, validation, organization	79. Genome Sciences Centre, BC Cancer Agency, Vancouver Canada.
128	White	White	Elizabeth	Heinz	TF primary data analysis/management	1: Dataset generation, validation, organization	6. Institute for Genomics & Systems Biology, Argonne National Laboratory, 9700 South Cass Ave, Lemont, IL, USA.
140	White	Ren	Zirong	Li	TF primary data analysis/management	1: Dataset generation, validation, organization	6. Institute for Genomics & Systems Biology, Argonne National Laboratory, 9700 South Cass Ave, Lemont, IL, USA.
149	White	Meyer	Folker	Meyer	TF primary data analysis/management	1: Dataset generation, validation, organization	51. Ludwig Institute for Cancer Research, University of California San Diego Biology Division, La Jolla, CA, USA.
151	White	Posakony	Steven	W Miller	Analysis of TF binding and histone modification data across development toward functional validation of promoters and enhancers.	1: Dataset generation, validation, organization, 3: Chromatin marks, 5: TFS and Hotspots	55. Division of Biological Sciences, Section of Cell & Developmental Biology, University of California San Diego, 9500 Gilman Drive, La Jolla, California 92093 USA.
153	White	White	Carolyn	A Morrison	ChIP data production	1: Dataset generation, validation, organization	53. Mathematics and Computer Science Division, Argonne National Laboratory, 9700 South Cass Ave, Lemont, IL, USA.
159	White	White	Douglas	A Schefner	Part of team that was responsible for data analysis and data organization	1: Dataset generation, validation, organization	5. Institute for Genomics & Systems Biology, Department of Human Genetics, The University of Chicago, 920 E. 58th Street, Chicago, Illinois 60637, USA.
160	White	White	Lionel	Senderowicz	Ab production	1: Dataset generation, validation, organization	22. Basic Sciences Division, Fred Hutchinson Cancer Research Center, 1100 Fairview Avenue North, Seattle, Washington 98109, USA.
161	White	White	Parantu	K Shah	Quality control and primary analysis	N/A	5. Institute for Genomics & Systems Biology, Department of Human Genetics, The University of Chicago, 920 E. 58th Street, Chicago, Illinois 60637, USA.
165	White	White	Sarah	Suchy	ChIP data production	1: Dataset generation, validation, organization	5. Institute for Genomics & Systems Biology, Department of Human Genetics, The University of Chicago, 920 E. 58th Street, Chicago, Illinois 60637, USA.
167	White	Grossman	Feng	Tian	TF primary data analysis/management	1: Dataset generation, validation, organization	5. Institute for Genomics & Systems Biology, Department of Human Genetics, The University of Chicago, 920 E. 58th Street, Chicago, Illinois 60637, USA.
168	White	Bellen	Koen	J Venken	contributed Pacman clones for TF tagging	1: Dataset generation, validation, organization	48. Department of Molecular and Human Genetics, Howard Hughes Medical Institute, Baylor College of Medicine, Houston, TX, USA.
170	White	White	Robert	White	analysis of primary ChIP data	1: Dataset generation, validation, organization	57. Department of Molecular and Human Genetics, Baylor College of Medicine, One Baylor Plaza, Houston, TX, USA.
171	White	Meyer	Jared	Wilkenning	TF primary data analysis/management	1: Dataset generation, validation, organization	58. Department of Physiology, Development and Neuroscience, University of Cambridge, Cambridge, UK.
177	White	White	Jennifer	T Zieba	ChIP data production	1: Dataset generation, validation, organization	5. Institute for Genomics & Systems Biology, Department of Human Genetics, The University of Chicago, 920 E. 58th Street, Chicago, Illinois 60637, USA.
<b>Replication</b>							
6	MacAlpine	MacAlpine	Matthew	L Eaton	Integrative analysis in the context of replication data	3: Chromatin marks, 5: TFS and Hotspots	65. Department of Pharmacology and Cancer Biology, Duke University Medical Center, Durham, North Carolina 27710, USA.
43	MacAlpine	MacAlpine	Heather	K MacAlpine	Data production; generated ORC ChIP-seq data, replication timing, and early origin data	1: Dataset generation, validation, organization	65. Department of Pharmacology and Cancer Biology, Duke University Medical Center, Durham, North Carolina 27710, USA.
46	MacAlpine	Orr-Weaver	Jared	T Nordman	aCGH and RNA-seq analysis of differential replication in the larval fat body and midgut tissues.	Chromatin-based annotation of functional elements	66. Whitehead Institute, Cambridge, Massachusetts 02142, USA.
49	MacAlpine	MacAlpine	Sara	K Powell	generated data for copy number analysis of cell lines	CNV section	65. Department of Pharmacology and Cancer Biology, Duke University Medical Center, Durham, North Carolina 27710, USA.
55	MacAlpine	Orr-Weaver	Noa	Sher	Worked on salivary gland CGH arrays	Data generation of replication data	66. Whitehead Institute, Cambridge, Massachusetts 02142, USA.
81	MacAlpine	Orr-Weaver	Terry	L Orr-Weaver	Co-PI for replication group, analyzed data on differential DNA replication in differentiated tissues, contributed sentences to manuscript.	Replication section within chromatin based annotation of functional elements	66. Whitehead Institute, Cambridge, Massachusetts 02142, USA.
97	MacAlpine	MacAlpine	David	M MacAlpine	PI for replication group, data production, analysis, and writing	3: Chromatin marks, 4: Chromatin states, 5: TFS and Hotspots, Intro/Discussion, Leadership/Organization	65. Department of Pharmacology and Cancer Biology, Duke University Medical Center, Durham, North Carolina 27710, USA.
113	MacAlpine	MacAlpine	Leyna	C DeNapoli	production	replication	65. Department of Pharmacology and Cancer Biology, Duke University Medical Center, Durham, North Carolina 27710, USA.
114	MacAlpine	MacAlpine	Quyen	Ding	production	replication	65. Department of Pharmacology and Cancer Biology, Duke University Medical Center, Durham, North Carolina 27710, USA.
120	MacAlpine	Orr-Weaver	Thomas	Eng	production	1: Dataset generation, validation, organization	66. Whitehead Institute, Massachusetts Institute of Technology, Cambridge, Massachusetts 02142, USA.
134	MacAlpine	Orr-Weaver	Helena	Kashevsky	production	1: Dataset generation, validation, organization	66. Whitehead Institute, Massachusetts Institute of Technology, Cambridge, Massachusetts 02142, USA.

139	MacAlpine	Orr-Weaver	Sharon	Li	production	1. Dataset generation, validation, organization	66. Whitehead Institute, Massachusetts Institute of Technology, Cambridge, Massachusetts USA
155	MacAlpine	MacAlpine	Joseph	A Prinz	production pipeline and DCC submission	1. Dataset generation, validation, organization	65. Department of Pharmacology and Cancer Biology, Duke University Medical Center, Durham, NC, USA
<b>Small RNAs</b>							
15	Lai	Lai	Nicolas	Robine	Analyzed small RNA data obtained.	2: Coding and non-coding genes	44. Sloan-Kettering Institute, 1275 York Avenue, Box 252, New York, New York 10065, USA
18	Lai	Lai	Eugene	Berezikov	Identification of novel miRNAs.	2: Coding and non-coding genes	46. Hubrecht Institute, Royal Netherlands Academy of Arts and Sciences & University Medical Center Utrecht, Utrecht, The Netherlands
33	Lai	Lai	Qi	Dai	Small RNA library construction	2: Coding and non-coding genes	44. Sloan-Kettering Institute, 1275 York Avenue, Box 252, New York, New York 10065, USA
47	Lai	Lai	Katsutomo	Okamura	Small RNA library construction	2: Coding and non-coding genes	44. Sloan-Kettering Institute, 1275 York Avenue, Box 252, New York, New York 10065, USA
96	Lai	Lai	Eric	C Lai	Directed small RNA analysis, assembled Fig2C and Table S2, helped write paper.	2: Coding and non-coding genes	44. Sloan-Kettering Institute, 1275 York Avenue, Box 252, New York, New York 10065, USA
126	Lai	Lai	Gregory	J Hannon	small RNA sequencing	2: Coding and non-coding genes	45. Cold Spring Harbor Laboratory, Howard Hughes Medical Institute, Cold Spring Harbor, New York 11724, USA
156	Lai	Lai	Michelle	Rooks	small RNA sequencing	2: Coding and non-coding genes	45. Cold Spring Harbor Laboratory, Howard Hughes Medical Institute, Cold Spring Harbor, New York 11724, USA
130	Lai	Lai	Martin	Hirst	small RNA sequencing	2: Coding and non-coding genes	79. Genome Sciences Centre, BC Cancer Agency, Vancouver Canada
146	Lai	Lai	Marco	Marra	small RNA sequencing	2: Coding and non-coding genes	79. Genome Sciences Centre, BC Cancer Agency, Vancouver Canada
156	Lai	Lai	Richard	Moore	small RNA sequencing	2: Coding and non-coding genes	79. Genome Sciences Centre, BC Cancer Agency, Vancouver Canada
176	Lai	Lai	Yongjun	Zhao	small RNA sequencing	2: Coding and non-coding genes	79. Genome Sciences Centre, BC Cancer Agency, Vancouver Canada
<b>Nucleosomes</b>							
39	Henikoff	Henikoff	Jorja	G Henikoff	Gene ends analysis of NimbleGen data.	1: Dataset generation, validation, organization	44. Sloan-Kettering Institute, 1275 York Avenue, Box 252, New York, New York 10065, USA
51	Henikoff	Ahmad	Akiko	Sakai	performed salt fractations, mapping of mononucleosome/H2Av, discussed methods/results	3: Chromatin marks, 4: Chromatin states	23. Department of Biological Chemistry and Molecular Pharmacology, Harvard Medical School, 240 Longwood Avenue, Boston, Massachusetts 02115, USA
67	Henikoff	Ahmad	Kami	Ahmad	Generation of salt fraction chromatin profiles for Drosophila cell lines	4: Chromatin states	46. Hubrecht Institute, Royal Netherlands Academy of Arts and Sciences & University Medical Center Utrecht, Utrecht, The Netherlands
94	Henikoff	Henikoff	Steven	Henikoff	PI, concalls	4: Chromatin states	22. Basic Sciences Division, Fred Hutchinson Cancer Research Center, 1100 Fairview Avenue North, Seattle, Washington 98109, USA
106	Henikoff	Henikoff	Terri	D Bryson	production	1: Dataset generation, validation, organization	22. Basic Sciences Division, Fred Hutchinson Cancer Research Center, 1100 Fairview Avenue North, Seattle, Washington 98109, USA
<b>Data Coordination Center</b>							
12	Stein (DCC)	Stein	Bradley	I Arshinoff	Analyzed genome coverage scaling.	6: Genome coverage/conservation	74. Department of Molecular Genetics, University of Toronto, 27 King's College Circle, Toronto, Ontario M5S 1A1, Canada, 73. Ontario Institute for Cancer Research, 101 College Street, Suite 800, Toronto Ontario M5G 0A3, Canada
16	Stein (DCC)	Lewis	Nicole	L Washington	Pipeline design, curation, data organization; Figure S1 and Table S1	1: Dataset generation, validation, organization, Intro/Discussion	68. Genomics Division, Lawrence Berkeley National Laboratory, 1 Cyclotron Road, Mailstop 64-121, Berkeley, California 94720 USA
22	Stein (DCC)	Micklem	Adrian	Carr	Conservation analysis/statistics; modMine	6: Genome coverage/conservation	47. Department of Genetics and Cambridge Systems Biology Centre, University of Cambridge, Downing Street, Cambridge, CB2 3EH, UK
36	Stein (DCC)	Stein	Xin	Feng			
48	Stein (DCC)	Stein	Marc	D Perry	Collecting metadata, validating data sets, ensuring entries for ChIP-chip/-seq data were correct	1: Dataset generation, validation, organization, Figure 1b	73. Ontario Institute for Cancer Research, 101 College Street, Suite 800, Toronto Ontario M5G 0A3, Canada
79	Stein (DCC)	Kent	William	J Kent	Managed UCSC folks who developed initial ruby on rails pipeline and comparative genomics alignments and display.	1: Dataset generation, validation, organization, 6: Genome coverage/conservation	Other, 445 Palm St
89	Stein (DCC)	Lewis	Suzanna	E Lewis	Data set provisioning	Supplement/Methods, Leadership/Organization	68. Genomics Division, Lawrence Berkeley National Laboratory, 1 Cyclotron Road, Mailstop 64-121, Berkeley, California 94720 USA
90	Stein (DCC)	Micklem	Gos	Micklem	Conservation analysis/ stats, DCC/modMine, supervised Cambridge DCC contributions	1: Dataset generation, validation, organization, 6: Genome coverage/conservation	47. Department of Genetics and Cambridge Systems Biology Centre, University of Cambridge, Downing Street, Cambridge, CB2 3EH, UK
98	Stein (DCC)	STEIN	Lincoln	D Stein	Heavy editing of draft; redid figure 1b; supervised Brad; organized companion pages	6: Genome coverage/conservation, Leadership/Organization	73. Ontario Institute for Cancer Research, 101 College Street, Suite 800, Toronto Ontario M5G 0A3, Canada, 74. Department of Molecular Genetics, University of Toronto, 27 King's College Circle, Toronto, Ontario M5S 1A1, Canada
103	Stein (DCC)	Kent	Galt	Barber	Wrote initial modencode submission pipeline front-end	web infrastructure	85. Center for Biomolecular Science and Engineering, School of Engineering, University of California Santa Cruz (UCSC), Santa Cruz, CA 95064, USA
107	Stein (DCC)	Micklem	Aurelien	Chateigner	Preliminary work on conservation analysis	6: Genome coverage/conservation	47. Department of Genetics and Cambridge Systems Biology Centre, University of Cambridge, Downing Street, Cambridge, CB2 3EH, UK
109	Stein (DCC)	Kent	Hiram	Clawson	5-way worm alignment	6: Genome coverage/conservation	69. Department of Biomolecular Engineering, University of California, Santa Cruz, California 95064, USA
111	Stein (DCC)	Micklem	Sergio	Contrino	modMine	1: Dataset generation, validation, organization	47. Department of Genetics and Cambridge Systems Biology Centre, University of Cambridge, Downing Street, Cambridge, CB2 3EH, UK
124	Stein (DCC)	Micklem	Francois	Guillier	DCC infrastructure/ modMine	1: Dataset generation, validation, organization	47. Department of Genetics and Cambridge Systems Biology Centre, University of Cambridge, Downing Street, Cambridge, CB2 3EH, UK
129	Stein (DCC)	Kent	Angle	S Hinrichs	multiz alignments, phastCons conservation, ported UCSC display code for alignments to Gbrowse	1: Dataset generation, validation, organization, 6: Genome coverage/conservation	85. Center for Biomolecular Science and Engineering, School of Engineering, University of California Santa Cruz (UCSC), Santa Cruz, CA 95064, USA
136	Stein (DCC)	Stein	Ellen	T Kephart	assisted in IDR analysis that was ultimately not included	1: Dataset generation, validation, organization	73. Ontario Institute for Cancer Research, 101 College Street, Suite 800, Toronto Ontario M5G 0A3, Canada
143	Stein (DCC)	Stein	Paul	Lloyd	companion page / data processing	Supplement/Methods, companion page	73. Ontario Institute for Cancer Research, 101 College Street, Suite 800, Toronto Ontario M5G 0A3, Canada
144	Stein (DCC)	Micklem	Rachel	Lyne	modMine. Preliminary work on conservation analysis	6: Genome coverage/conservation	47. Department of Genetics and Cambridge Systems Biology Centre, University of Cambridge, Downing Street, Cambridge, CB2 3EH, UK
148	Stein (DCC)	Stein	Sheldon	McKay	modENCODE genome browser	1: Dataset generation, validation, organization	67. Cold Spring Harbor Laboratory, 1 Bungtown Road, Cold Spring Harbor, New York 11542 USA
152	Stein (DCC)	Marra	Richard	A Moore	small RNA sequencing	1: Dataset generation, validation, organization, 3: Chromatin marks	79. Genome Sciences Centre, BC Cancer Agency, Vancouver Canada
157	Stein (DCC)	Micklem	Kim	M Rutherford	Helped develop the DCC InterMine data warehouse	1: Dataset generation, validation, organization	4. Department of Genetics, University of Cambridge, CB2 3EH, UK; Cambridge Systems Biology Centre, Tennis Court Road, Cambridge CB2 1QR, UK
158	Stein (DCC)	Stein	Peter	Ruzanov	worked on visualization of modENCODE data in DCC Browser	1: Dataset generation, validation, organization	73. Ontario Institute for Cancer Research, 101 College Street, Suite 800, Toronto Ontario M5G 0A3, Canada
163	Stein (DCC)	Micklem	Richard	Smith	modMine	1: Dataset generation, validation, organization	47. Department of Genetics and Cambridge Systems Biology Centre, University of Cambridge, Downing Street, Cambridge, CB2 3EH, UK
164	Stein (DCC)	Lewis	E.O.	Stinson	Created and maintained most of the submission pipeline	1: Dataset generation, validation, organization	68. Genomics Division, Lawrence Berkeley National Laboratory, 1 Cyclotron Road, Mailstop 64-121, Berkeley, California 94720 USA
174	Stein (DCC)	Stein	Zheng	Zha	submitted more than 200 GEO submissions	1: Dataset generation, validation, organization	73. Ontario Institute for Cancer Research, 101 College Street, Suite 800, Toronto Ontario M5G 0A3, Canada
<b>Comparative Transcription</b>							
29	Oliver	Oliver	Carlo	G Artieri	Generated list of novel transcribed regions based on RNA-seq data. Analyzed their patterns of expression in heads of multiple <i>Drosophila</i> species	1: Dataset generation, validation, organization, 2: Coding and non-coding genes	71. Section of Developmental Genomics, Laboratory of Cellular and Developmental Biology, National Institute of Diabetes and Digestive and Kidney Diseases, National Institutes of Health, Bethesda MD 20892, USA
42	Oliver	Oliver	Renhua	Li	data management and data submission to DCC	1: Dataset generation, validation, organization	8. Division of Extramural Research, National Human Genome Research Institute, National Institutes of Health, 5635 Fishers Lane, Suite 4076, Bethesda, Maryland 20892-9305, USA, , Developmental Genomics Section, LDCB, NIDDK, National Institutes of Health, Bethesda, MD 20892
44	Oliver	Oliver	John	H Malone	Evolutionary analysis of gene expression in <i>Drosophila</i> species	1: Dataset generation, validation, organization, 2: Coding and non-coding genes, Supplement/Methods	71. Section of Developmental Genomics, Laboratory of Cellular and Developmental Biology, National Institute of Diabetes and Digestive and Kidney Diseases, National Institutes of Health, Bethesda MD 20892, USA
57	Oliver	Oliver	David	Sturgill	Helped analyze the developmental timecourse data	1: Dataset generation, validation, organization, 2: Coding and non-coding genes	71. Section of Developmental Genomics, Laboratory of Cellular and Developmental Biology, National Institute of Diabetes and Digestive and Kidney Diseases, National Institutes of Health, Bethesda MD 20892, USA
91	Oliver	Oliver	Brian	Oliver	RNA-Seq	1: Dataset generation, validation, organization, 2: Coding and non-coding genes, 6: Genome coverage/conservation, Intro/Discussion, Leadership/Organization	71. Section of Developmental Genomics, Laboratory of Cellular and Developmental Biology, National Institute of Diabetes and Digestive and Kidney Diseases, National Institutes of Health, Bethesda MD 20892, USA
132	Oliver	Oliver	Lichun	Jiang	RNA-Seq QC experiments	1: Dataset generation, validation, organization	71. Section of Developmental Genomics, Laboratory of Cellular and Developmental Biology, National Institutes of Health, Bethesda MD 20892, USA
147	Oliver	Oliver	Nicolas	R Mattiuzzo	Maternal, paternal expression, RNA-Seq QC experiments	1: Dataset generation, validation, organization	71. Section of Developmental Genomics, Laboratory of Cellular and Developmental Biology, National Institutes of Health, Bethesda MD 20892, USA
<b>RNA families</b>							
27	Other	Berger	Sebastian	Will	Analysis of structural properties in ncRNAs, co-authored main section 2c, contributed panel to figure 2, supplementary text and supplementary figure.	2: Coding and non-coding genes	3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA
69	Other	Berger	Bonnie	Berger	Analysis of structural properties in ncRNAs, co-authored main section 2c, contributed panel to figure 2, supplementary text and supplementary figure.	2: Coding and non-coding genes	86. Department of Mathematics, Massachusetts Institute of Technology, Cambridge MA 02139, USA, 3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA
<b>Program Management</b>							
63	Other	NHGRI	Elise	A Feingold	Project Management	Leadership/Organization	8. Division of Extramural Research, National Human Genome Research Institute, National Institutes of Health, 5635 Fishers Lane, Suite 4076, Bethesda, Maryland 20892-9305, USA
64	Other	NHGRI	Peter	J Good	Project Management	Leadership/Organization	8. Division of Extramural Research, National Human Genome Research Institute, National Institutes of Health, 5635 Fishers Lane, Suite 4076, Bethesda, Maryland 20892-9305, USA
65	Other	NHGRI	Mark	S Guyer	Project Management	Leadership/Organization	8. Division of Extramural Research, National Human Genome Research Institute, National Institutes of Health, 5635 Fishers Lane, Suite 4076, Bethesda, Maryland 20892-9305, USA
66	Other	NHGRI	Rebecca	F Lowdon	Project Management	Leadership/Organization	8. Division of Extramural Research, National Human Genome Research Institute, National Institutes of Health, 5635 Fishers Lane, Suite 4076, Bethesda, Maryland 20892-9305, USA

Table S8 - Complete Author List with Author Contributions