

Figure S1: Pseudocode of McEnhancer

```

input      : sequences of known DHSs regulating genes belonging to specific
               expression clusters
               sequences of all DHSs in a window of +/-50 kb around TSSs of
               each gene in cluster  $c$ .
output    : DHSs regulating genes in cluster  $c$  are assigned a yes label.
1 begin
2   for  $i \in \text{clusters}$  do
      initialize: set  $D_i^{lc}$  = known DHSs for cluster  $c$ 
                  set  $D_i^l$  = known DHSs for cluster  $i$  (negative/null model)
                  set  $D_i^{uc}$  = all other DHSs in +/-50 kb window around
                  TSSs for genes in cluster  $c$ .
      // E-Step
3      $\text{positiveIMM}$  = train a 3-order IMM from  $D_i^{lc}$  ;
4      $\text{negativeIMM}_i$  = train a 3-order IMM on known DHSs for cluster
       $i$  ( $D_i^l$  set);
6     repeat
7       for ( $DHS^u \in D_i^{uc}$ ) do
8         calculate loglikelihood  $DHS_{pos}^u$  from positive-IMM;
9         calculate loglikelihood  $DHS_{neg}^u$  from negative-IMM;
10        if ( $DHS_{pos}^u > DHS_{neg}^u$ ) and
11          ( $DHS_{pos}^u - DHS_{neg}^u > \text{rejectClassThreshold}$ ) :
12          | add DHS to  $D_i^{lc}$  ;
13          | remove DHS from  $D_i^{uc}$  ;
14        end
      // M-Step
15      $\text{positiveIMM}$  = update positive-IMM parameters based on
      updated  $D_i^{lc}$  set ;
16     for ( $DHS^l \in D_i^l$ ) do
17       calculate loglikelihood  $DHS_{pos}^l$  from positive-IMM;
18       calculate loglikelihood  $DHS_{neg}^l$  from negative-IMM;
19       if ( $DHS_{pos}^l > DHS_{neg}^l$ ) :
20       | keep DHS in  $D_i^{lc}$  ;
21       else:
22       | remove DHS from  $D_i^{lc}$  ;
23       | add DHS to  $D_i^{uc}$  ;
24     end
25     until  $D_i^{uc}$  does not change or maximum number of iterations
      exceeded;
26   end
27    $\text{getDHSsSelectedAtLeast10Times}(D_i^{lc})$ 
28 end

```

Figure S2: Number of labeled DHS-gene pairs. No. of known distal DHS-gene pairs overlapping REDfly and VT, which are used for model initialization.

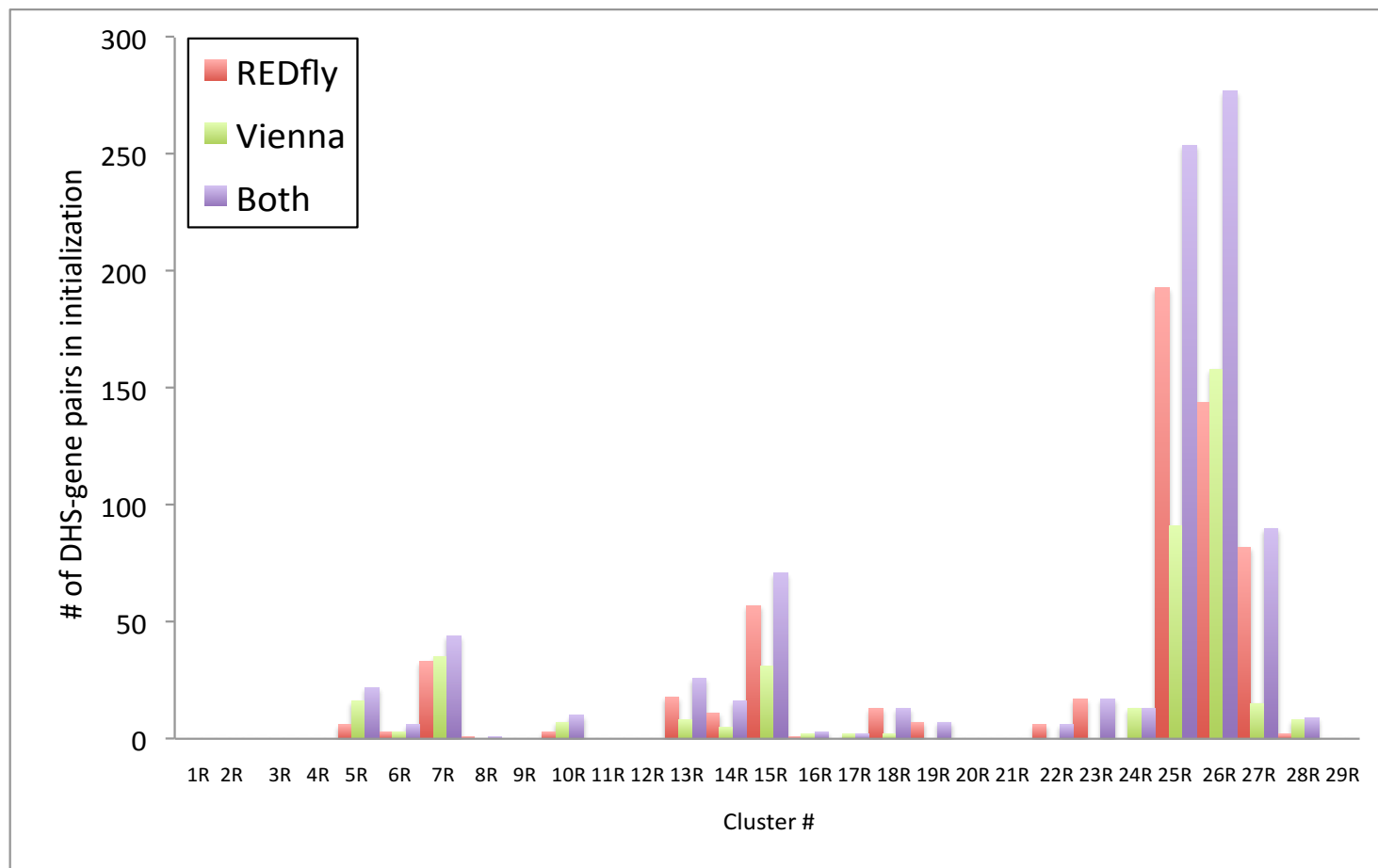


Figure S3: No. of genes that are within 100kb of each other per cluster

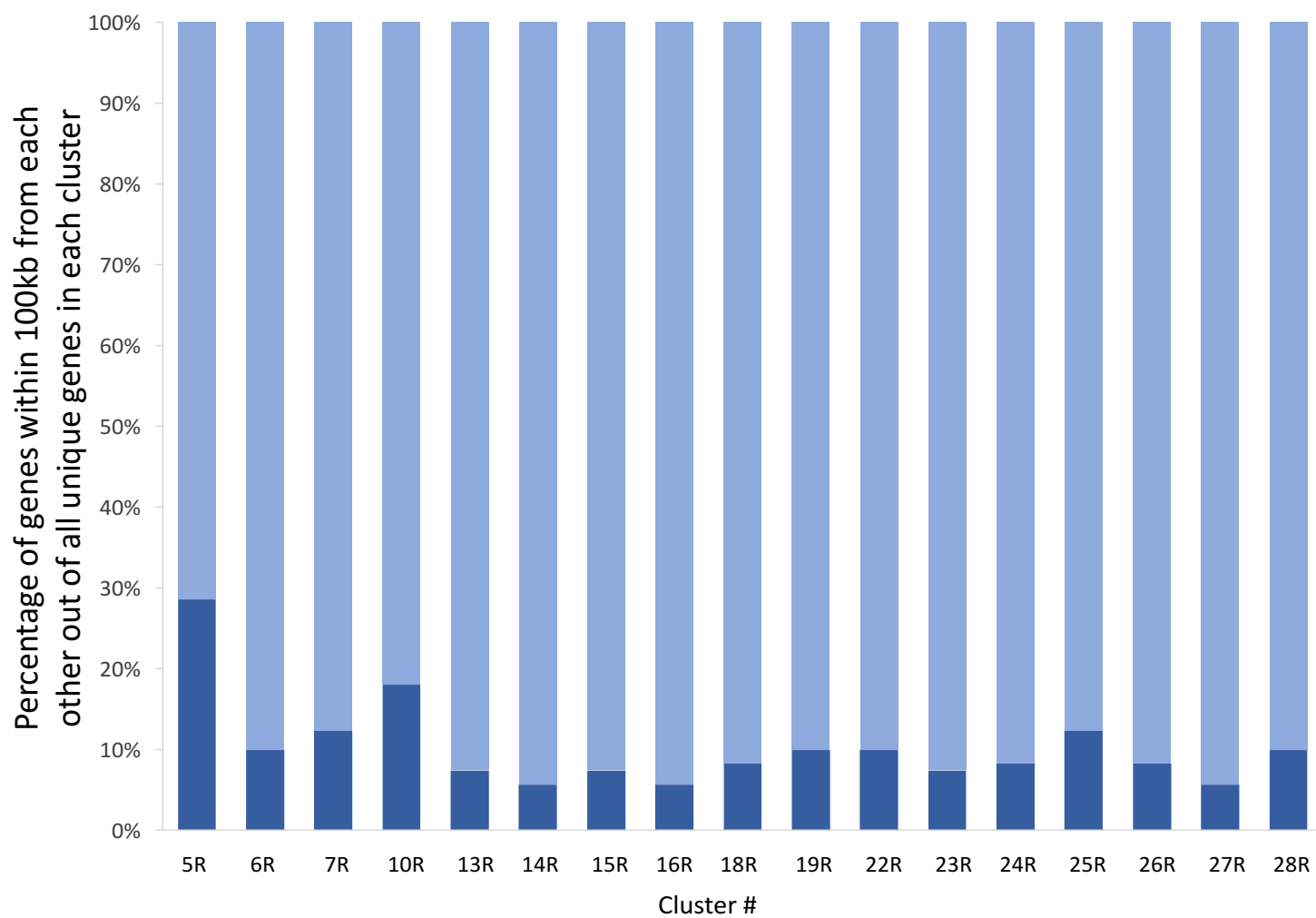


Figure S4: Computational validation of McEnhancer predictions by initialization with REDfly exclusively. (A) Overlap of selected DHSs when initializing McEnhancer by REDfly only, with the full set of VT DHSs assigned to each cluster. (B) Distribution of enhancer-gene distances in the VT set, and of the recovered subset when initializing on REDfly only.

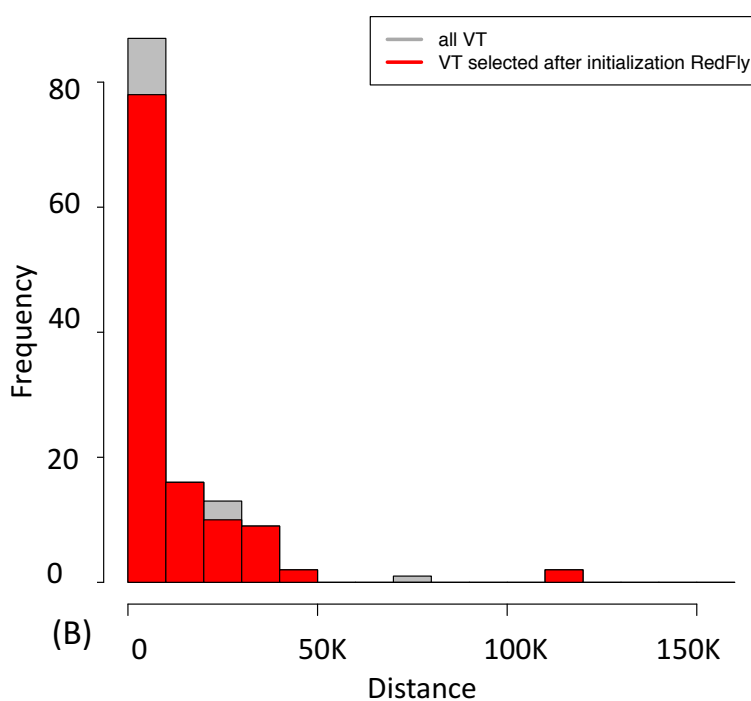
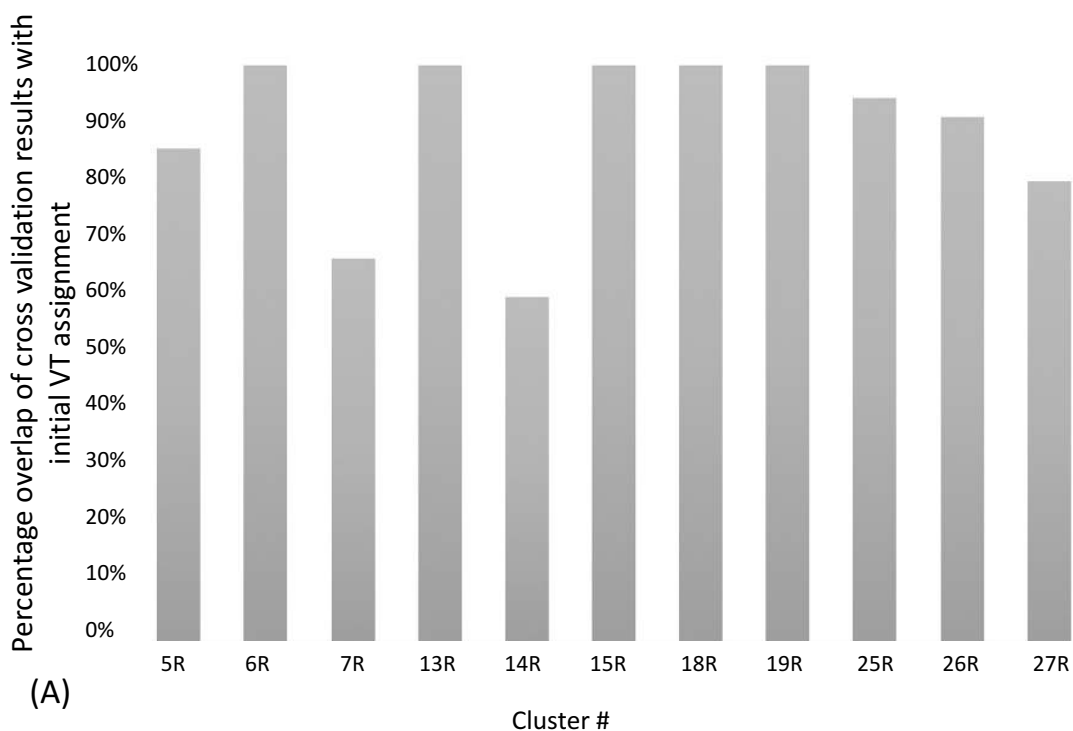


Figure S5: Average DHS accessibility for McEnhancer predictions and overlap with significant Hi-C fragments. (A) The average DHS accessibility for the McEnhancer selected DHSs and other control groups. DHS peak signal was averaged across each of the four groups and centered over DHS midpoint. (B) Validation of predicted DHS-gene pairs by overlapping with significantly identified fragment-promoter Hi-C interactions. We used the default setting of HIPPIE for setting significance threshold for interaction to be P-value < 0.0001. Percentages represented by this overlap with respect to filtered predicted DHS-gene pairs are shown by the dotted line.

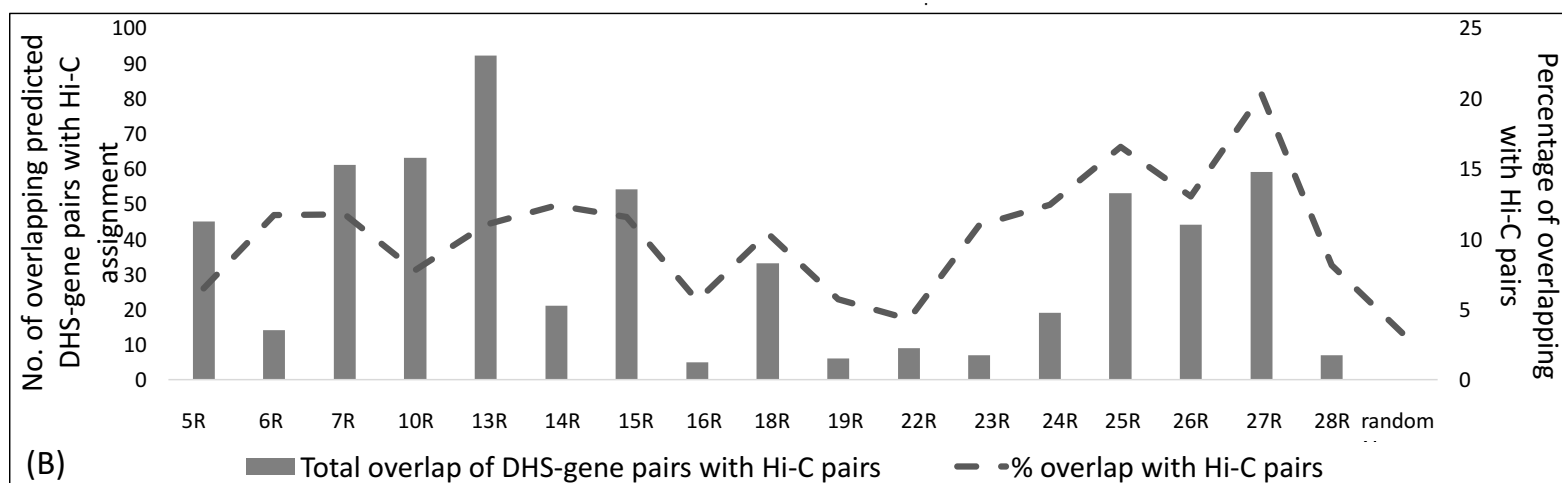
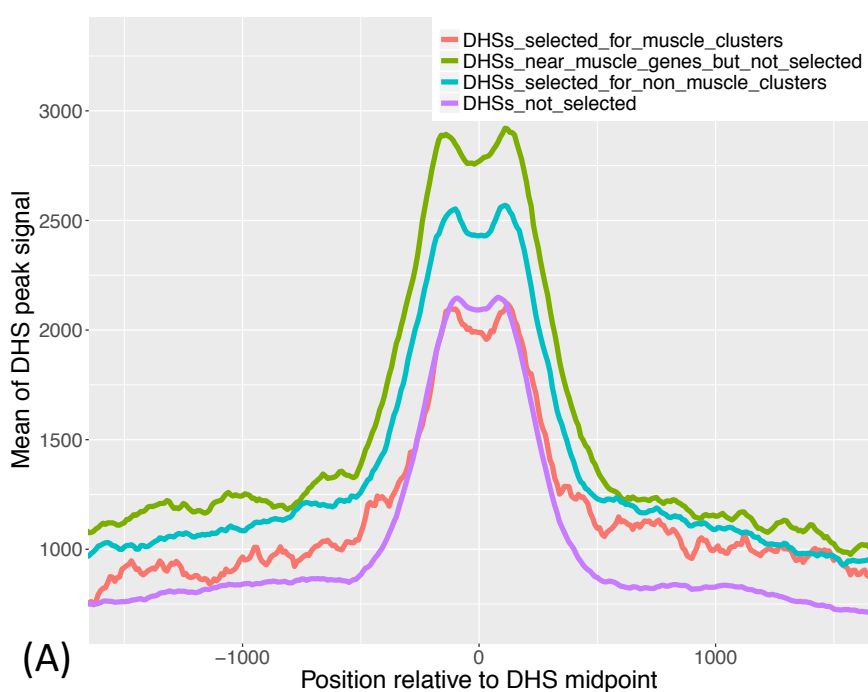


Figure S6: TFs Heatmap showing enriched TFs in each gene cluster.

Using AME motif enrichment algorithm, enriched TFs based on reported binding preferences were identified, and their corresponding Pvalue score are shown on the heatmap (upper panel). Overlapping 5-mers and their matched TFs are shown on the bottom panel, with a black squares means “match” and a white one “no match”.

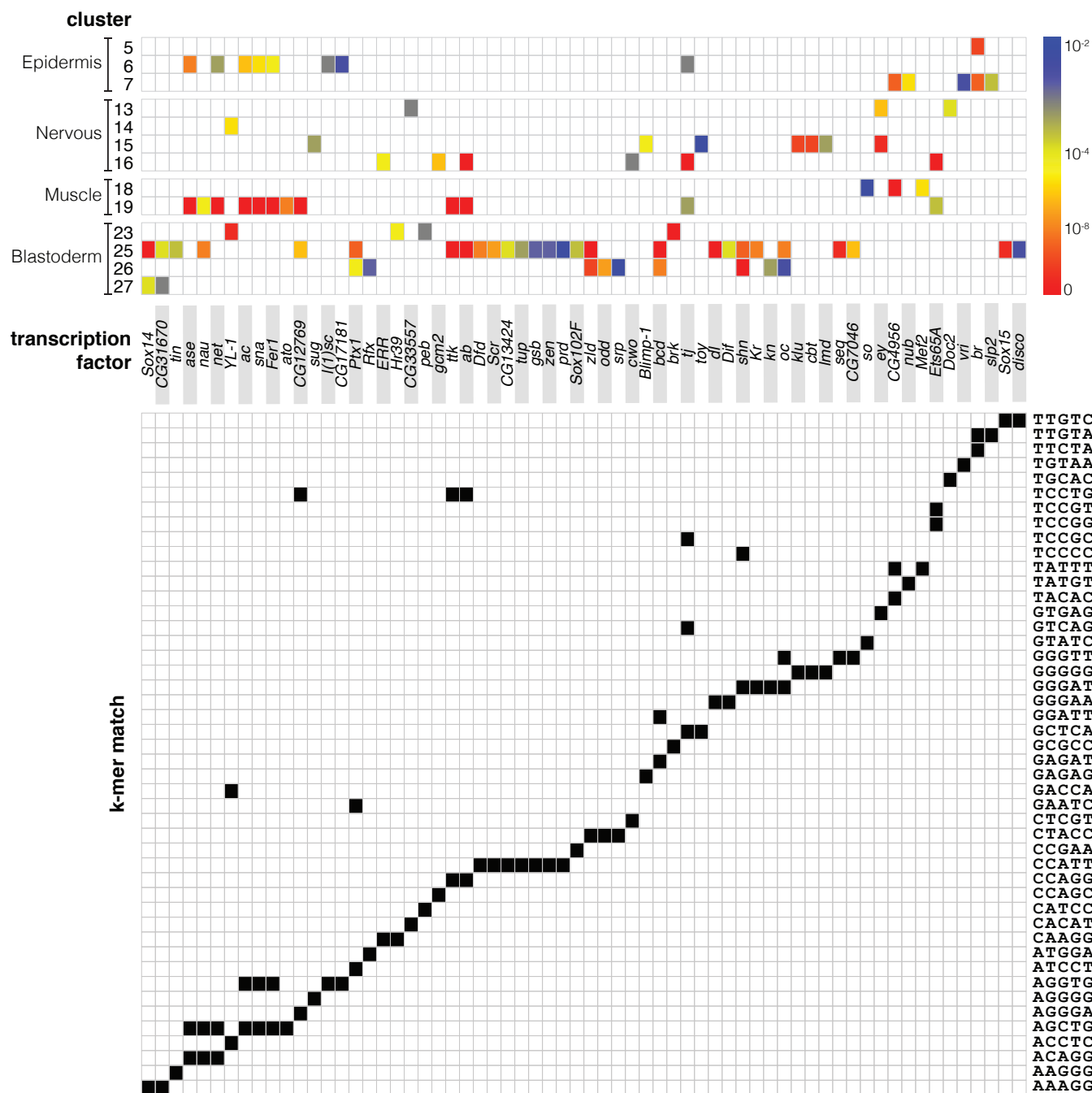


Figure S7: Experimental validation for selected candidates. Figure structured as Figure 8. Genomic and in-vivo reporter data is shown for CRM6053 (A), CRM5481 (B), CRM3775 (C), CRM4515 (D), and candidate 3 (E). Embryo developmental stages are indicated, orientation is anterior left. All embryos are lateral views ventral down, except B'-st.15, which is a dorsal view and D'-st.9, which is a ventral view.

