

Additional file 2: Supplementary materials

These supplementary materials accompany the manuscript “Genome-wide upstream motif and DNA-binding analysis of *Cryptosporidium parvum* genes clustered by expression profile” by Oberstaller et al. Supplementary Figures S1-S8 appear in this document. Supplementary Tables 1-4 appear in the spreadsheet Additional file 1.

Results

Motif 5, which is rich in G's and A's, is similar to the 5'-GAGA-3' motif identified in *Drosophila* [1], as well as the NTBA (NucleoTide Biosynthesis-A) motif identified upstream of the nucleotide biosynthetic genes in *T. gondii* [2]. No proteins similar to the known GAGA-binding family of transcription factors are annotated in the *C. parvum* genome. The GAGA motif is found to be significantly overrepresented in the upstream regions of genes in 12 co-expressed gene clusters (totaling 127 genes) in the current study (Figure S1 A-D; Additional file 1, Table S3). Important GO terms enriched in these clusters are nucleosome assembly, DNA dependent transcription initiation, protein dephosphorylation and ubiquitin-dependent protein catabolic process.

Motifs 13, 16 and 20 are similar to the CCAAT-box (5'-CCAAT-3') motif [3]. The *C. parvum* genome contains three putative CCAAT-binding transcription factors (Table 2). These motifs are overrepresented in the upstream regions of genes present in 21 gene clusters (228 genes) that are maximally expressed at any of the seven time points (Figure S2 A-D), with biological process enrichments in processes including ATP synthesis coupled proton transport, intracellular protein transport, oxidative phosphorylation, cytoskeletal organization and microtubule-based movement.

Motifs 10, 12, 15, 17, 18, and 19 comprise the family we term “Unknown set 1” and do not appear similar to any known *cis*-regulatory motifs. Several of these motifs are overrepresented upstream of single clusters (Figure S3 A-D). Motif 17 was found only in the promoter region of the genes present in cluster 6, which is enriched with ribosomal proteins. GO enrichment analysis of the genes in cluster 6 revealed highly significant enrichment of GO terms associated with gene expression, translation, translational elongation and tRNA aminoacylation. Motifs 18 and 19 are each conserved in the upstream region of clusters, 21 and 143 respectively. The genes in cluster 21 do not show any specific GO-enrichment. Genes in cluster 143 are enriched with ubiquitin-dependent protein catabolic process and protein dephosphorylation GO terms.

Motifs 9 and 24 comprise the “Unknown set 2” motif family, with the consensus sequence 5'-[C/T][C/T]T[A/G]CA-3'. Unknown set 2 motifs are found upstream of 15 clusters maximally expressed at any of the seven time points (Figure S4 A-D). Motifs 21, 22 and 25 are unrelated, do not have overt similarity to any known *cis*-regulatory motifs, and are overrepresented upstream of 3, 5, and 12 clusters respectively. Again, clusters containing these overrepresented motifs are expressed at any of the tested time points (Figure S5-S7 A-D). Major biological process GO terms found enriched in these clusters include nucleosome assembly (motif 21), translation (motif 22), regulation of gene expression, DNA-dependent transcription initiation, post-translational protein modification and protein dephosphorylation (motif 25). Finally, Figure S8 contains all 200 of the clustered expression profiles, 40 profiles per page.

Supplementary figures

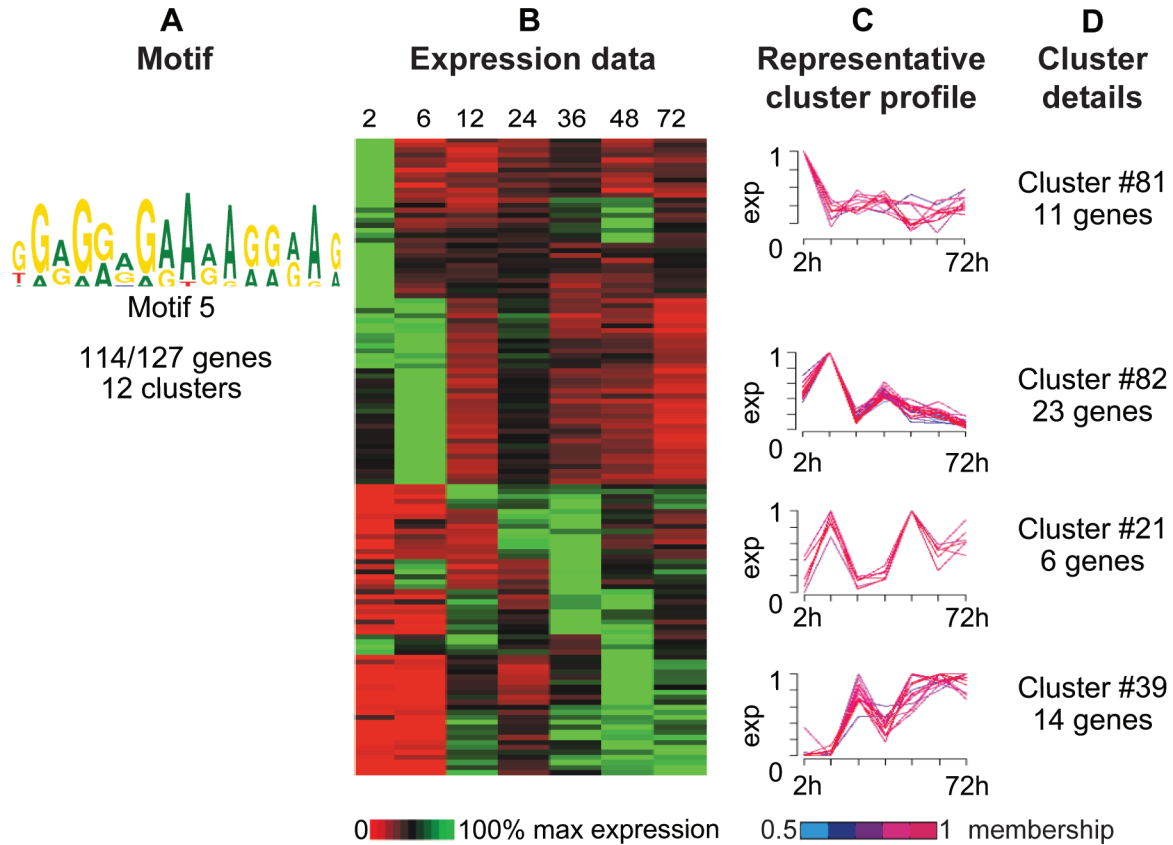


Figure S1. Data supporting identification of GAGA-like motifs. **A.** GAGA-like motifs. Motif name and total number of genes possessing each motif per total genes in all clusters where the motif is overrepresented are indicated. **B.** Expression data for 2, 6, 12, 24, 36, 48 and 72 hours post-infection for all genes from each cluster where the motif is overrepresented. Each row indicates a gene; rows are sorted first by cluster, then by peak expression at each time point. Gene IDs for genes associated with each cluster can be found in Additional file 1, Table S2. Expression is indicated on a scale of 0-100% of max for each gene. **C.** Four representative cluster profiles selected from the 12 clusters containing overrepresented GAGA-like motifs. Line colors for individual gene profiles indicate the membership values of that gene profile to the cluster ranging from 0.5 to 1. Each cluster profile is located next to the corresponding rows in the gene expression heatmap. **D.** Cluster number and total number of genes in each displayed representative cluster.

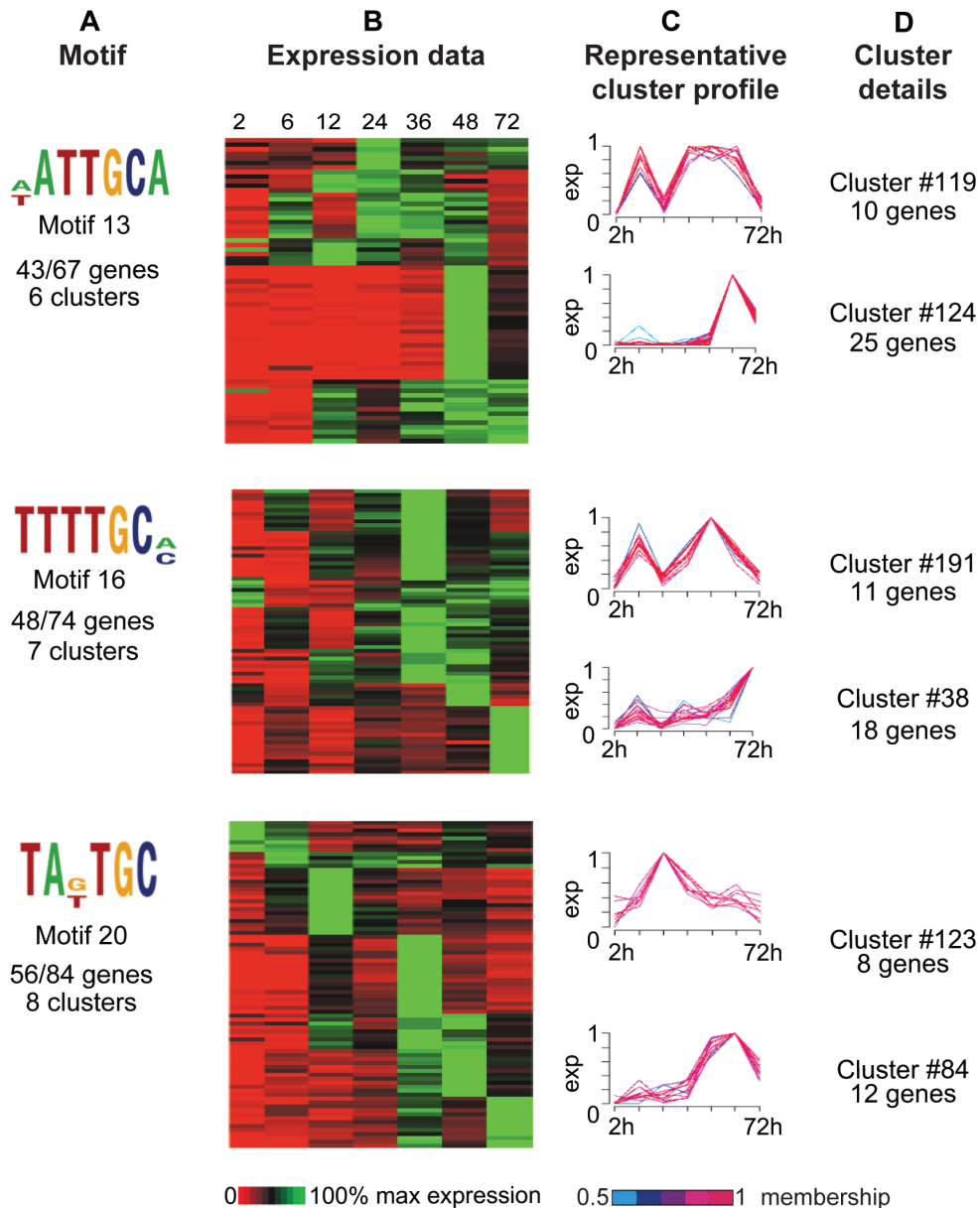


Figure S2. Data supporting identification of CAAT-box-like motifs. **A.** CAAT-box-like motifs (reverse complement). Motif name and total number of genes possessing each motif per total genes in all clusters where the motif is overrepresented are indicated. **B.** Expression data for 2, 6, 12, 24, 36, 48 and 72 hours post-infection for all genes from each cluster where the motif is overrepresented. Each row indicates a gene; rows are sorted first by cluster, then by peak expression at each time point. Gene IDs for genes associated with each cluster can be found in Additional file 1, Table S2. Expression is indicated on a scale of 0-100% of max for each gene. **C.** Six representative cluster profiles selected from the 21 clusters containing overrepresented CAAT-box-like motifs. Line colors for individual gene profiles indicate the membership values of that gene profile to the cluster ranging from 0.5 to 1. Each cluster profile is located next to the corresponding rows in the gene expression heatmap. **D.** Cluster number and total number of genes in each displayed representative cluster.

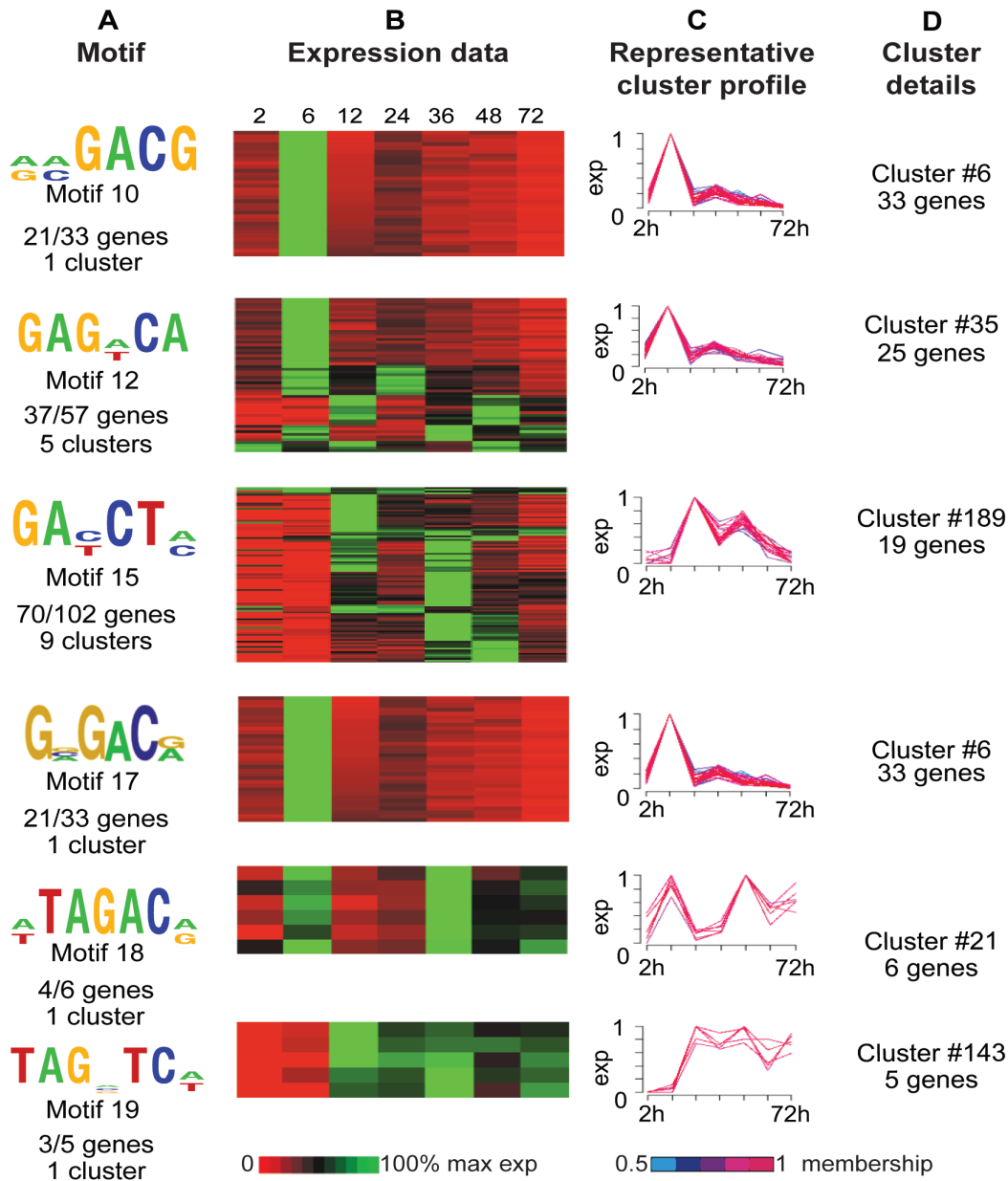


Figure S3. Data supporting identification of Unknown set 1 motifs. **A.** Unknown set 1 motifs. Motif name and total number of genes possessing each motif per total genes in all clusters where the motif is overrepresented are indicated. **B.** Expression data for 2, 6, 12, 24, 36, 48 and 72 hours post-infection for all genes from each cluster where the motif is overrepresented. Each row indicates a gene; rows are sorted first by cluster, then by peak expression at each time point. Gene IDs for genes associated with each cluster can be found in Additional file 1, Table S2. Expression is indicated on a scale of 0-100% of max for each gene. **C.** Six representative cluster profiles selected from the 18 clusters containing overrepresented Unknown set 1 motifs. Line colors for individual gene profiles indicate the membership values of that gene profile to the cluster ranging from 0.5 to 1. Each cluster profile is located next to the corresponding rows in the gene expression heatmap. **D.** Cluster number and total number of genes in each displayed representative cluster.

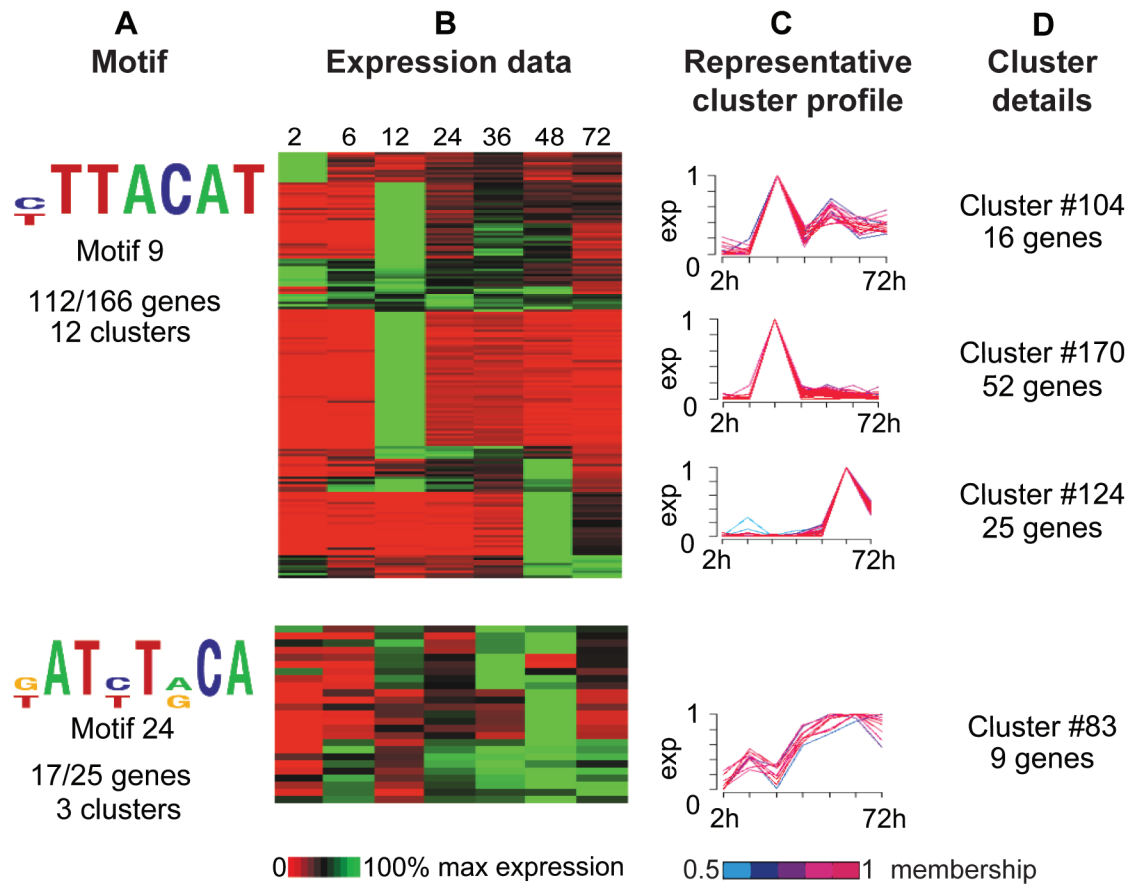


Figure S4. Data supporting identification of Unknown set 2 motifs. **A.** Unknown set 2 motifs. Motif name and total number of genes possessing each motif per total genes in all clusters where the motif is overrepresented are indicated. **B.** Expression data for 2, 6, 12, 24, 36, 48 and 72 hours post-infection for all genes from each cluster where the motif is overrepresented. Each row indicates a gene; rows are sorted first by cluster, then by peak expression at each time point. Gene IDs for genes associated with each cluster can be found in Additional file 1, Table S2. Expression is indicated on a scale of 0-100% of max for each gene. **C.** Four representative cluster profiles selected from the 15 clusters containing overrepresented Unknown set 2 motifs. Line colors for individual gene profiles indicate the membership values of that gene profile to the cluster ranging from 0.5 to 1. Each cluster profile is located next to the corresponding profile rows in the gene expression heatmap. **D.** Cluster number and total number of genes in each displayed representative cluster.

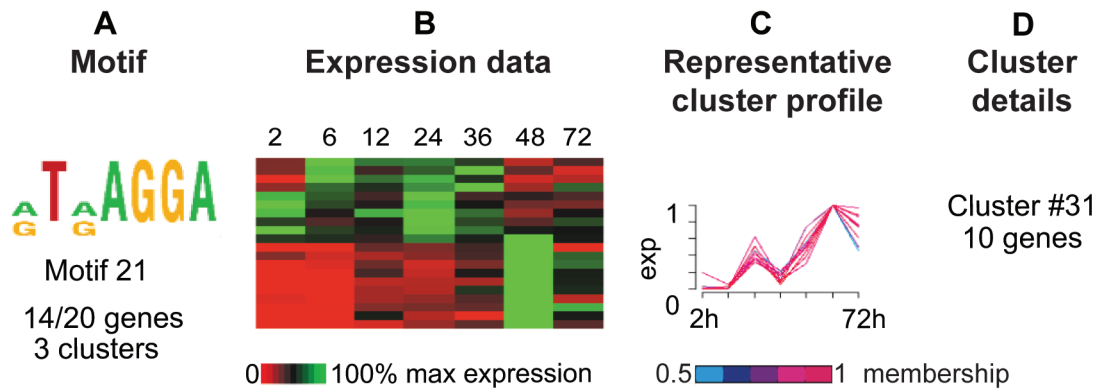


Figure S5. Data supporting identification of Unknown motif 21. **A.** Unknown motif 21. Total number of genes possessing the motif per total genes in all clusters where the motif is overrepresented is indicated. **B.** Expression data for 2, 6, 12, 24, 36, 48 and 72 hours post-infection for all genes from each cluster where the motif is overrepresented. Each row indicates a gene; rows are sorted first by cluster, then by peak expression at each time point. Gene IDs for genes associated with each cluster can be found in Additional file 1, Table S2. Expression is indicated on a scale of 0-100% of max for each gene. **C.** One representative cluster profile selected from the 3 clusters containing overrepresented Unknown motif 21. Line colors for individual gene profiles indicate the membership values of that gene profile to the cluster ranging from 0.5 to 1. Each cluster profile is located next to the corresponding rows in the gene expression heatmap. **D.** Cluster number and total number of genes in the displayed representative cluster.

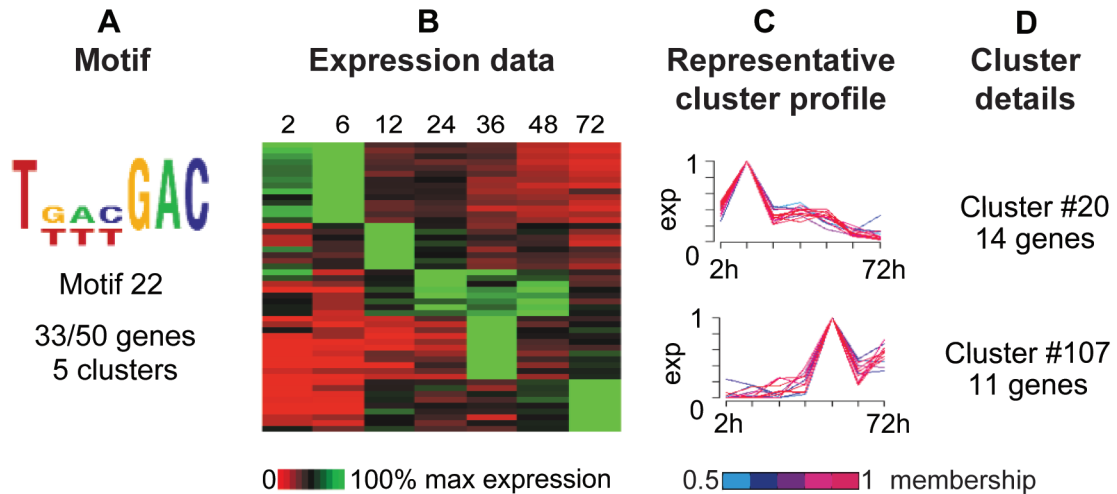


Figure S6. Data supporting identification of Unknown motif 22. **A.** Unknown motif 22. Total number of genes possessing the motif per total genes in all clusters where the motif is overrepresented is indicated. **B.** Expression data for 2, 6, 12, 24, 36, 48 and 72 hours post-infection for all genes from each cluster where the motif is overrepresented. Each row indicates a gene; rows are sorted first by cluster, then by peak expression at each time point. Gene IDs for genes associated with each cluster can be found in Additional file 1, Table S2. Expression is indicated on a scale of 0-100% of max for each gene. **C.** Two representative cluster profiles selected from the 5 clusters containing overrepresented Unknown motif 22. Line colors for individual gene profiles indicate the membership values of that gene profile to the cluster ranging from 0.5 to 1. Each cluster profile is located next to the corresponding rows in the gene expression heatmap. **D.** Cluster number and total number of genes in each displayed representative cluster.

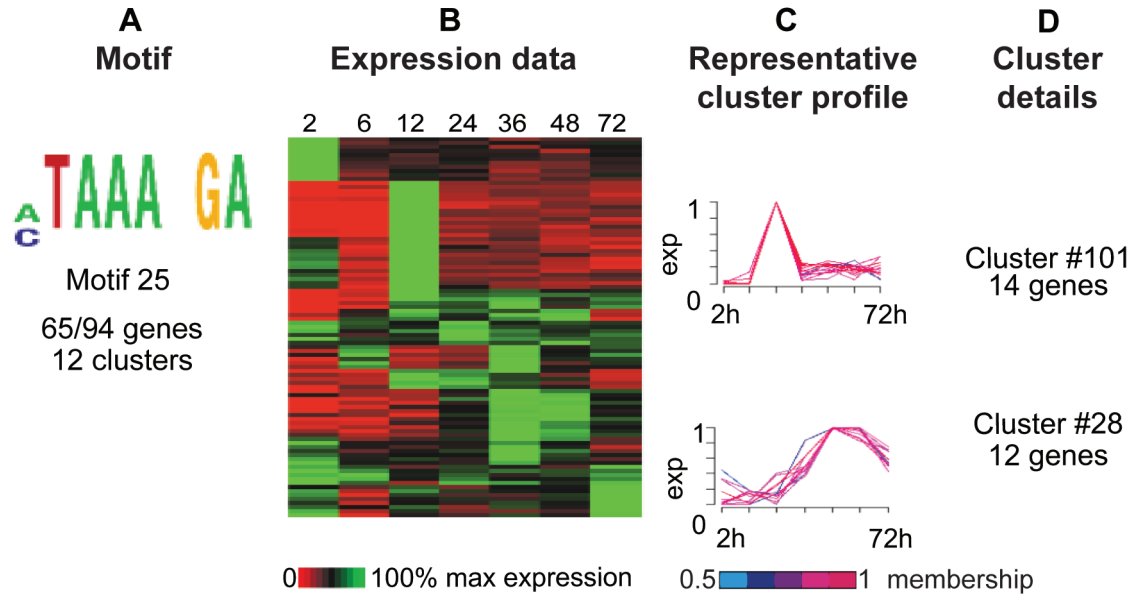
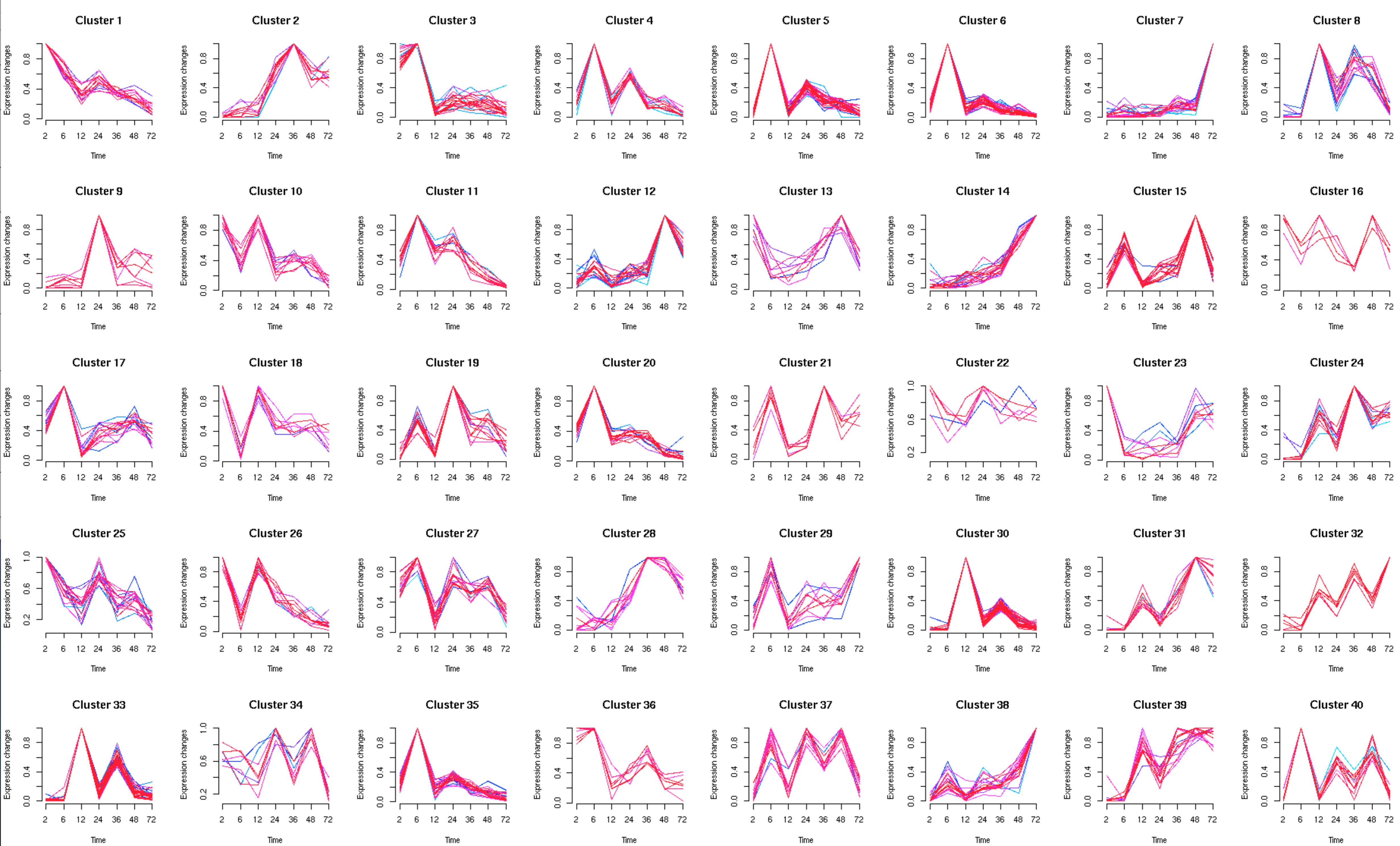
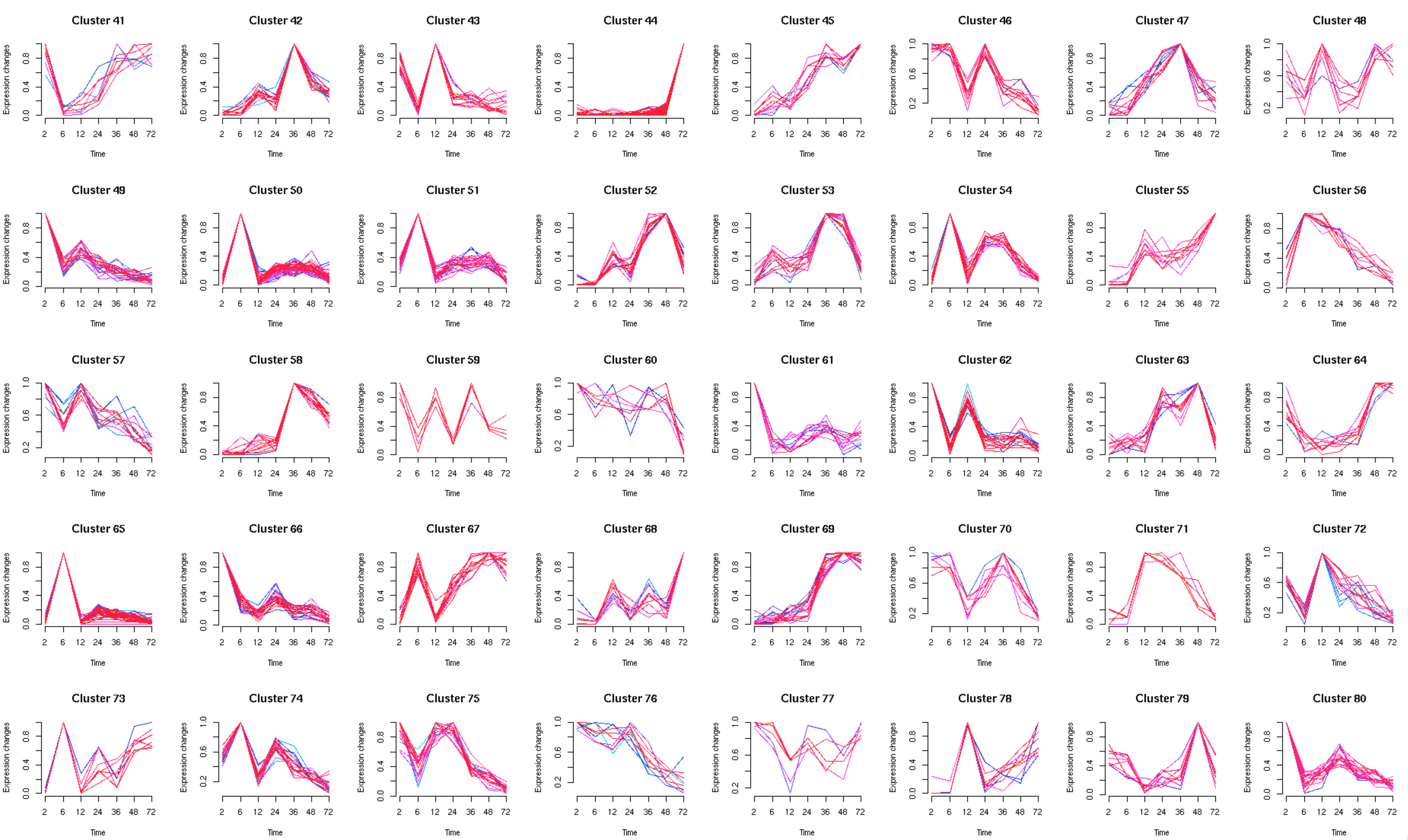
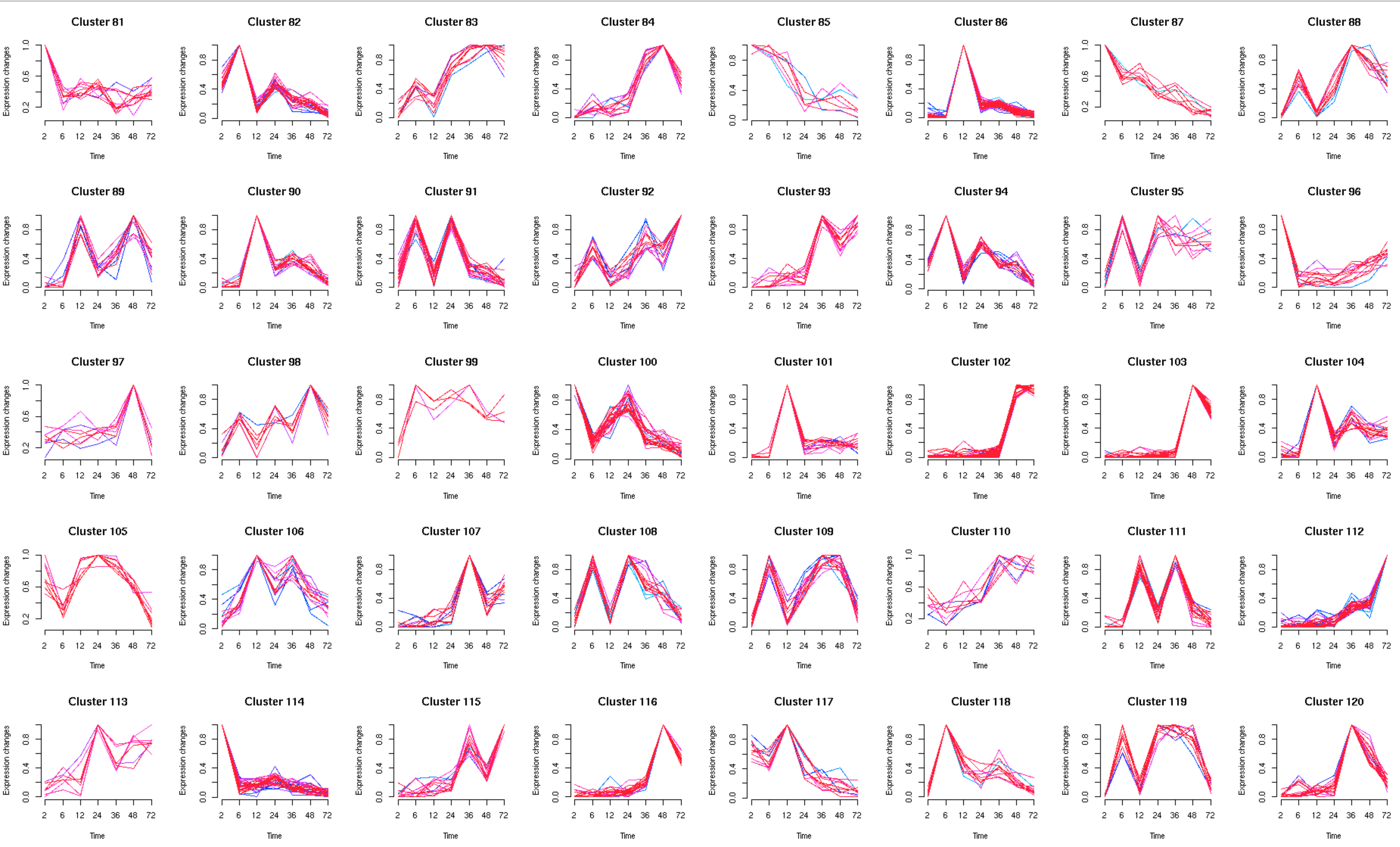


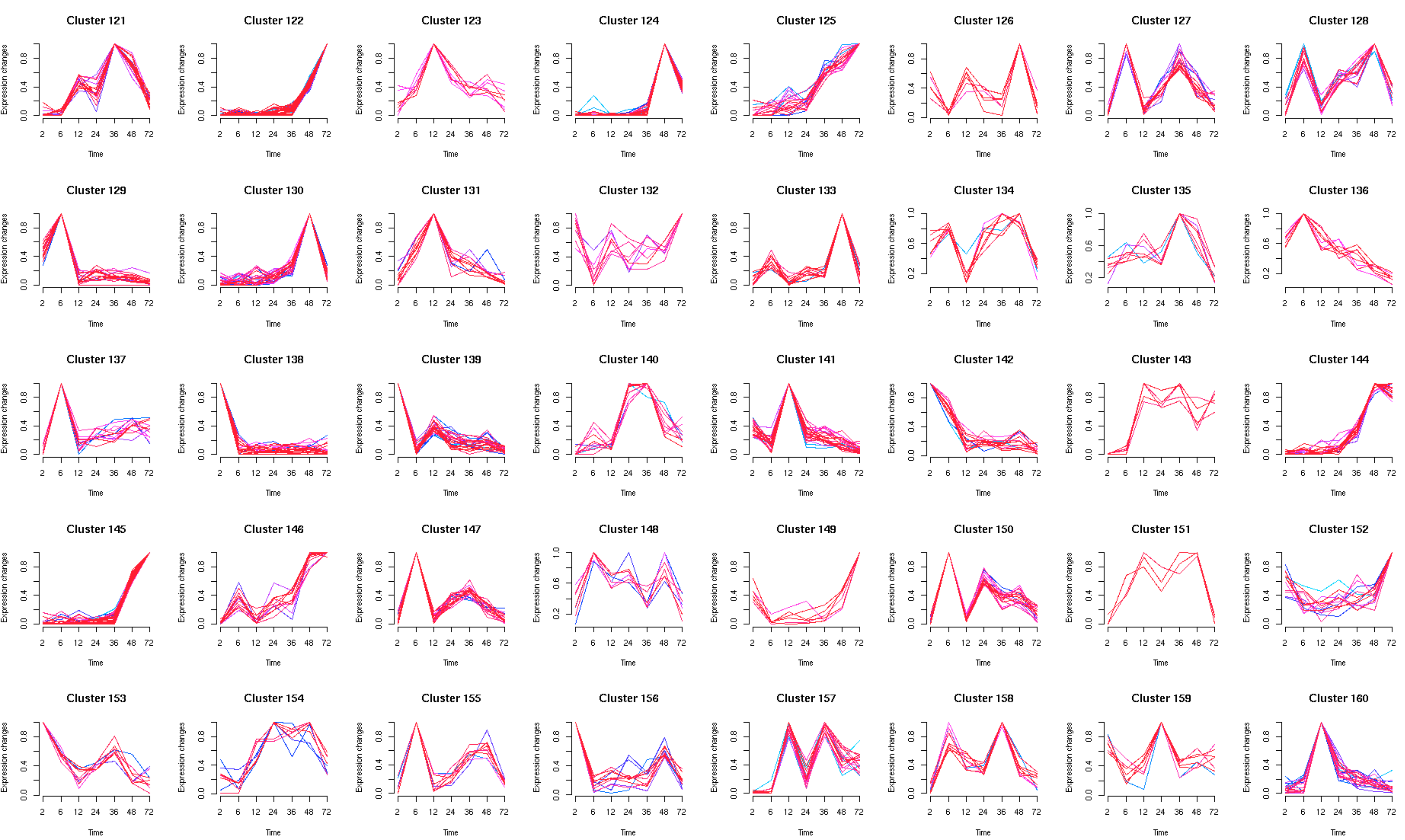
Figure S7. Data supporting identification of Unknown motif 25. **A.** Unknown motif 25. Total number of genes possessing the motif per total genes in all clusters where the motif is overrepresented is indicated. **B.** Expression data for 2, 6, 12, 24, 36, 48 and 72 hours post-infection for all genes from each cluster where the motif is overrepresented. Each row indicates a gene; rows are sorted first by cluster, then by peak expression at each time point. Gene IDs for genes associated with each cluster can be found in Additional file 1, Table S2. Expression is indicated on a scale of 0-100% of max for each gene. **C.** Two representative cluster profiles selected from the 12 clusters containing overrepresented Unknown motif 25. Line colors for individual gene profiles indicate the membership values of that gene profile to the cluster ranging from 0.5 to 1. Each cluster profile is located next to the corresponding rows in the gene expression heatmap. **D.** Cluster number and total number of genes in each displayed representative cluster.

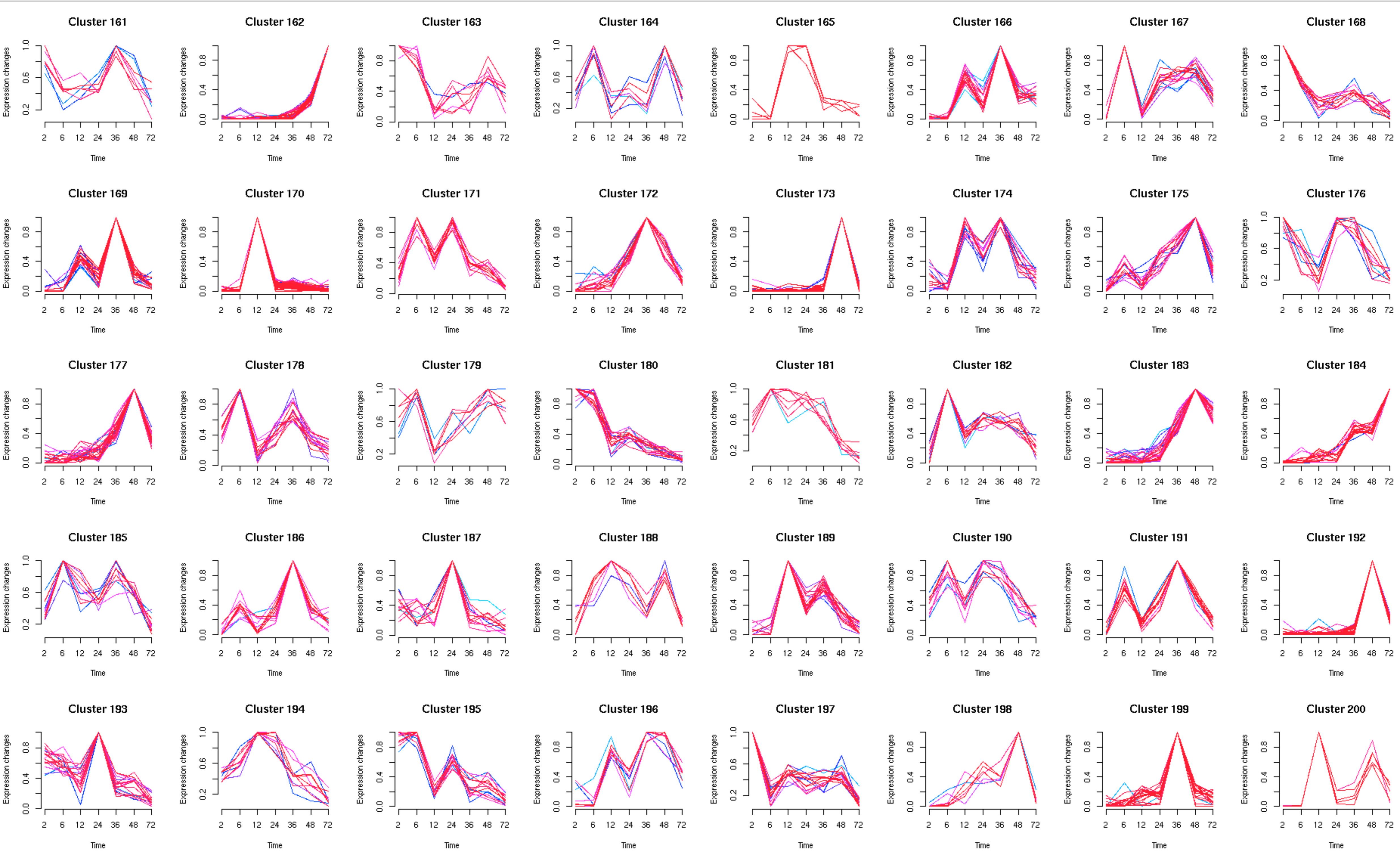
Figure S8. Clustered expression profiles, clusters 1-200. Normalized expression levels for each gene on a scale of 0 (no expression) to 1 (max expression) at 2, 6, 12, 24, 36, 48 and 72 hours post-infection for each cluster. Line colors for individual genes in the profiles indicate the membership values of that gene profile to the overall cluster. The membership values range from 0.5 to 1. Red = 0.9-1; bright pink = 0.8-0.89; purple = 0.7-0.79; dark blue = 0.6-0.69; light blue = 0.5-0.59.











References

1. Biggin MD, Tjian R: **Transcription factors that activate the Ultrabithorax promoter in developmentally staged extracts.** *Cell* 1988, **53**:699-711.
2. Mullapudi N, Joseph SJ, Kissinger JC: **Identification and functional characterization of cis-regulatory elements in the apicomplexan parasite *Toxoplasma gondii*.** *Genome Biol* 2009, **10**:R34.
3. Maity SN, de Crombrughe B: **Role of the CCAAT-binding protein CBF/NF-Y in transcription.** *Trends in Biochemical Sciences* 1998, **23**:174-178.