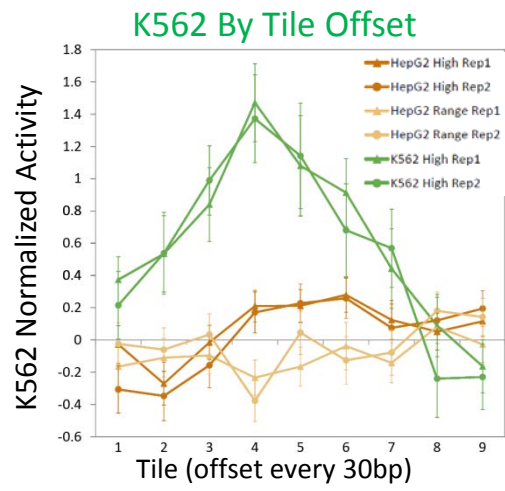
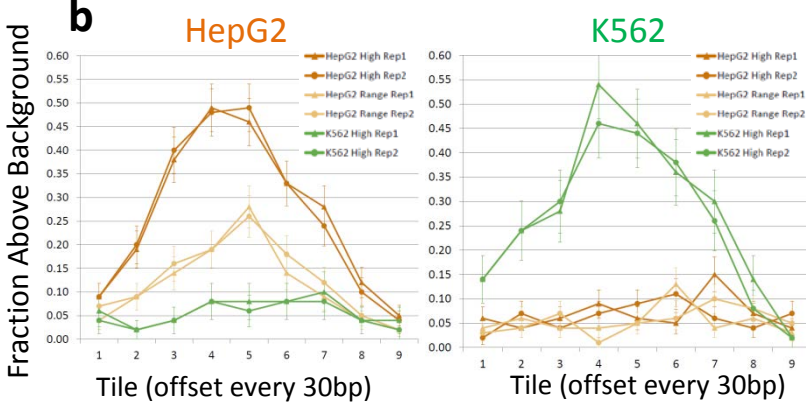
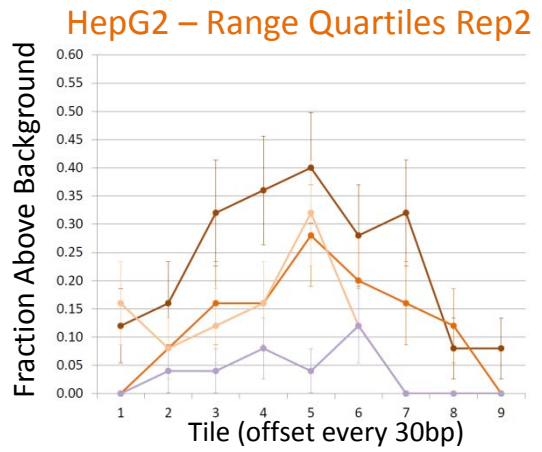
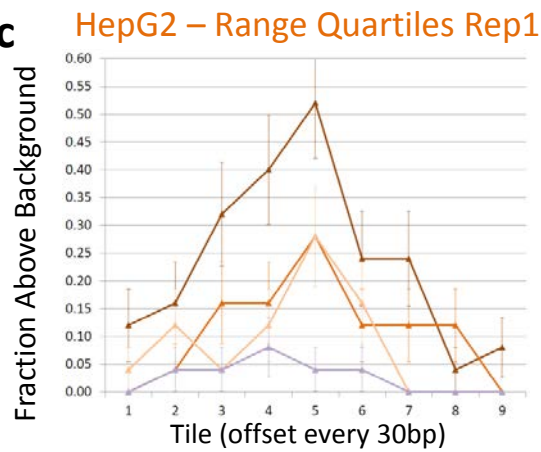
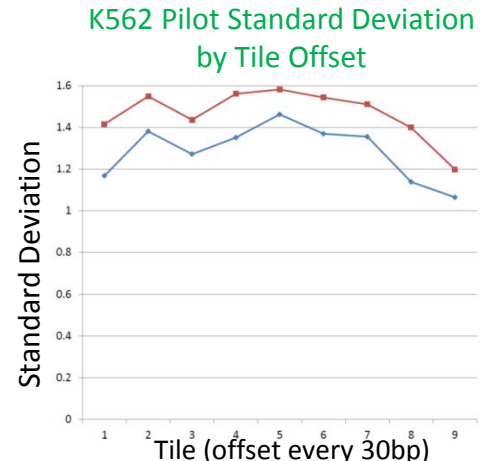
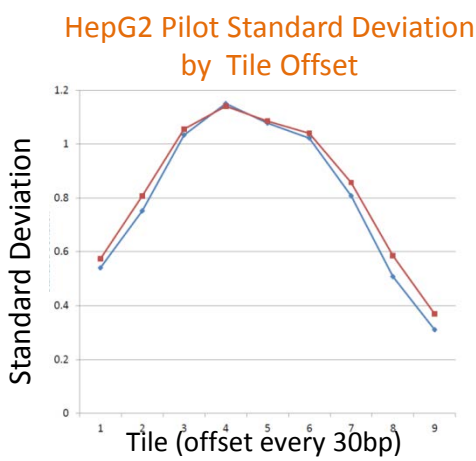


**Supplementary Figure 1 – Large-step Pilot Design** **(a)** Overview of the selection of regions and experiments. Two hundred fifty regulatory regions were selected to be tested with 100 being selected based on being in an HepG2 candidate enhancer state and having a high H3K27ac dip score, 100 being selected based on being in an HepG2 candidate enhancer state and covering a range of H3K27ac dip scores, and 50 for being in a K562 candidate enhancer state with a high H3K27ac dip score in K562 and in a low-activity state in HepG2 (see Methods). These regulatory regions were tested in both K562 and HepG2 using a SV40 promoter in replicate. **(b)** Chromatin state model used for selecting strong enhancer regions in HepG2 and K562, and low-activity regions in K562<sup>2</sup>. The frequency of each chromatin state mark is shown (scaled to be between 0 and 100). Note that for the scale-up model, we used a richer 25-state model to capture a higher diversity of chromatin states (shown in Supplementary Fig. 5). **(c)** The first two larger columns are the H3K27ac signal<sup>2</sup> in HepG2 and K562 cells respectively, while the next two larger columns are the DNase signal<sup>7</sup> in HepG2 and K562 cells respectively. Each larger column is separated by the white lines. Within each larger column is a heatmap of the signal across the 385-bp region based on the color scale shown at bottom. The first larger set of rows (separated by horizontal white lines) are those regions corresponding to the HepG2 tiled enhancer regions with a high dip score, the second are those regions corresponding to the HepG2 tiled enhancer regions with a range of dip scores, and the third set are the K562 based tiled enhancer regions. Within each set the regions are ordered in terms of decreasing dip score. **(d)** Heatmap of the reporter expression values. The first two larger columns show the results for the experiments in two HepG2 replicates and the next two columns show the experiments in the two K562 replicates. Within each larger column is the median log-base two ratio reporter expression values at each of the nine tile offsets normalized by taking the difference with the average value for the expression of the tiles in the outmost tile offsets (tiles #1 and #9) colored based on the indicated color scale shown at the bottom.

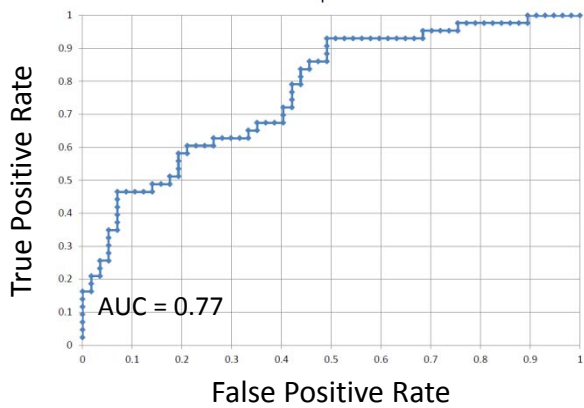
**a****b****c****d**

**Supplementary Figure 2 – Large-step Pilot Design Analysis.** (a) The same plot on average reporter expression by tile offset and group of tiled regions as in Fig. 2a but for the experiments conducted in K562 cells. (b) The fraction of reporter values that met the threshold at which only 5% of outmost tile offsets (tiles #1 and #9) reporter values did for (left) HepG2 and (right) K562 cells experiments. (c) The same as (b) except only showing the range values broken down into four quartiles and showing in separate graphs replicate 1 (left) and replicate 2 (right). (d) Standard deviation (y-axis) of the reporter expression values at each tile offset (x-axis) for HepG2 experiments (left) and K562 experiments (right) for each replicate (color).

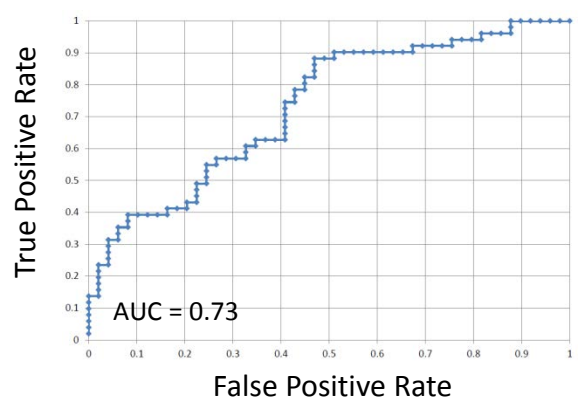
[Supplementary Figure 2 panel e is continued on next page]

**e**

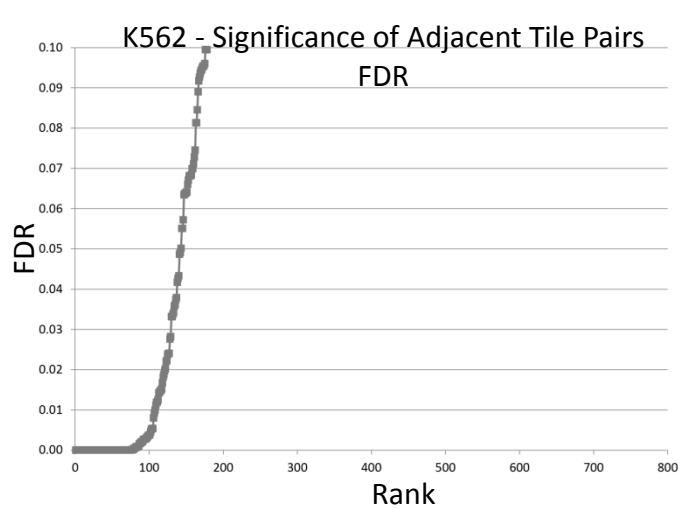
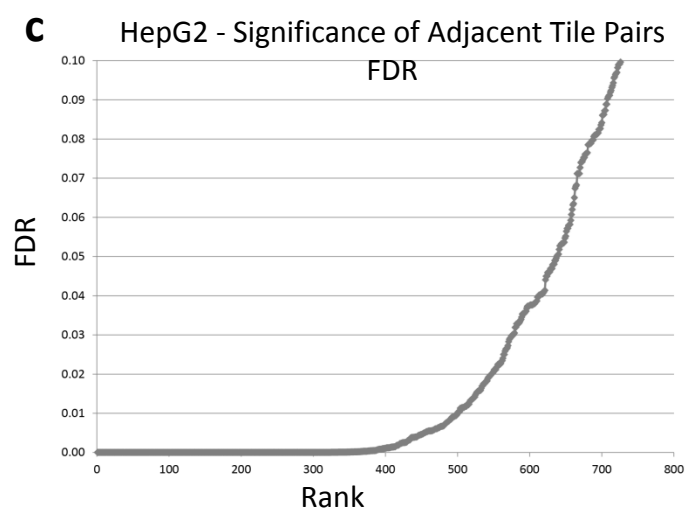
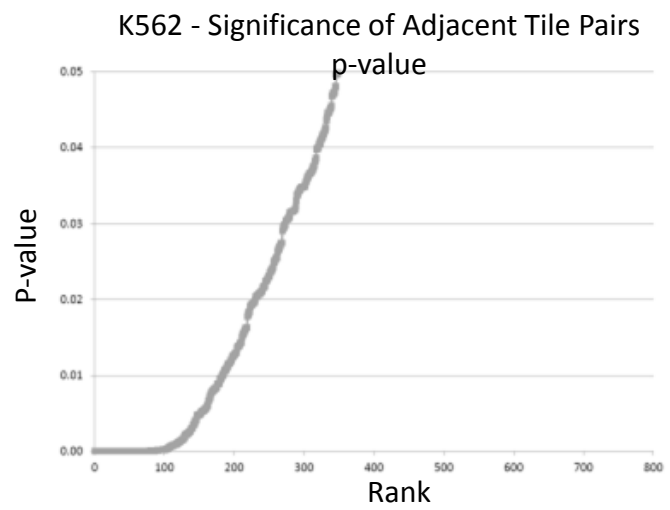
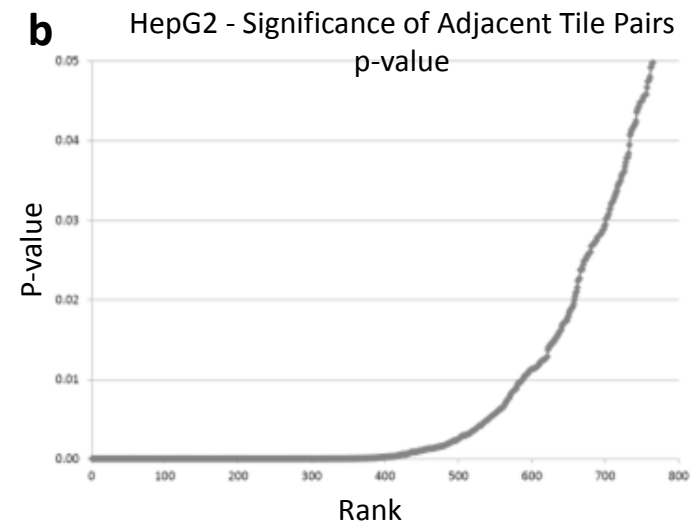
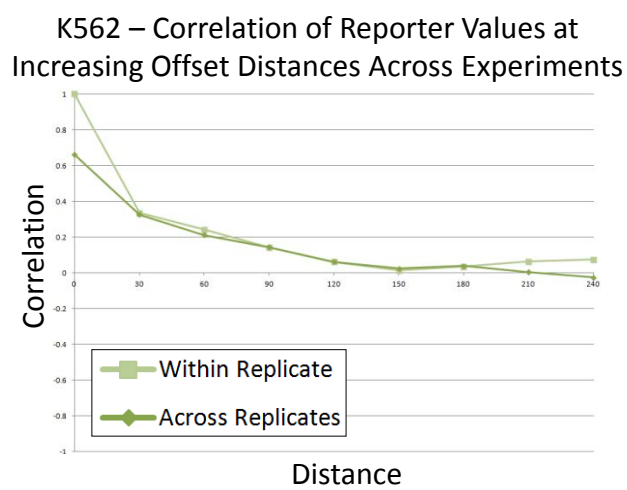
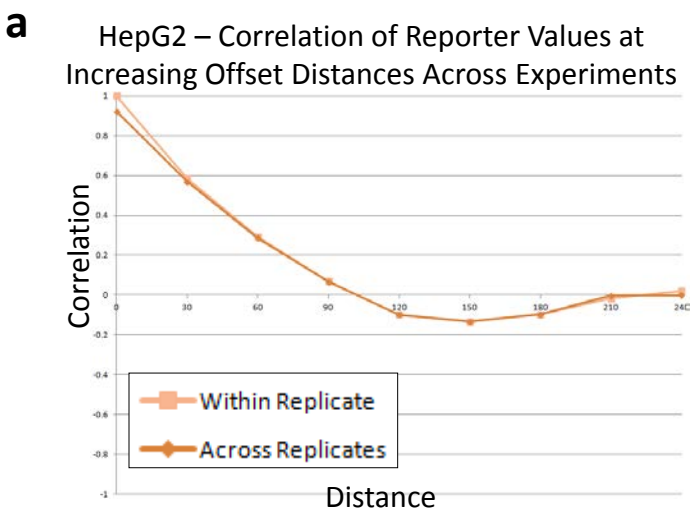
HepG2 – Range Rep1 Dip Score  
Predicting Above Background Regions



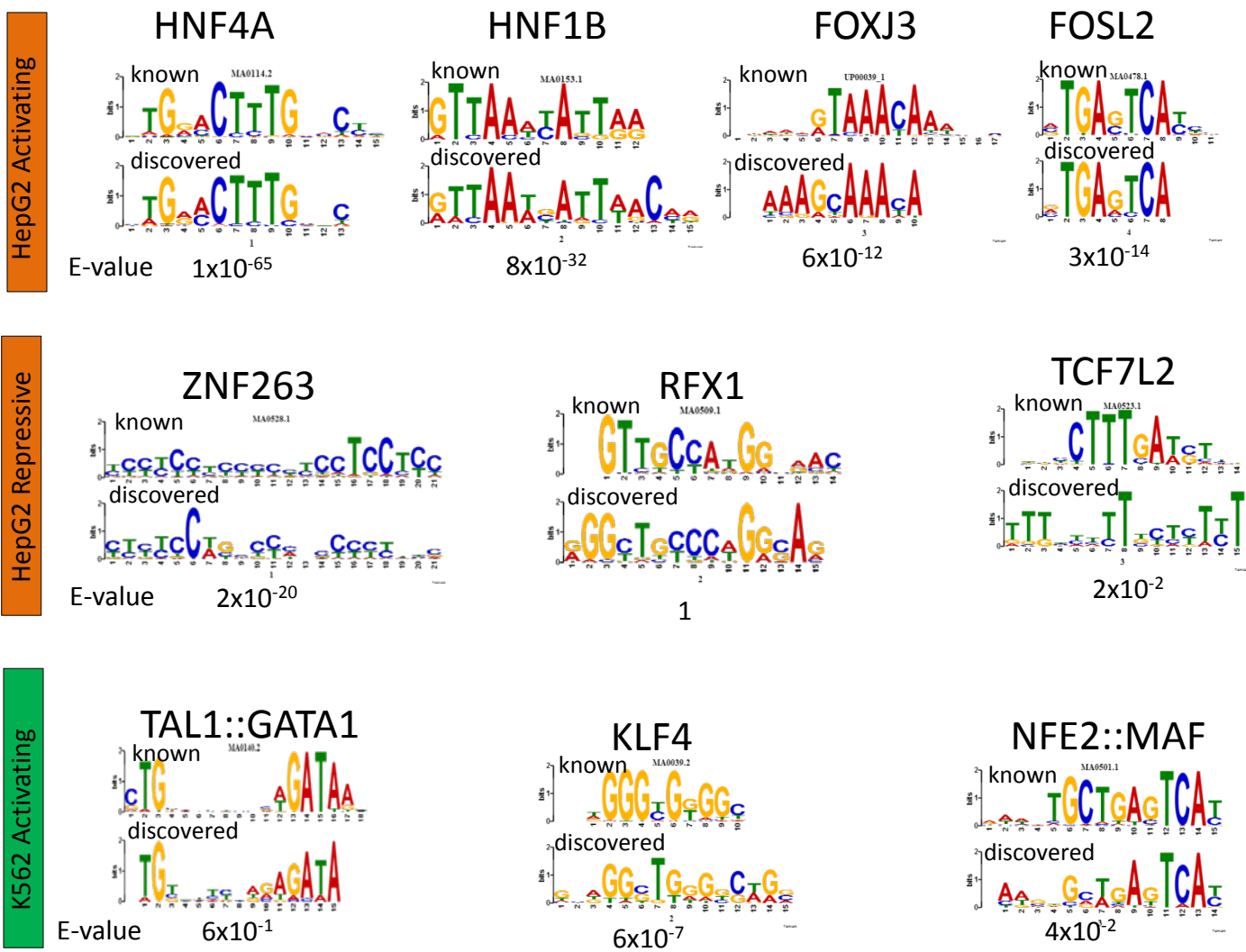
HepG2 – Range Rep2 Dip Score  
Predicting Above Background Regions



**Supplementary Figure 2 continued from previous page: (e)** For the 100 HepG2 tiled regions selected to span a range of H3K27ac dip scores, ROC curve for predicting regions with at least one tile above threshold, when ranking based on the dip score for the region for replicate 1 (left) and replicate 2 (right) experiments.



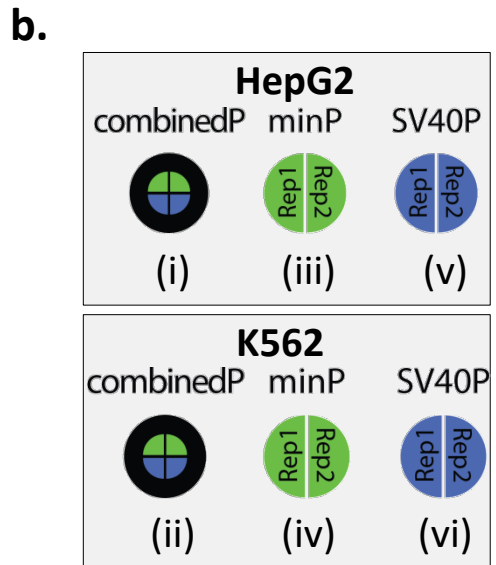
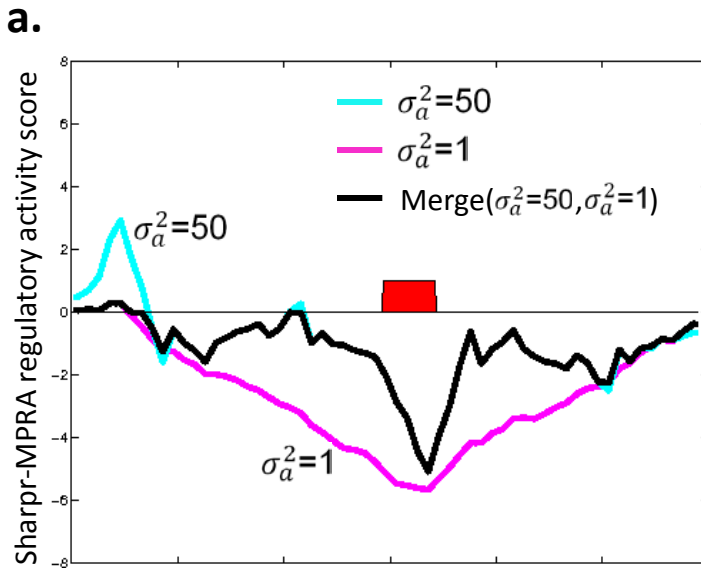
**Supplementary Figure 3 – Large-step Pilot Correlation of Reporter Values as a Function of Distance.** (a) Correlation of reporter expression values (y-axis) for each indicated offset (x-axis) both within (boxes) and across replicates (diamonds) in HepG2 (left) and K562 (right). (b) Cumulative number of adjacent Large-step pairs of tiles (x-axis) that showed a difference in expression at a given uncorrected p-value (y-axis) based on a Mann-Whitney Test (see Methods) in HepG2 (left) experiments and K562 (right) experiments. (c) The FDR values corresponding to the p-values shown in b.

**a****b**

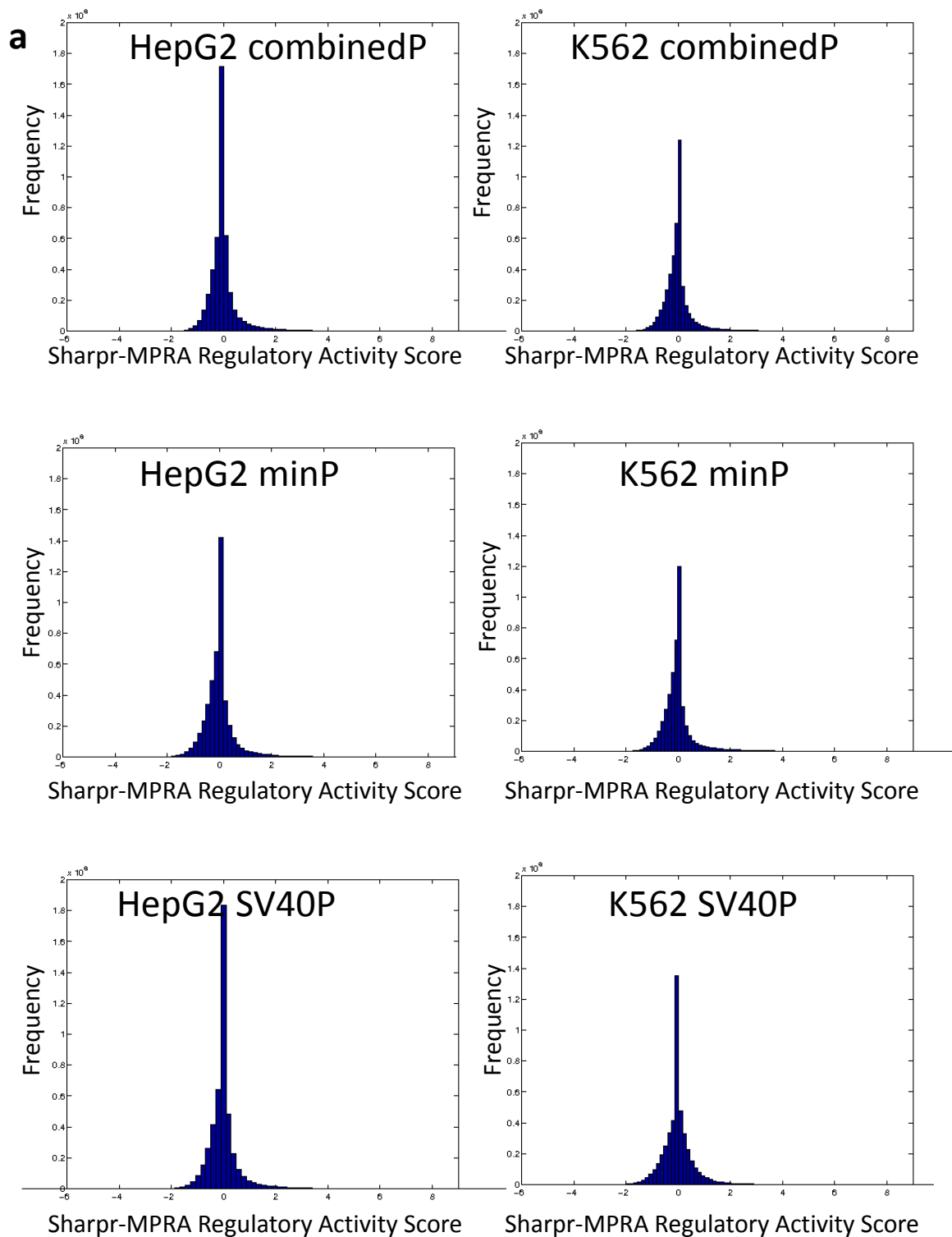
	HepG2 activating	HepG2 repressive	K562 activating	K562 repressive
GATA_known14	1.0	1.0	6.4	1.0
LMO2_2	1.0	1.0	4.4	1.0
TAL1_disc1	1.0	1.0	3.0	1.0
CPHX_1	1.0	1.0	2.7	1.0
JDP2_2	1.0	1.0	2.4	1.0
NFE2L2_3	1.0	1.0	2.0	1.0
HNF4_known9	4.6	0.3	1.0	1.0
NR2F6_2	3.7	1.0	1.0	1.0
RXRA_known10	3.4	0.9	1.0	1.0
PPARA_4	3.1	1.0	1.0	1.0
HNF1A_4	2.8	1.0	1.0	1.0
HNF1_4	2.6	1.0	1.0	1.0
TCF7L2_disc1	2.4	1.0	1.0	1.0
AP1_known4	2.3	1.0	1.6	1.0
SMARCDisc1	2.2	1.0	1.1	1.0
CEBPB_known1	2.1	1.0	1.0	1.0
TLX2_2	2.0	1.0	1.0	1.0
FOXJ2_4	2.0	1.0	1.0	1.0
EGR1_disc2	2.0	1.0	1.0	1.0
RFX2_3	1.0	3.0	1.0	1.0
RFX5_known9	1.0	2.9	1.0	1.0

**Supplementary Figure 4 – Large-step Pilot Design Motifs.** (a) Motifs discovered based on running MEME software<sup>72</sup> on the identified at a 5% false discovery rate (top) HepG2 activating, (center) HepG2 repressive, and (bottom) K562 activating 30-bp sequences extended 10-bp to include the common sequence (see Methods). Only up to the top four discovered motifs that had a reported E-value of 1 or less are shown. Motifs are ordered left to right as returned by the MEME program. There were no motifs discovered for the K562 negative set that met the criterion. Each discovered motif is matched to a best matching known motif using TOMTOM<sup>73</sup> that is shown above the discovered motif. (b) Motifs (rows) enriched in at least one of the four sets (HepG2 Activating, HepG2 Repressive, K562 Activating, and K562 Repressive) at least 2 fold based on the program of Ref<sup>13</sup> extended to use a background of all nucleotides tested. The columns show the enrichment of the motif in the four sets considered using the same program and background. Only motifs with at least three instances in the background are included and if multiple motifs corresponding to the same transcription factor met the threshold only the one with the highest enrichment is listed.



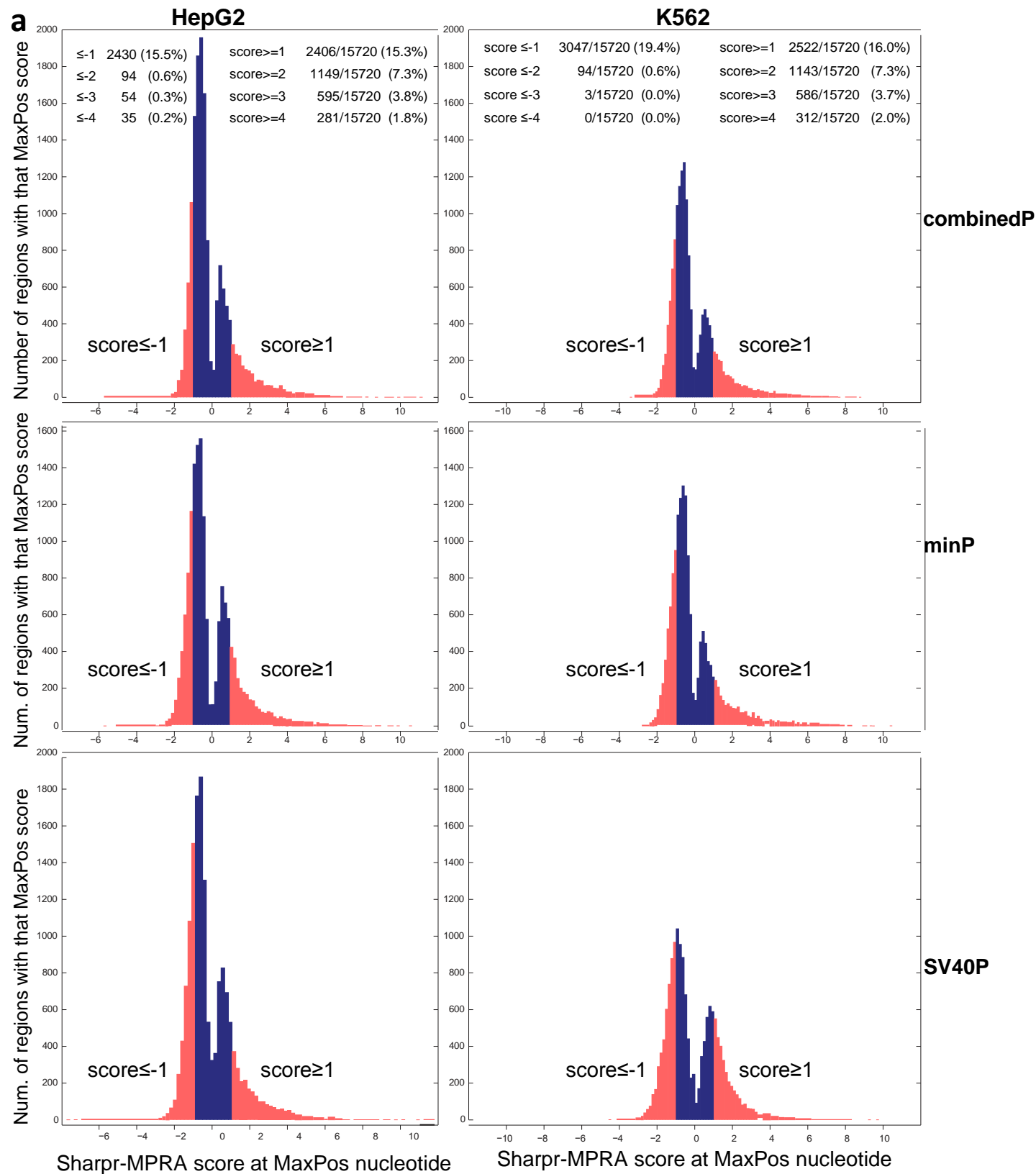


**Supplementary Figure 6 – Inference using multiple variance priors, and for multiple replicates, promoters, and cell types. (a) Example Inference with Two Different Variance Parameter Priors.** Effect of the variance prior on the output in the example in **Fig. 3b** right. If the variance prior is larger (turquoise line) the output can capture finer features, but can potentially overfit the data as seen by the activation peak on left unsupported by a transcription factor binding site prediction based on CENTIPEDE<sup>8</sup> (shown in red) in close proximity to the repressive peak. On the other hand if a smaller variance prior is used (pink line), then the output is more consistently activating or repressive, but can underfit providing a lower resolution output that can call additional nucleotides activating or repressive that are not. The strategy used was to apply the inference with both a relatively smaller variance prior ( $\sigma_a^2=1$ ) and a relatively larger variance prior ( $\sigma_a^2=50$ ) and merge the output of two inferences at each position, using the one with smaller absolute value between the two when their signs agreed, and using 0 when their signs disagreed. **(b) Summary of Sharpr-MPRA regulatory activity scores for each region.** The resulting track of applying the two variance parameters is applied in six settings for each regulatory region, resulting in six tracks of Sharpr-MPRA scores, each of 4.6 million nucleotides, which we release as genome browser tracks. The six released tracks correspond to: (i,ii) one combined promoter (combinedP) track for each cell type (black), incorporating a total of four experiments per region, including both minP replicates and both SV40P replicates; (iii,iv) one track for the minimal promoter (minP) in each cell type (green), combining the two replicates for the minP promoter; and (v,vi) one track for the strong promoter (SV40P) in each cell type (blue), combining the two replicates for the SV40P promoter. In addition to these six tracks, we ran SHARPR on each replicate separately for each cell type, and used the resulting inferences to evaluate the reproducibility of our inferences between individual replicates (**Fig. S10a**). For a subset of 44 genomic regions, two different barcode tilings overlaps perfectly, for all 295 nucleotides within them. For an additional 212 regions, multiple barcode tilings partially overlap, resulting in 256 regions with exact or partial overlap, of which 245 have two overlapping barcode sets, and 9 have three overlapping barcode sets. In forming the released browser tracks, scores for a base were averaged when the base was overlapped multiple times.



**Supplementary Figure 7a – Sharpr-MPRA regulatory activity score distribution. (a) Nucleotide-level score distribution for 4.6 million nucleotide positions.** The distribution of regulatory activity score at the nucleotide level for the combinedP promoter score (top row), minP score (middle row), and SV40P score (bottom row) in HepG2 cells (left) and K562 cells (right). The HepG2 distribution (top left) was also shown in gray in **Fig. 3c**. **(b,c)** Next two pages.



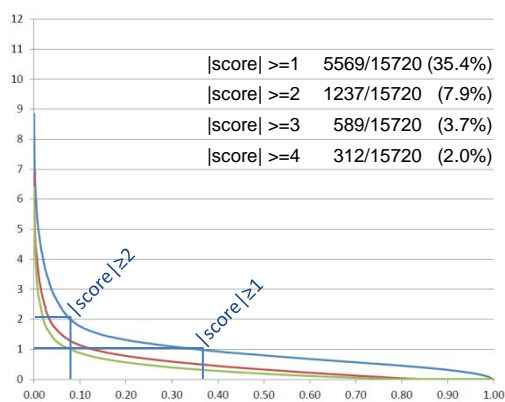
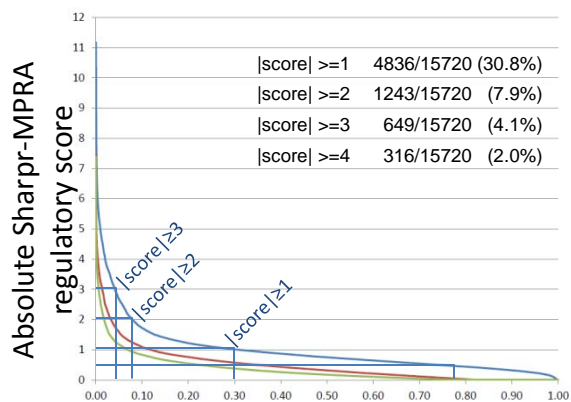


**Supplementary Figure 7b – Sharpr-MPRA regulatory activity score distribution. (a)** Previous page. **(b) MaxPos score distribution at the region level.** Score distribution for all 15,720 tiled region, using the score of a single nucleotide for each, MaxPos, the nucleotide with maximum absolute activity score, shown based on the combined minP and SV40P data (top), minP only (middle row), and SV40P only (bottom), each for HepG2 (left) and K562 (right). Pink bars highlight the subset of regulatory regions with MaxPos ≥ 1 (primarily activating) and MaxPos ≤ -1 (primarily repressive) that are plotted in **Supplementary Data Files 6a-d.** **(c)** Next Page.

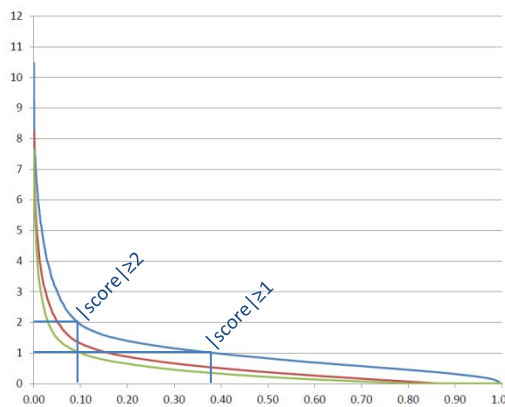
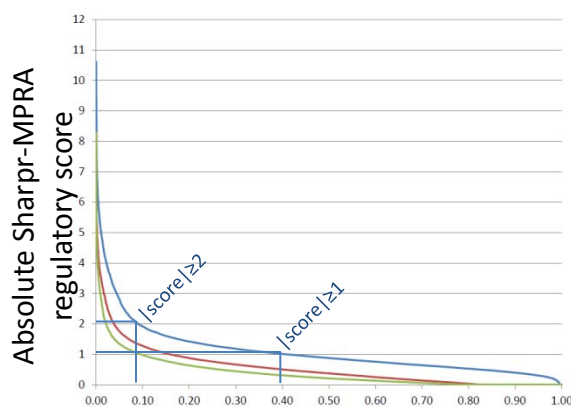
HepG2

K562

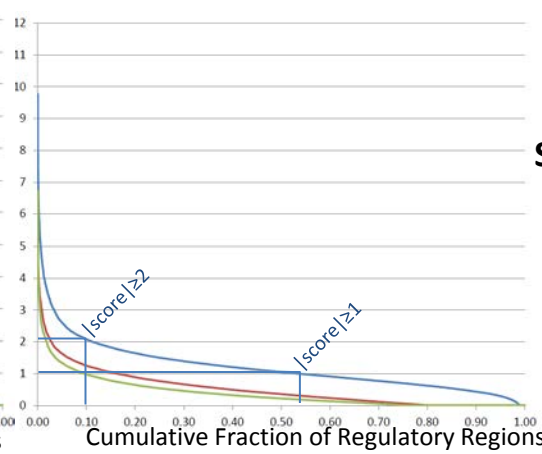
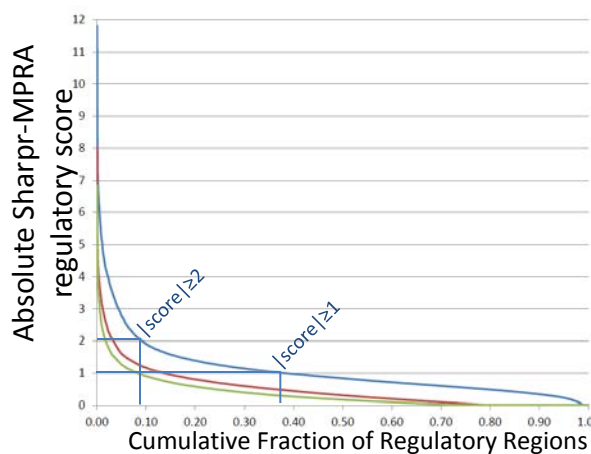
c



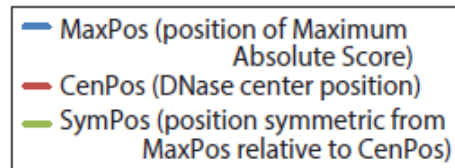
combinedP



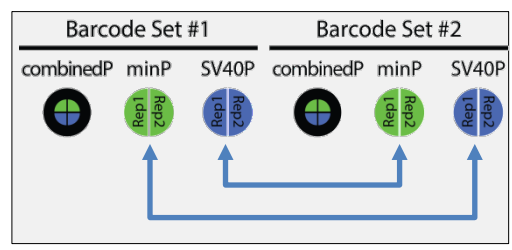
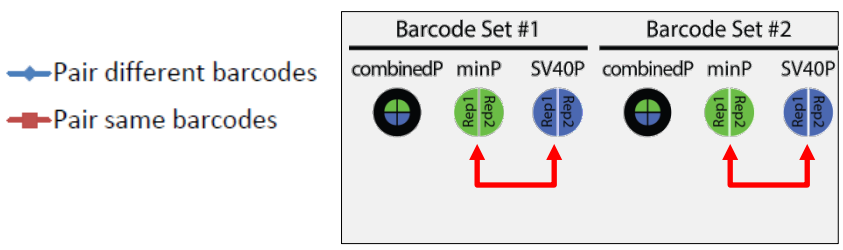
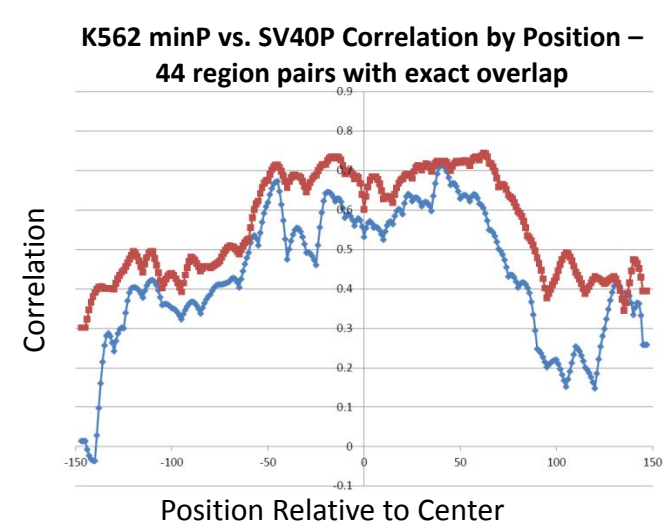
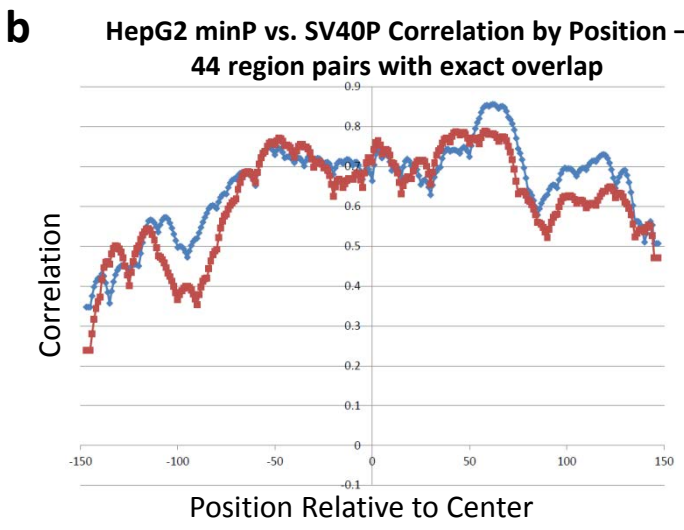
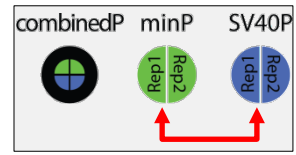
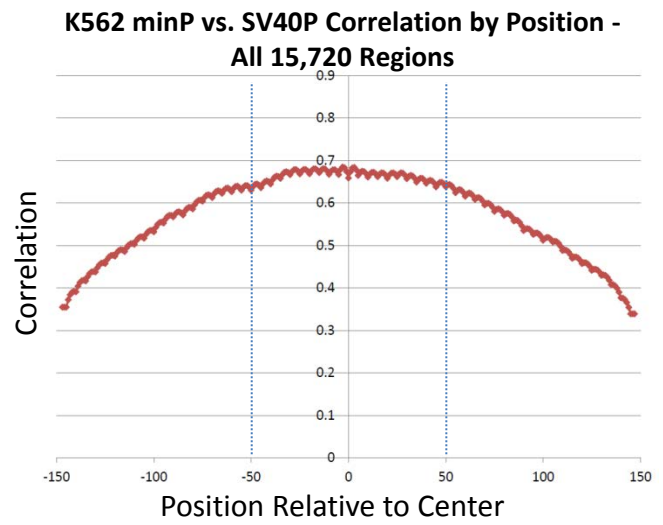
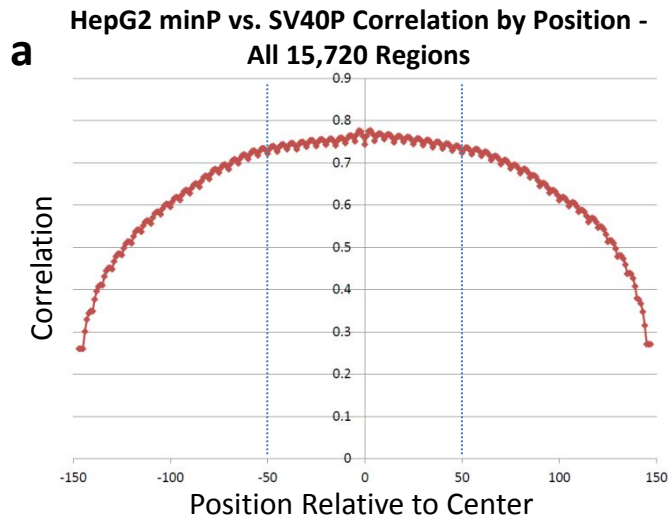
minP



SV40P

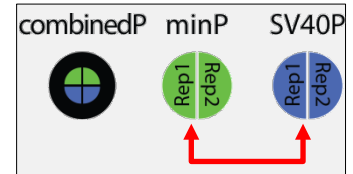
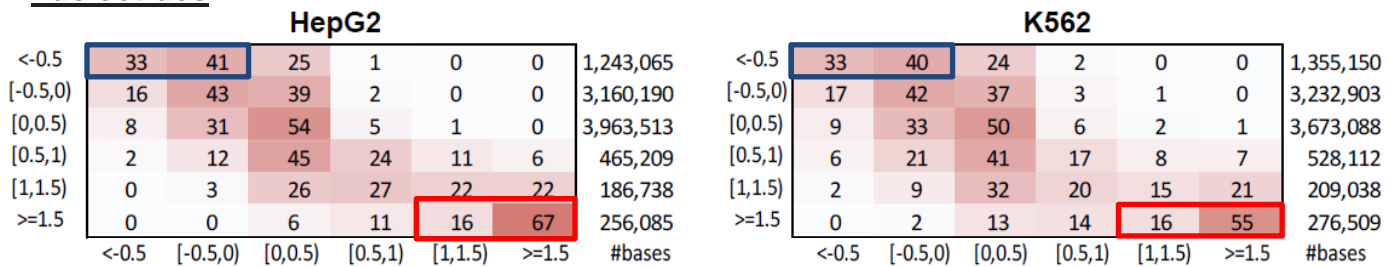


**Supplementary Figure 7c – Sharpr-MPRA Regulatory Activity Score Distribution.** (a,b) Previous pages. (c) **Absolute MaxPos, CenPos, and SymPos score distributions.** Distribution of the absolute Sharpr-MPRA regulatory score (y-axis) for the 15,720 regions at the nucleotide position of maximum absolute score (MaxPos, blue), the center nucleotide position (CenPos, red), and the symmetric nucleotide positions (SymPos, green), each ranked from highest to lowest absolute score (x-axis) for HepG2 (left) and K562 (right), using: combined minP and SV40P (combinedP) score (top row); the minP score only (middle row); and the SV40P score (bottom row). Blue tickmarks indicate the fraction of MaxPos nucleotides with absolute scores above the indicated values and are also listed for specific values based on the combinedP data (top). MaxPos, CenPos, and SymPos nucleotide positions are illustrated in the example of the bottom-left panel of Fig. 3b.

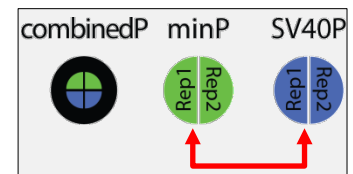
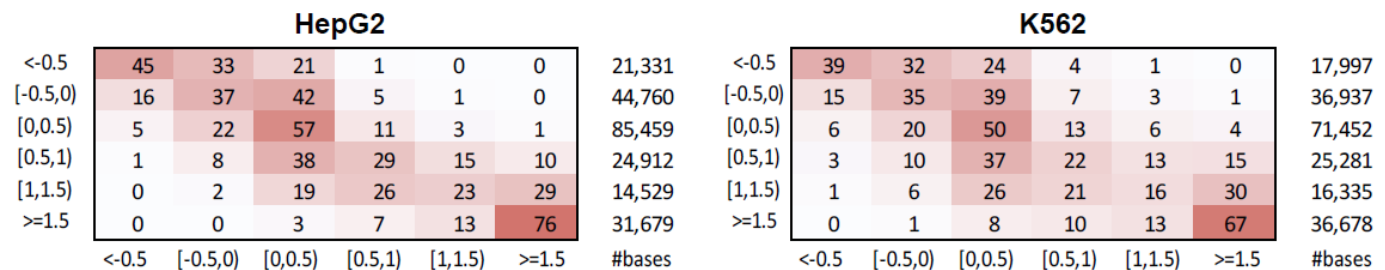


**Supplementary Figure 8 – Correlation between minP and SV40P inferred regulatory activity score by nucleotide position. (a)** Correlation between the inferred regulatory activities based on the minimal and SV40P promoter data (y-axis) as a function of nucleotide position relative to the DNase peak center position in the tiling (x-axis) for HepG2 (left) and K562 (right). We find higher correlation closer to the center where each nucleotide is covered by more reporter constructs (Fig. 3a), and where more meaningful regulation is likely to occur (see Fig. 2a and Supplementary Fig. 2). **(b)** Comparison of correlation (y-axis) between minP and SV40P using the same barcode set (red curves) and using different barcode sets (blue curves) for 88 pairs of inferences in the 44 regions that were tiled twice using the exact same set of reporter sequences for HepG2 (left) and K562 (right).

**a. Correspondence between SV40P and minP scores: all 15,720 tiled regions, all nucleotides**

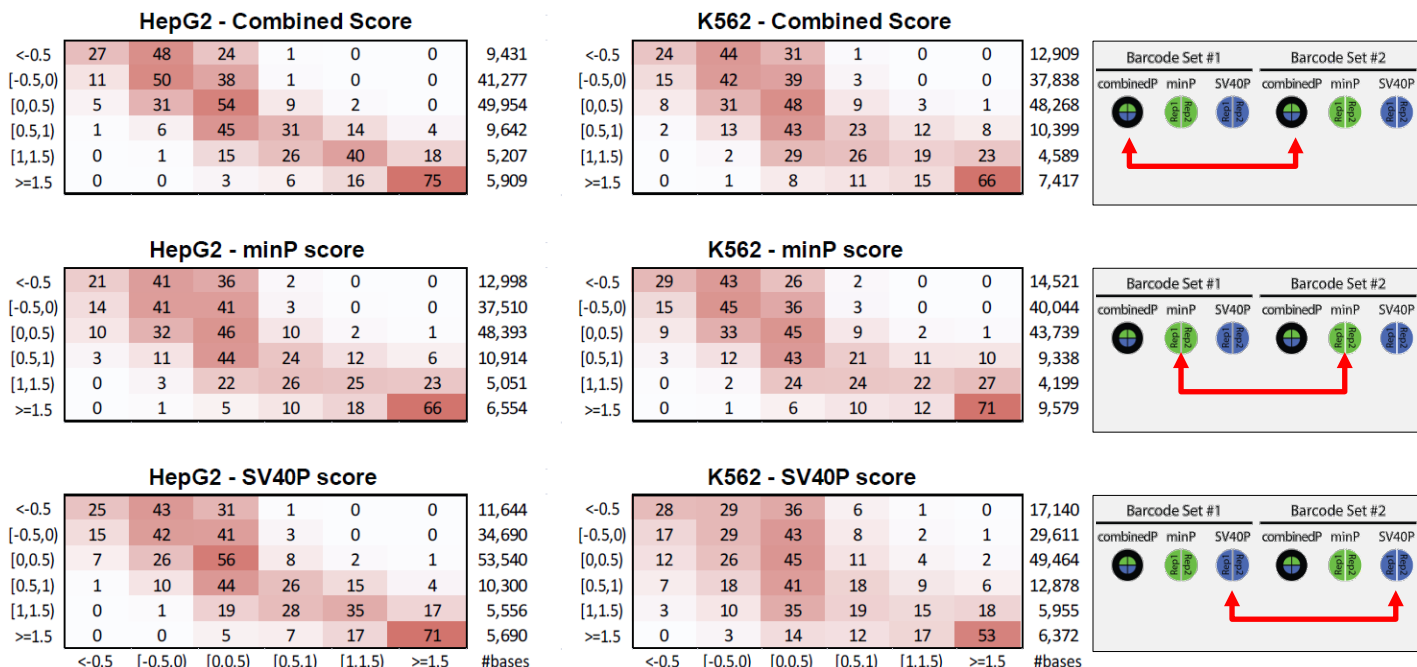


**b. Correspondence between SV40P and minP scores: all 15,720 tiled regions, CENTIPEDE nucleotides**

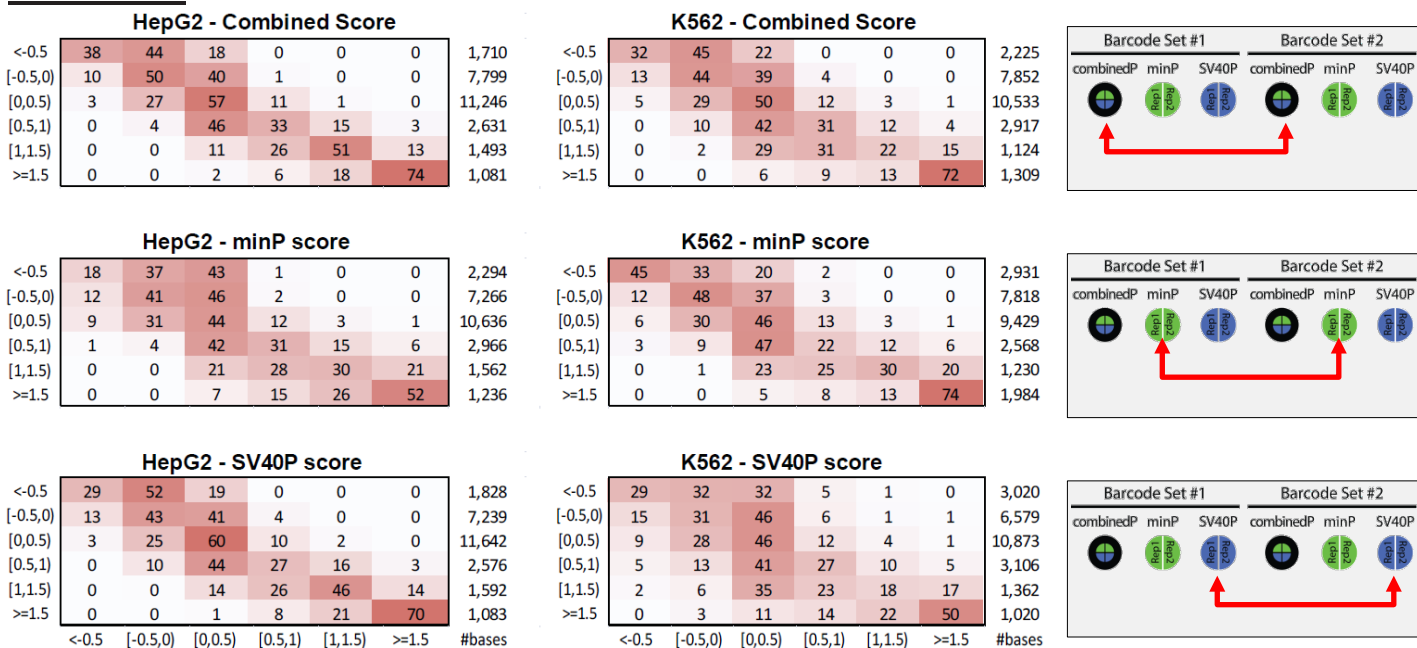


**Supplementary Figure 9 (a,b) – Nucleotide-level Sharpr-MPRA score correspondence at individual bases across promoter and barcode sets. (a)** Comparison of Sharpr-MPRA regulatory activity score across minP and SV40 experiments for all nucleotide positions assigned to each indicated score range for HepG2 (left) and K562 (right), as shown by the red arrow in the diagram. For each row, left-most column indicates the score range, right-most column (#bases) indicates total number of nucleotide positions in that score range, and six colored columns indicate percentage of bases assigned to each score range bin. Each comparison was considered in both directions between promoter types, and each base is counted twice, once for its minP score, and once for its SV40P score. For example, 83% of nucleotides with scores  $\geq 1.5$  in one promoter type show scores  $\geq 1$  in the other promoter type for HepG2 (71% for K562) (red boxes), and 74% of nucleotides with scores  $< -0.5$  in one promoter show scores  $< 0$  in the other promoter for HepG2 (73% for K562) (blue boxes). **(b)** Same as (a) except restricted to those nucleotides overlapping a CENTIPEDE base in the cell type. **(c-f)** Next pages

### c. Score correspondence between barcode sets: 256 exact- or partial-overlap regions, all nucleotides

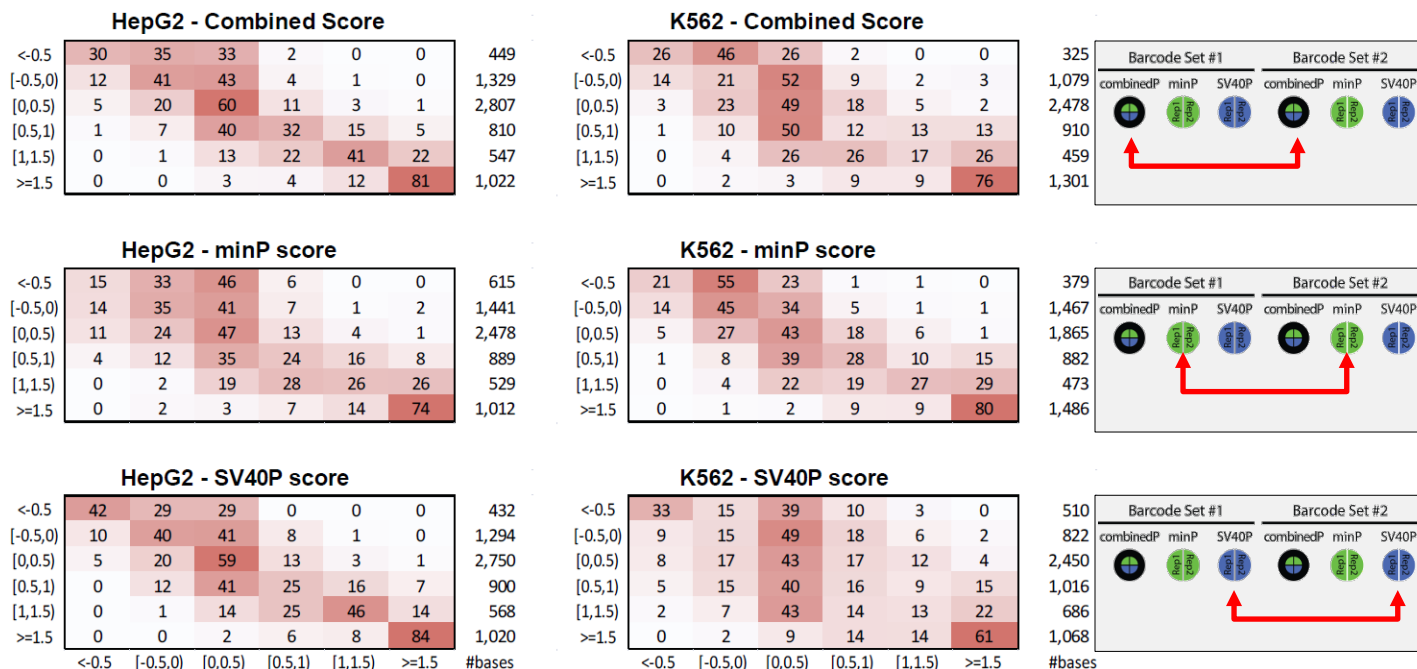


### d. Score correspondence between barcode sets: 44 exact-overlaps region pairs, all nucleotides

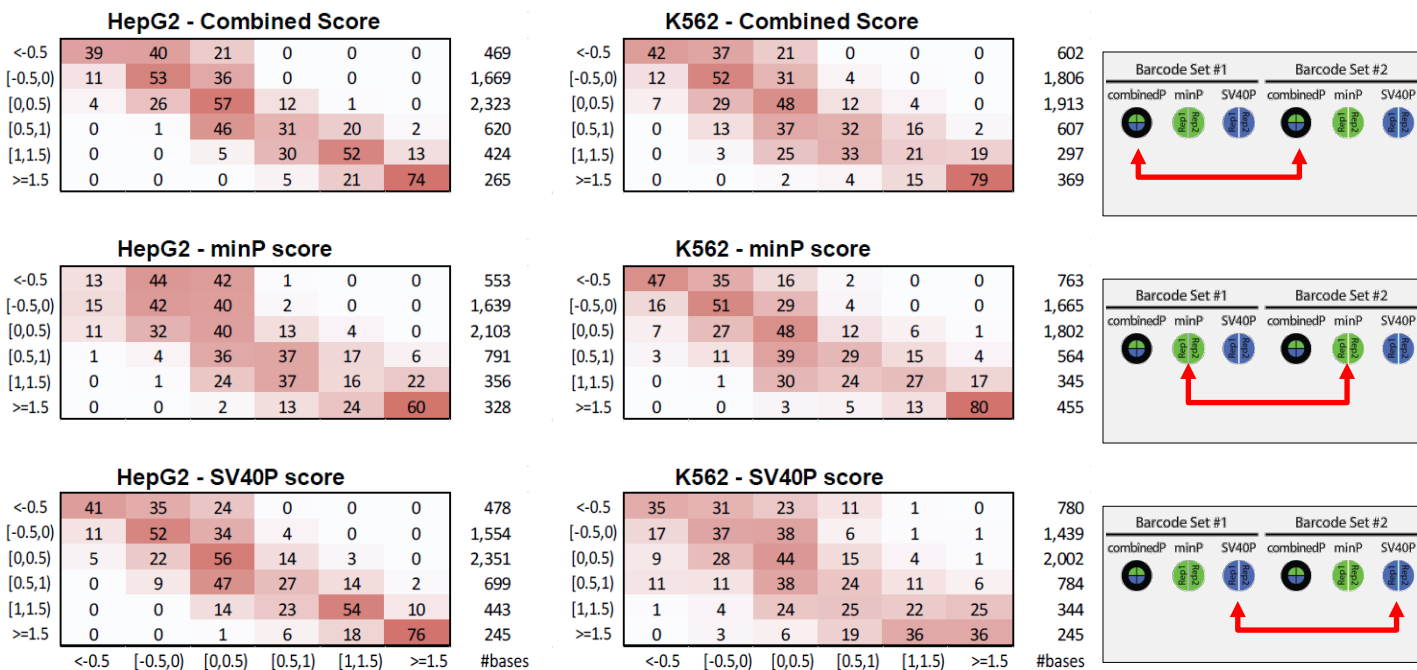


**Supplementary Figure 9 (c,d) – Sharpr-MPRA score correspondence at individual bases across promoter and barcode sets.** (a,b) Previous page. (c) Comparison of Sharpr-MPRA scores at individual bases across different barcode sets for all 256 multiply-tiled regions, including both exact and partial overlaps. Comparisons are shown for the minP and SV40P combined data (top), minP data only (middle row), and SV40P data only (bottom row), with format similar to panel a. (d) Same as (c) except for only the subset of 44 regions that were covered by two sets of barcodes with exact tile overlap. (e,f) Next page.

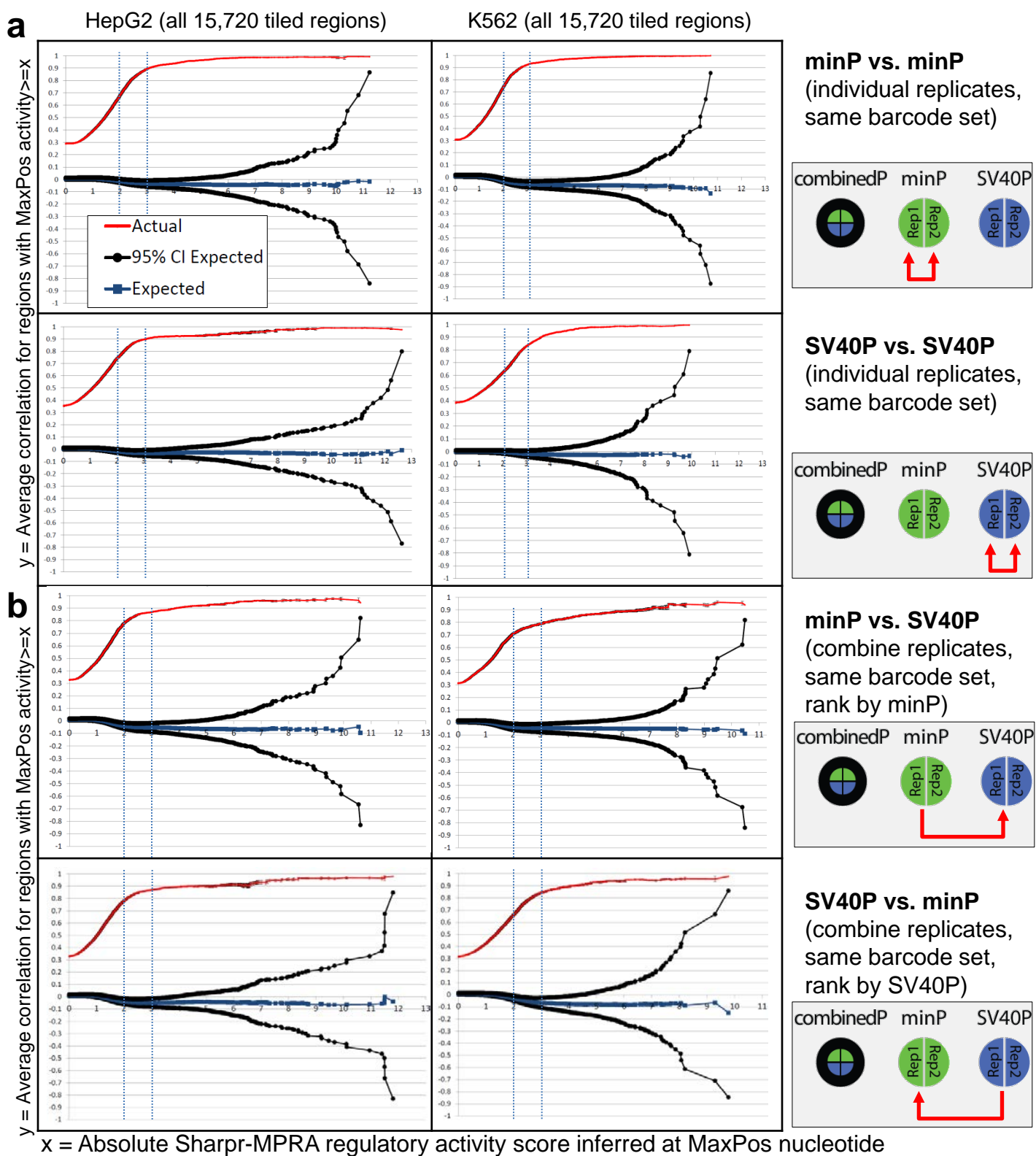
## e. Score correspondence between barcode sets: 256 exact- or partial-overlap regions, CENTIPEDE nucleotides



## f. Score correspondence between barcode sets: 44 exact-overlap regions, CENTIPEDE nucleotides



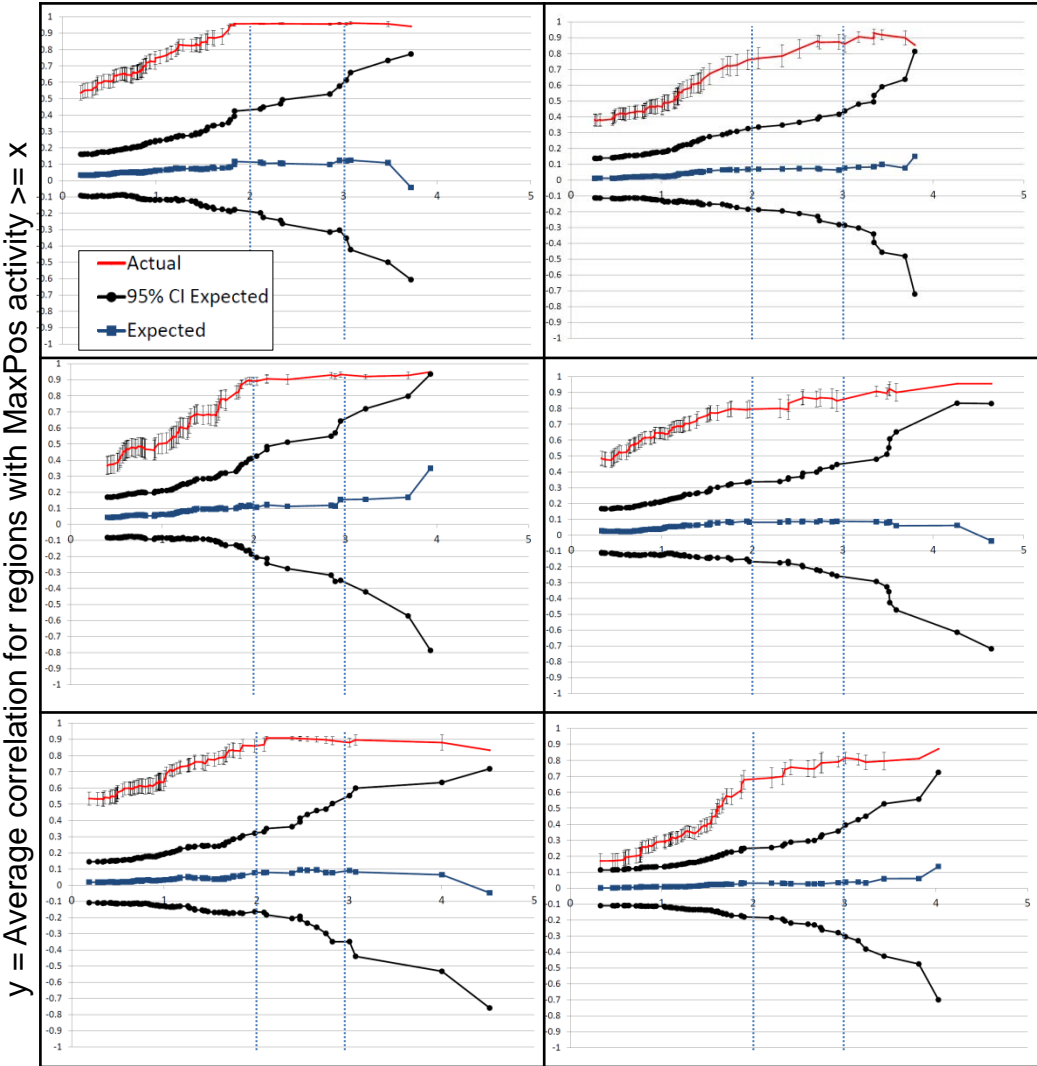
**Supplementary Figure 9 (e,f) – Sharpr-MPRA score correspondence at individual bases across promoter and barcode sets. (a-d) Previous pages. (e) Comparison of Sharpr-MPRA scores inferred across different barcode sets for all 256 multiply-tiled regions, including both exact and partial overlaps similar to (c), except restricted to those nucleotides overlapping a CENTIPEDE base in the cell type. (f) The same as (e), except for only the subset of 44 regions that were covered by two sets of barcodes with exact tile overlap.**



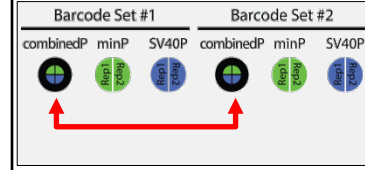
**Supplementary Figure 10a,b – Sharpr-MPRA within-region reproducibility correlation analysis. (a)** Average within region Pearson correlation in activity scores across individual replicates of minP and SV40P experiments (y-axis) for regions meeting or exceeding varying maximum absolute score values (x-axis) for HepG2 (left) and K562 (right), with comparison performed indicated by the red arrow in the corresponding diagram. Each region was included twice, once based on the MaxPos value from each replicate. Observed correlation (red curve) are compared to the expected value (blue) and 95% confidence interval (black) of 10,000 row-wise permutations of regions. Error bars indicate standard error. **(b)** Average Pearson correlation between minP and SV40P similar to panel (a). Each region was included only once and the MaxPos score on the x-axis is based on minP for the first row and SV40P for the second row. **(c)** Next two pages.

**C** HepG2 (88 MaxPos with exact overlap tiling)

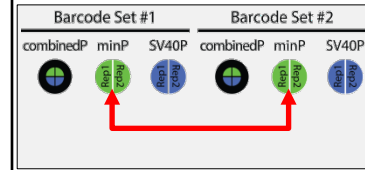
K562 (88 MaxPos with exact overlap tiling)



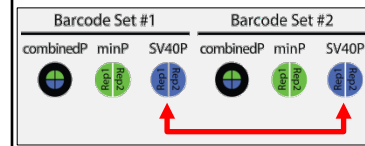
**Barcode set #1 vs. barcode set #2 (combinedP scores)**



**Barcode set #1 vs. barcode set #2 (minP scores)**



**Barcode set #1 vs. barcode set #2 (SV40P scores)**

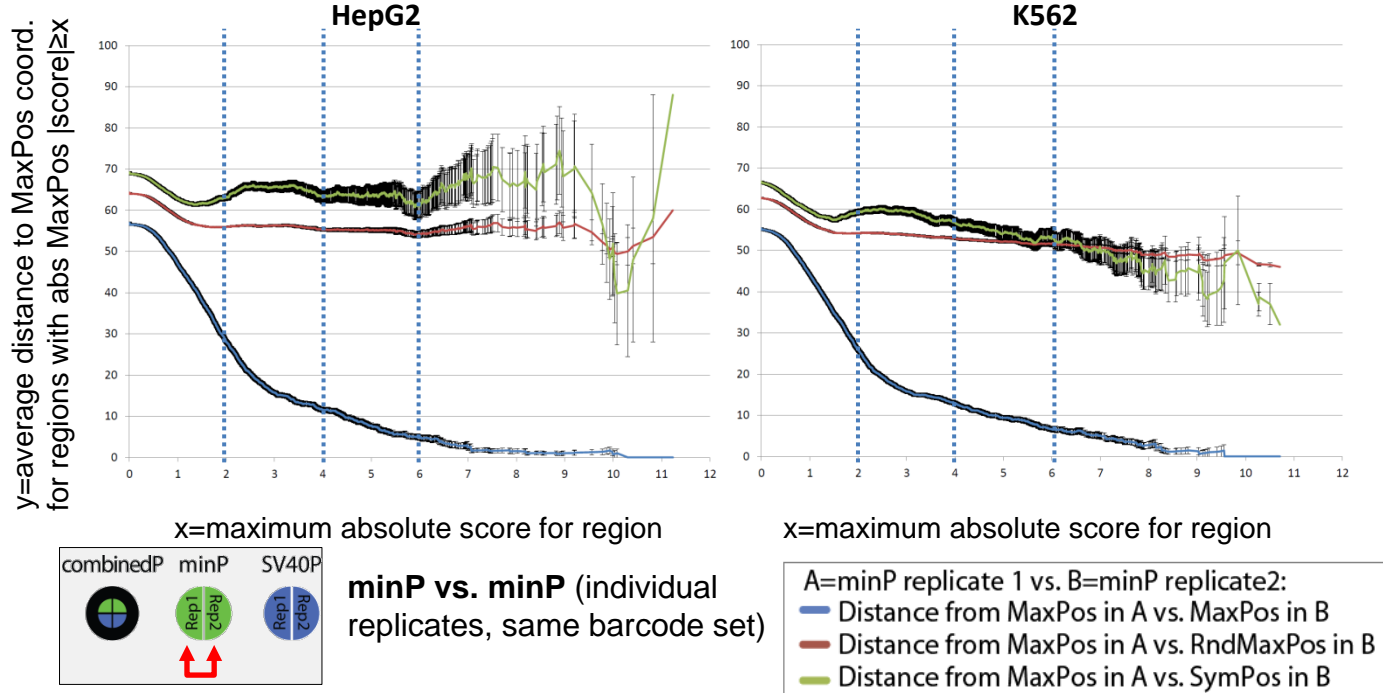


x = Absolute Sharpr-MPRA regulatory activity score inferred at MaxPos nucleotide

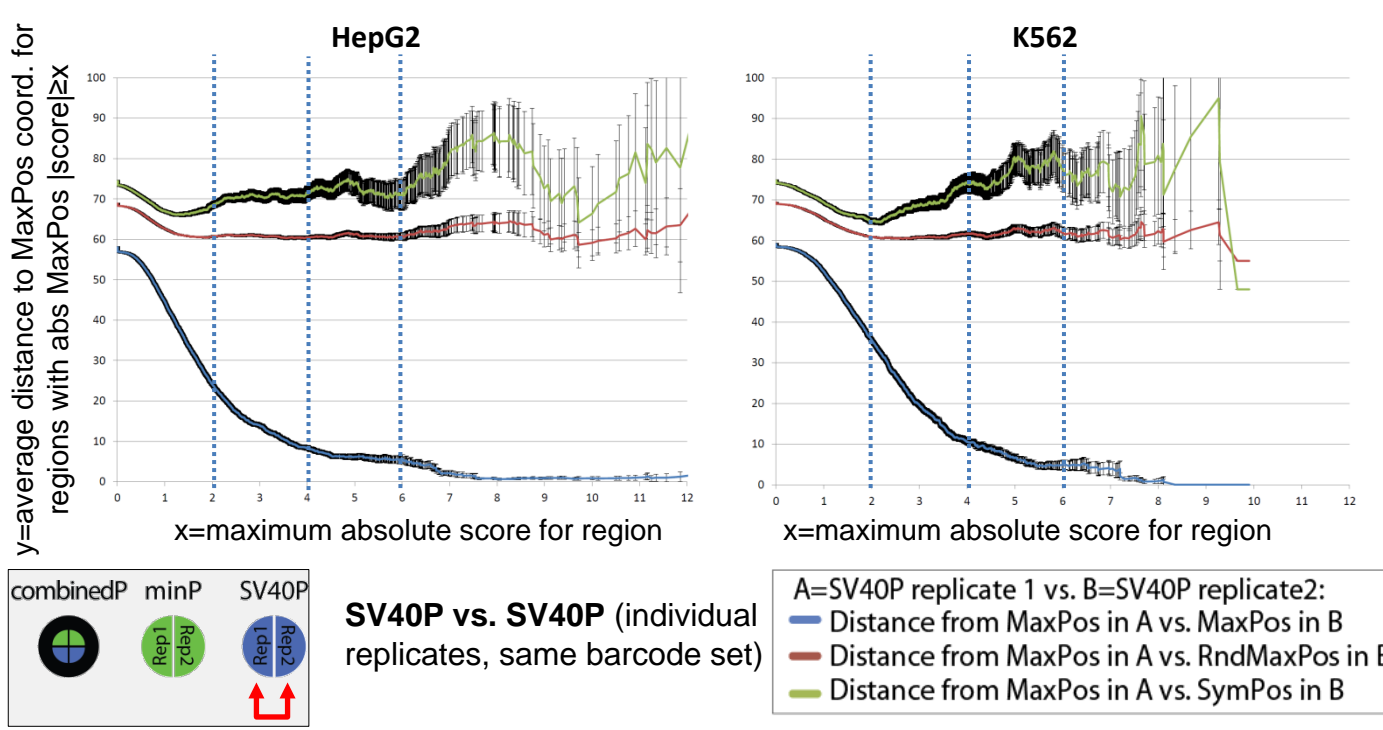
**Supplementary Figure 10c. Sharpr-MPRA within-region reproducibility correlation analysis. (a-b)** See previous page. **(c)** Similar to (a,b), but with all comparisons performed across different barcode sets of the same experiment. Top: Barcode comparison using combined minP and SV40P data. Middle row: Barcode comparison using minP data. Bottom: Barcode comparison using the SV40P data. Each panel shows 88 points corresponding to two MaxPos values for each of the 44 regions with exact overlap tiling.



**a. MaxPos position agreement between minP replicates**

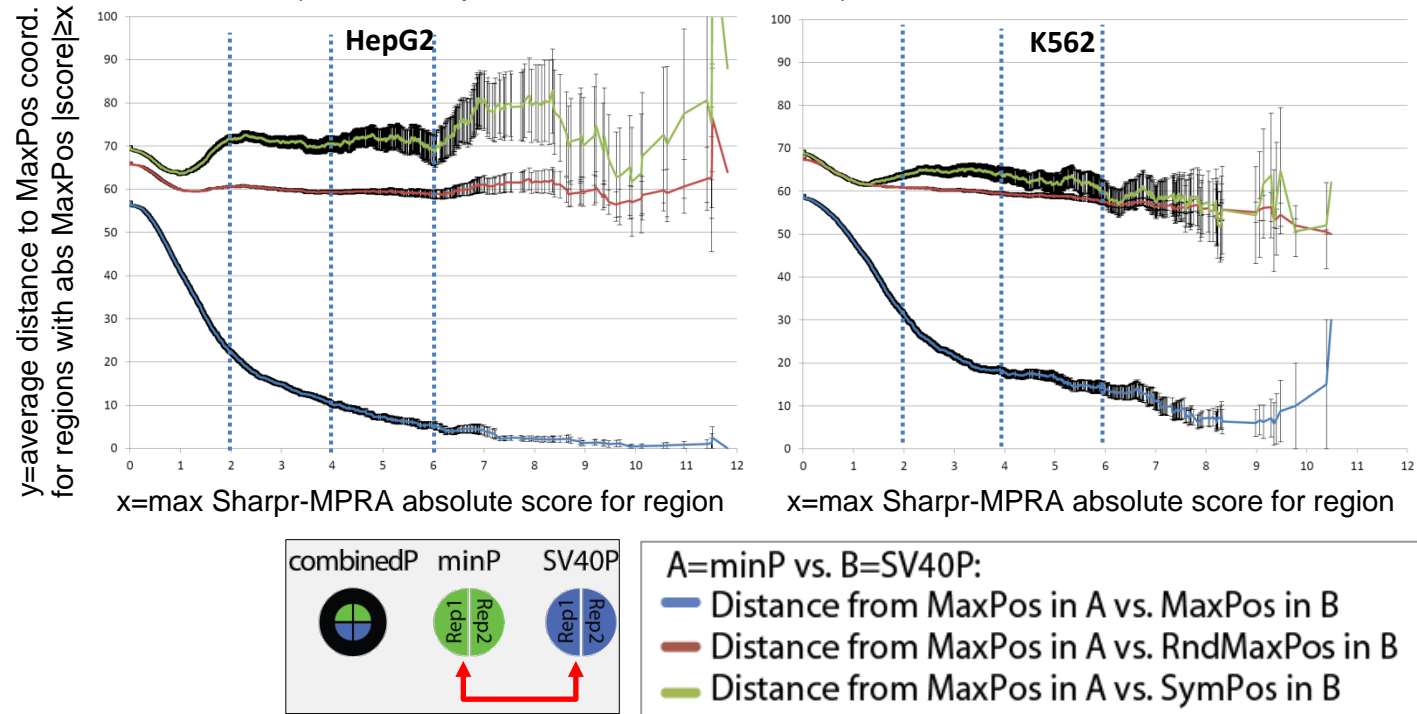


**b. MaxPos position agreement between SV40P replicates**

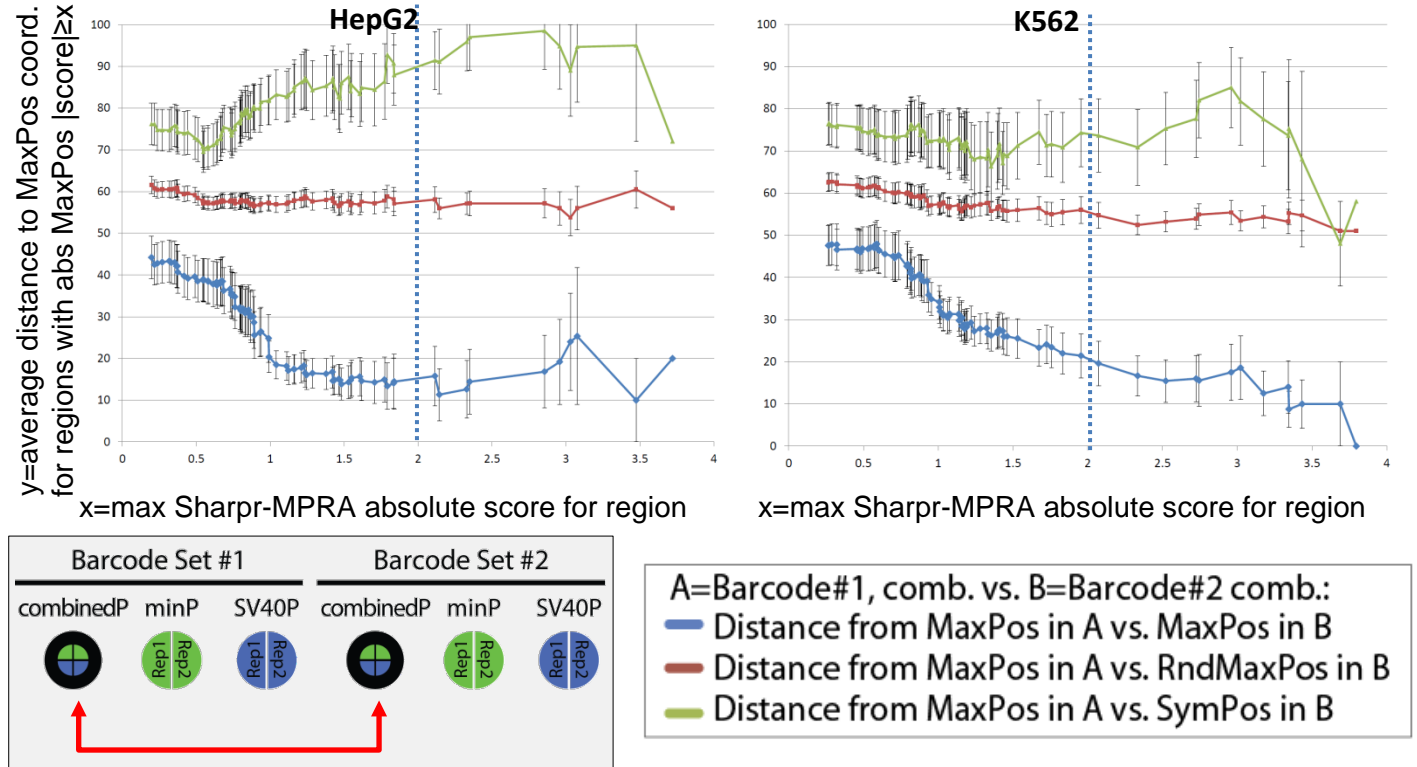


**Supplementary Figure 11a,b. Agreement in position of maximum absolute activity (MaxPos) between replicates, promoter types, and barcodes.** (a) Average absolute distance (y-axis) between the position of maximum absolute activity (MaxPos) in one minP replicate and the other for those regions meeting or exceeding varying maximum absolute score values (x-axis) for HepG2 (left) and K562 (right) (MaxPos, blue line). Each distance was included twice, once based on each of the two MaxPos values in the pair. This is compared to the expected distance between MaxPos in one minP-replicate and a position in the other minP replicate if sampled randomly from the distribution of all MaxPos positions across all regions (RndMaxPos, red line), and also compared to the distance between MaxPos in one minP replicate and the symmetric position in minP replicate (SymPos, green line). This shows that even when MaxPos regions are highly off-centered, there is significantly higher agreement in the MaxPos position than expected by chance. (b) Same comparison as (a) for SV40P individual replicates for HepG2 (left) and K562 (right). (c,d) Next page.

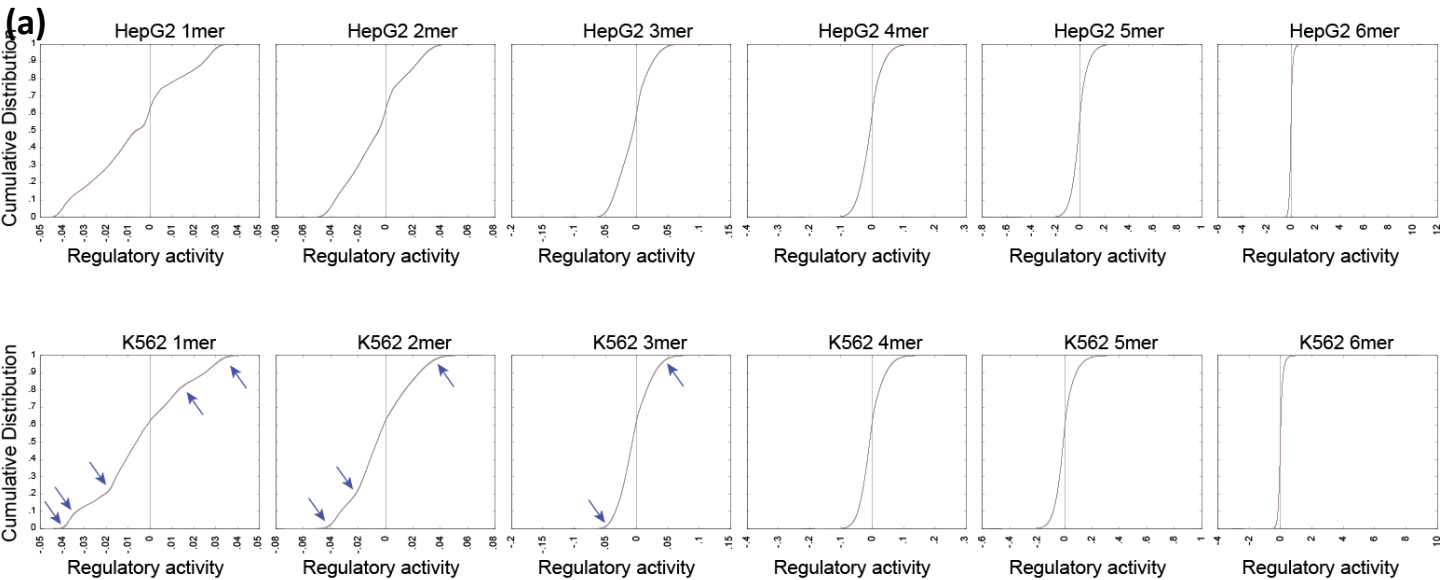
**c. minP vs. SV40P (combine replicates, same barcode set)**



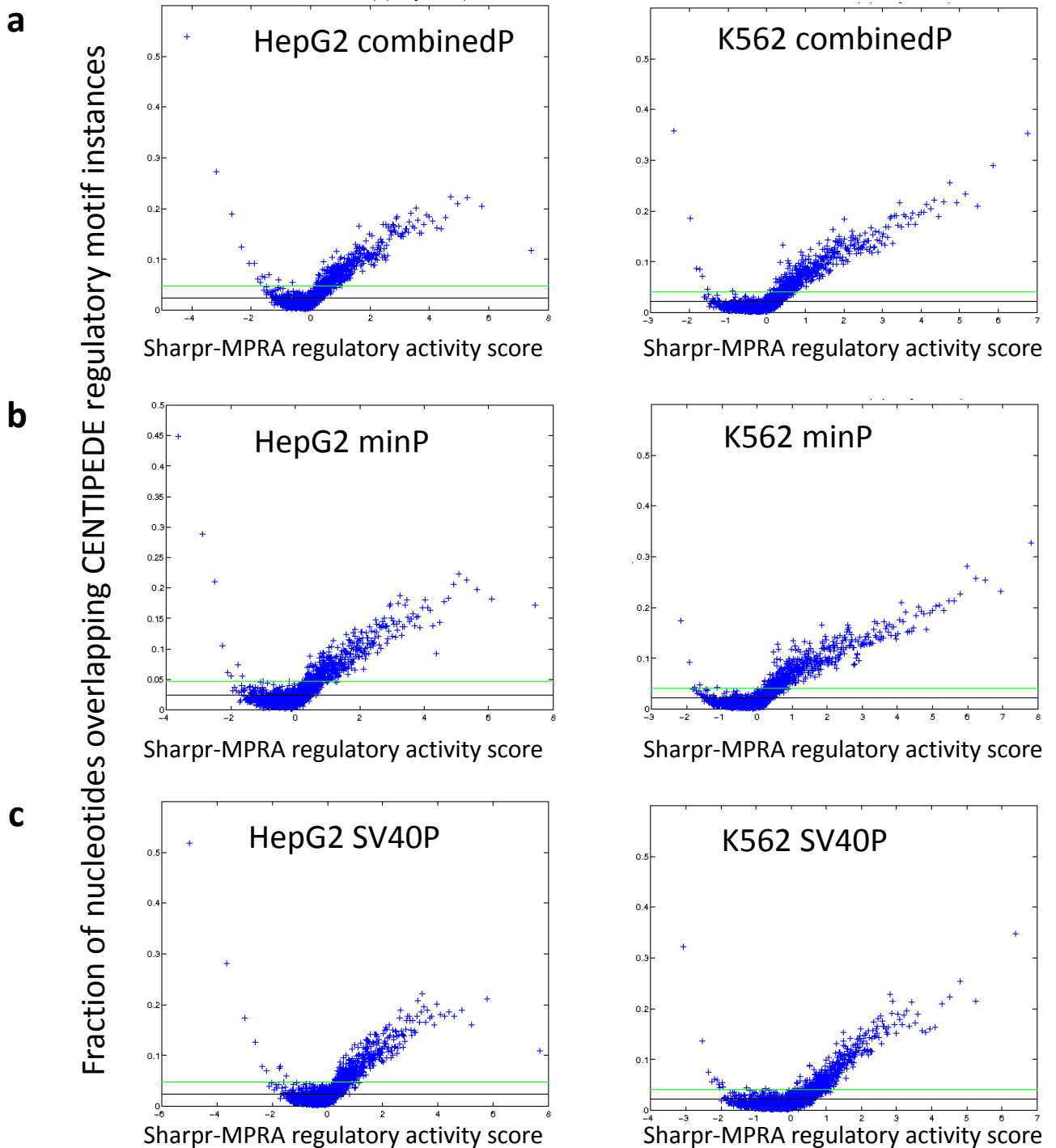
**d. Barcode set #1 vs. barcode set #2, in 44 exact-overlap regions (combinedP score)**



**Supplementary Figure 11c,d. Agreement in position of maximum absolute activity (MaxPos) between replicates, promoter types, and barcodes.** (a,b) Previous page. (c) Same comparison as (a), but the comparisons are between SV40P and minP. (d) Same comparisons as (a), but comparing the combined minP and SV40P data for two different barcode sets for the 44 regions tiled with exact overlap from two barcode sets.

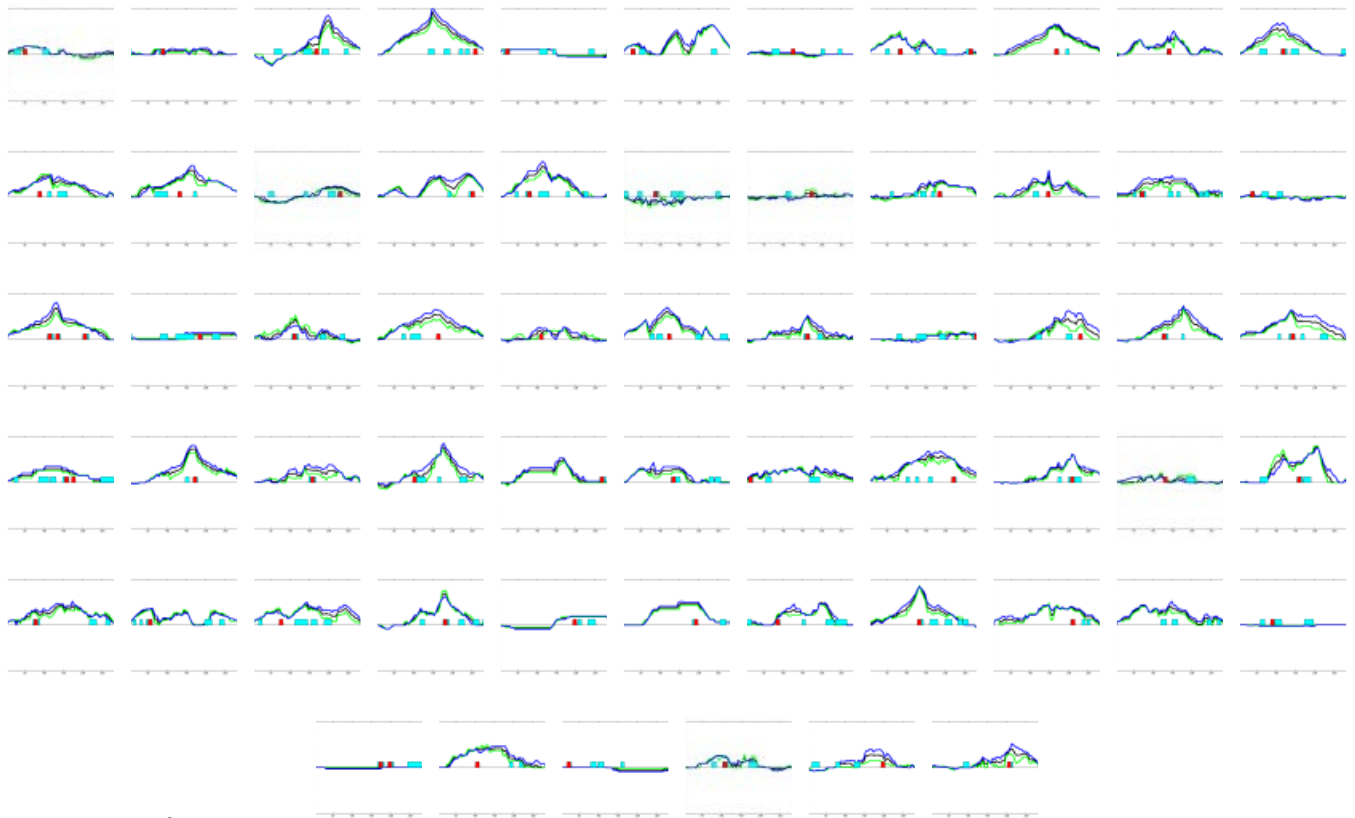


**Supplementary Figure 12 –Effect of the k-mer sequence within the 10-nucleotide barcode sequence on the Sharp-MPRA regulatory activity scores. (a)** Cumulative distribution of the average inferred activity (red) over each combination of k-mer sequence in the barcode, presence in each of the 31 tile positions, and the 295 inferred activity positions. Also shown is the expected (blue) and 95% confidence intervals (black) based on 400 randomizations in which all the barcodes were randomly reassigned to reporter sequences. In all HepG2 plots, there is no visible difference between the observed, expected, and 95% confidence intervals. Arrows denote the visible differences where the observed distribution (red) deviates from the expected distribution (blue) and the 95% confidence intervals (black). **(b)** Same as (a), but showing the difference between observed (red curve in panel a) and expected (blue curve in panel (a)) for each k-mer length and each cell type. Arrows denote the difference between the distribution of observed and expected activity for short k-mers in K562. In particular, AA and AT di-nucleotides were over-represented in the last two positions of barcodes with low reporter activity, suggesting a technical bias in reverse transcription initiation, rather than a biological role of 3'UTR motifs, as the last two barcode positions are the first to be reverse-transcribed. Note that the highlighted differences are on the order of 0.01, compared to activity scores ranging between 0.5-8 for positions that show motif enrichments (see for example **Fig. 3c**, or **Supplementary Fig. S13**).



**Supplementary Figure 13 – Overlap with CENTIPEDE predicted transcription factor binding sites as a function of Sharpr-MPRA regulatory activity score.** Extended version of Fig. 3c left, but showing the results for both K562 and HepG2 and the combined data, minP only, and SV40P only data. Each point represents the average of 927 nucleotide positions in each of 5,000 quantiles. Horizontal black line shows the expected overlap averaged across all 295 nucleotide positions of each region, and the green line shows the expected overlap fraction at the center nucleotide position. These are shown for: **(a)** the combined minP and SV40P data for HepG2 (left) and K562 (right); **(b)** the minP data only for HepG2 (left) and K562 (right); **(c)** the SV40P data only for HepG2 (left) and K562 (right).

## a GABPA – HEPG2

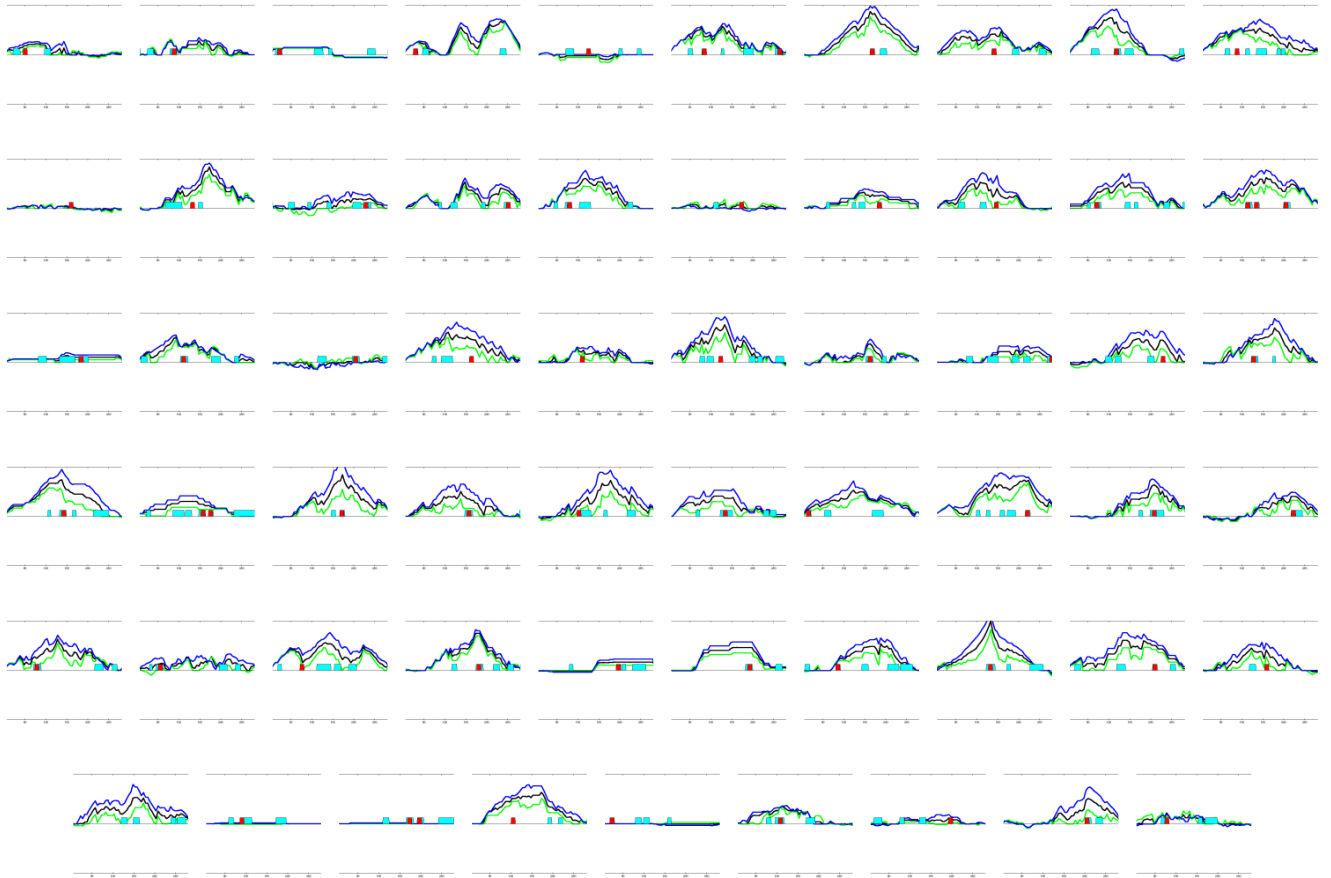


## b REST/NRSF – HEPG2



**Supplementary Figure 14 – Individual Examples. (a)** The figure shows relative Sharpr-MPRA regulatory activity scores overlapping all cell type specific binding sites for GABPA in HepG2 cells predicted based on the CENTIPEDE method<sup>8</sup> contained in regions tested. In red are the GABPA motifs, while in cyan are predicted binding sites in HepG2 for other regulators. The Sharpr-MPRA score based on the SV40P data is shown in green, the minP data in blue, and the combination of SV40P and minP in black. **(b)** Same as (a) except for the predicted binding sites of the repressor REST in HepG2 cells. **(c,d)** Next page.

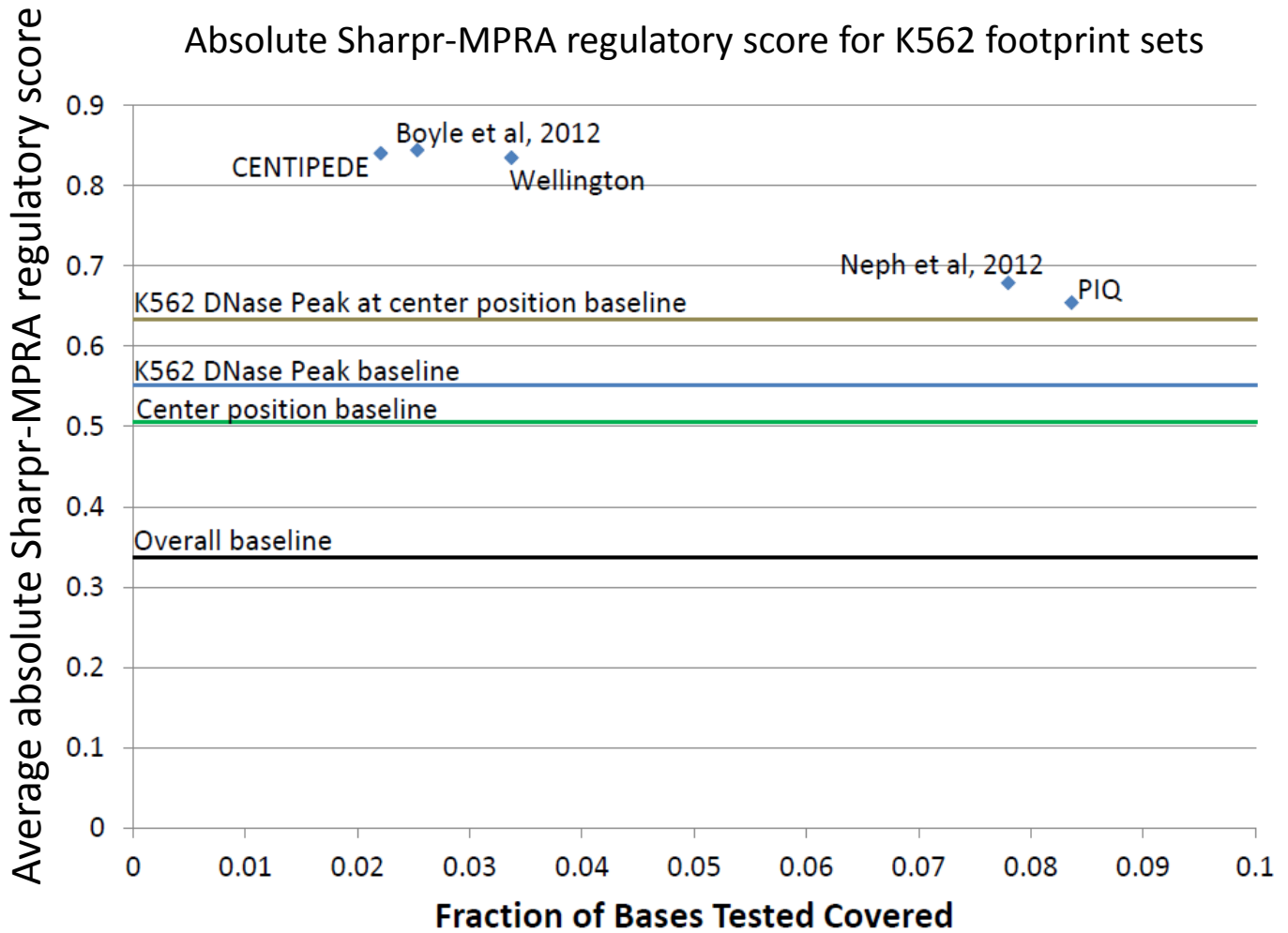
### c GABPA – K562



### d REST/NRSF – K562



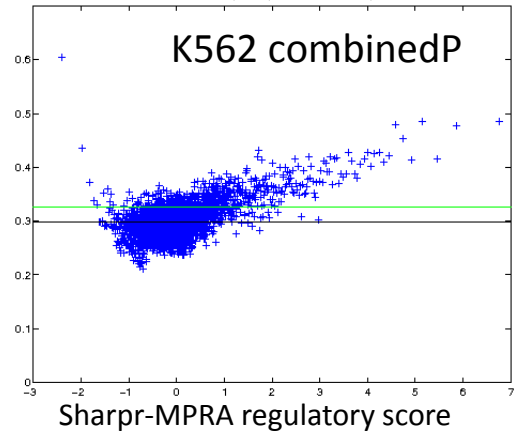
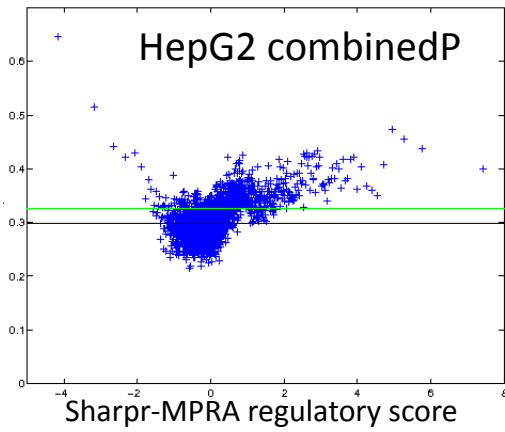
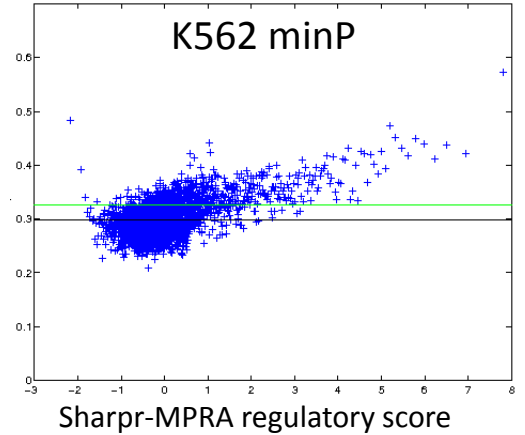
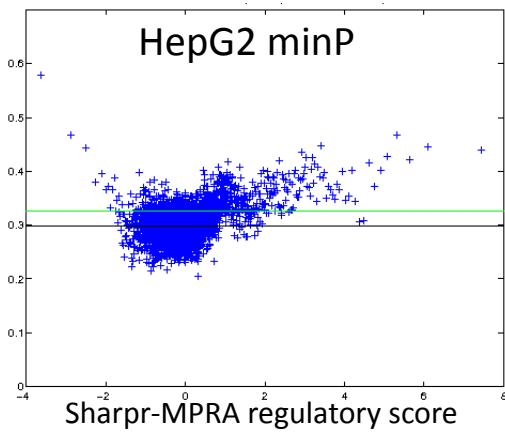
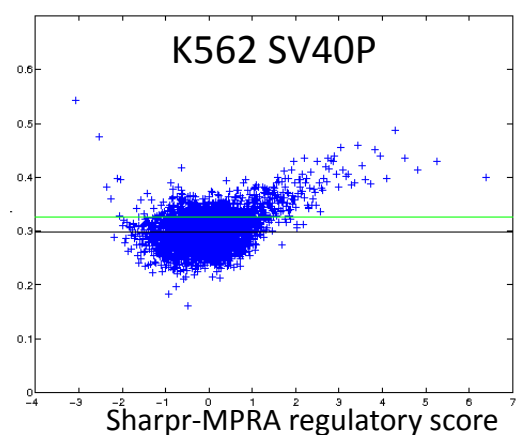
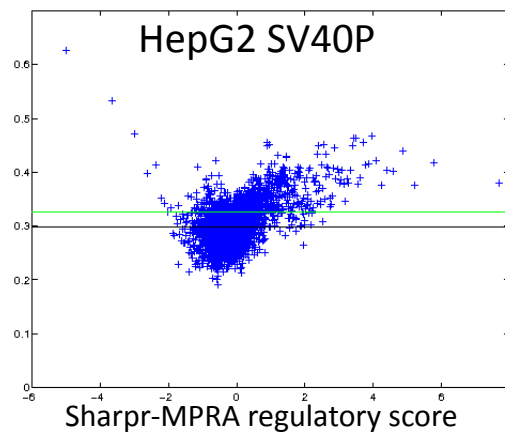
**Supplementary Figure 14 – Individual Examples. (a,b)** Previous page. **(c,d)** The same as (a,b) except for K562 cells.



**Supplementary Figure 15 – Comparison of Different Footprint Sets.** The plot evaluates in K562 cells the average absolute regulatory activity (y-axis) for positions overlapping predictions of locations of transcription factor binding sites previously predicted based on DNase footprint information in K562 cells by five different methods. Two of the five methods also use motif information, CENTIPEDE<sup>8</sup> and PIQ<sup>41</sup>, while the other three methods are motif independent, Wellington<sup>40</sup> and the methods of Boyle et al<sup>5</sup> and Neph et al<sup>6</sup>. The x-axis shows the fraction of nucleotides each of these footprint sets overlap, showing the two footprint sets that overlap more nucleotides tested<sup>6,41</sup> had a relatively lower average absolute activity compared to the other three. All five sets had greater absolute activity than four different baselines: (1) at the center position restricted to regions overlapping a K562 DNase peak, (2) at any position overlapping a K562 DNase peak, (3) at the center position over all regions, and (4) over all positions in all regions tested.

**a**

Fraction overlapping ENCODE regulatory motif instances

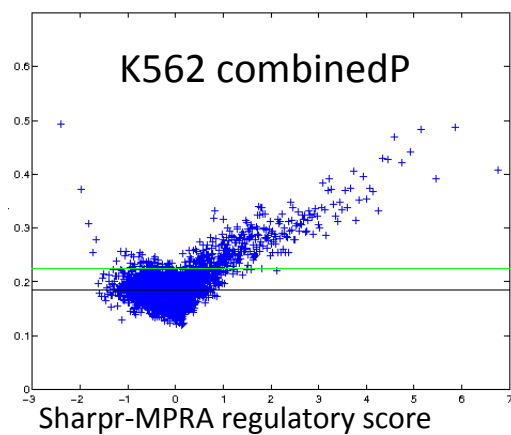
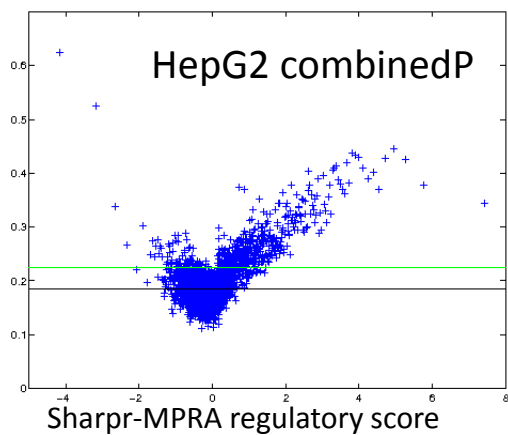
**b****c**

**Supplementary Figure 16 – Overlap with motif instance predictions as a function of Sharpr-MPRA regulatory activity score.** These are analogous plots to **Supplementary Fig. 13** except they are for the set of nucleotides covered by a motif instance prediction from Ref. 13, which does not use conservation or make cell type specific predictions. The plot is for Sharpr-MPRA regulatory activity scores in HepG2 (left) and K562 (right) cells based on: **(a)** minP and SV40P combined data; **(b)** minP data only; and **(c)** SV40P data only.

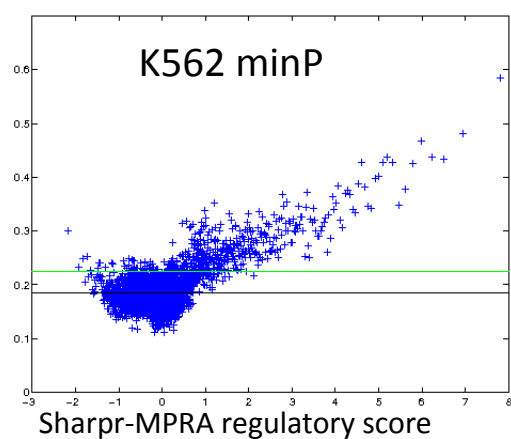
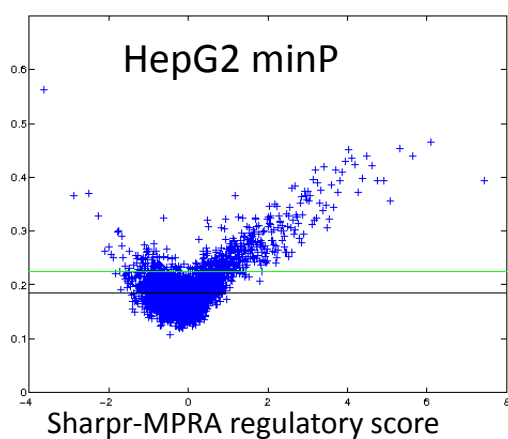


Fraction overlapping evolutionarily-conserved nucleotides

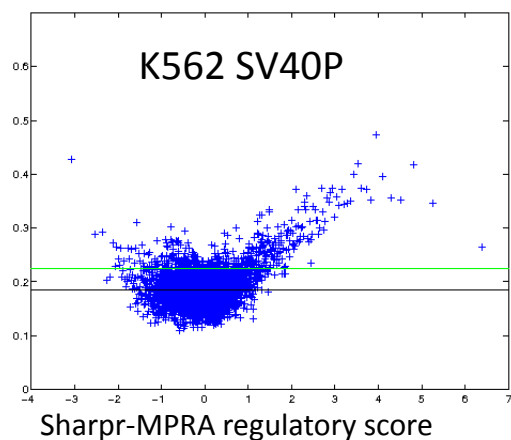
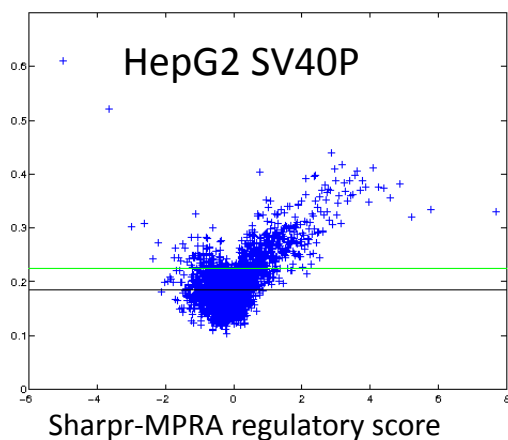
**a**



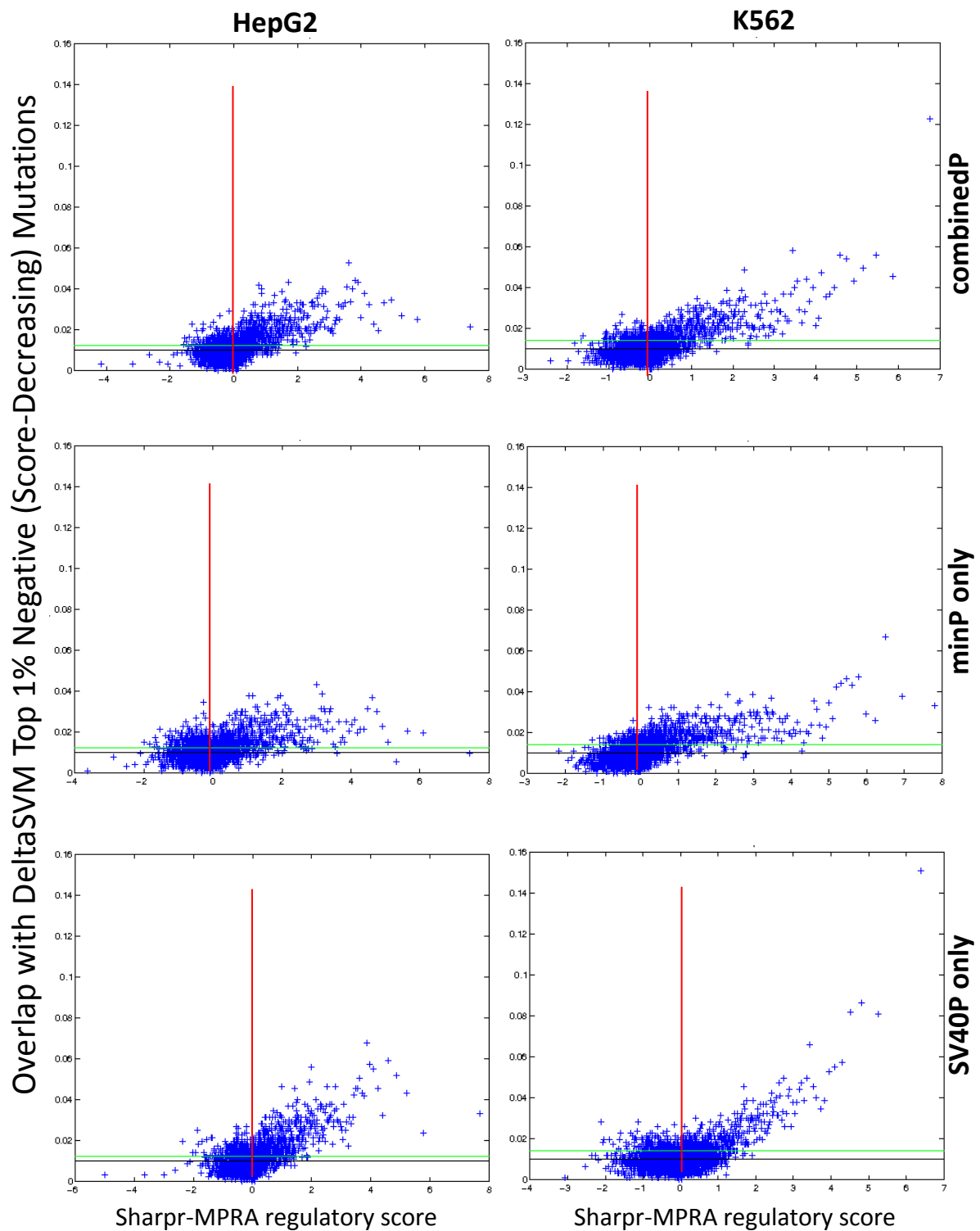
**b**



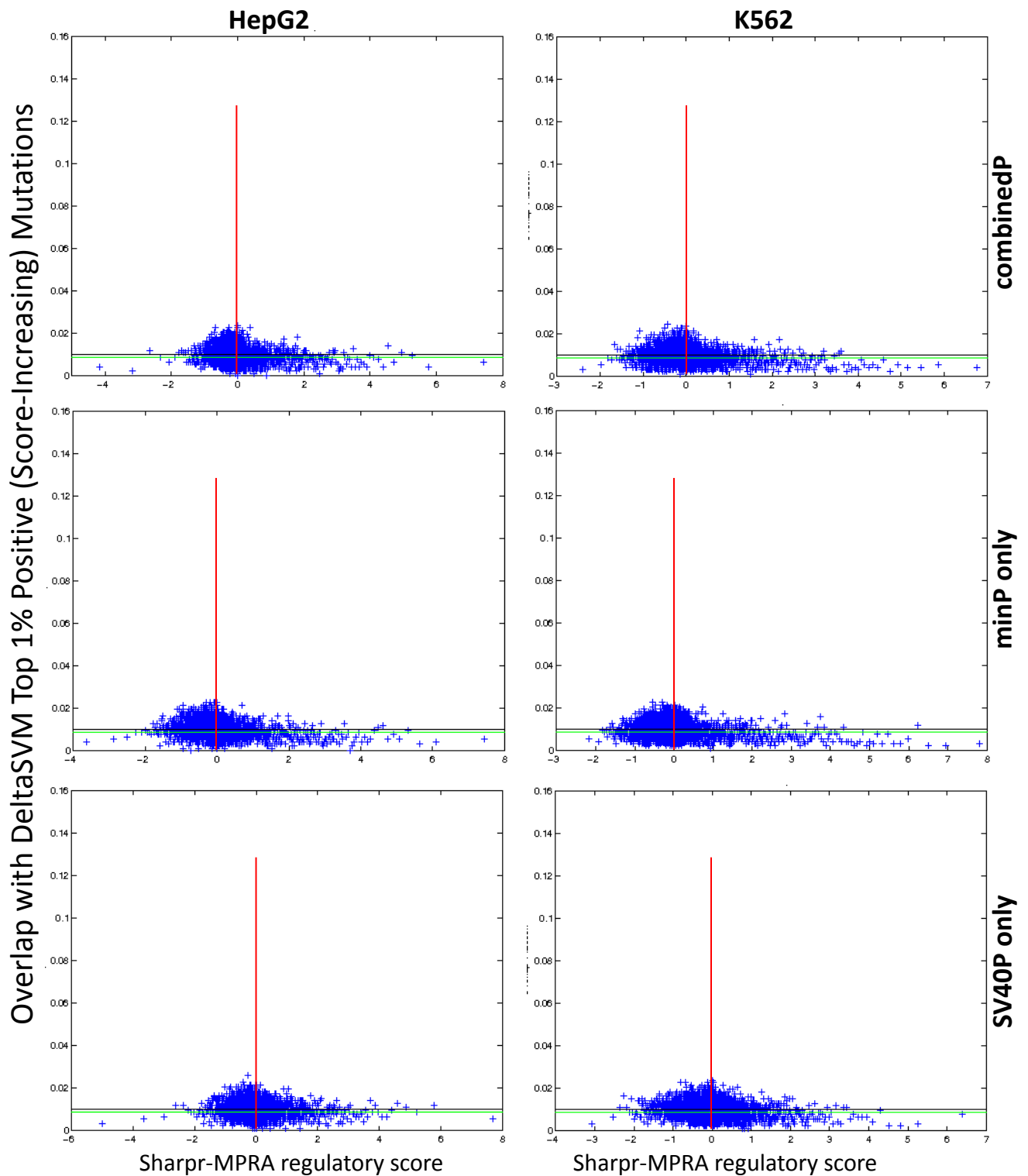
**c**



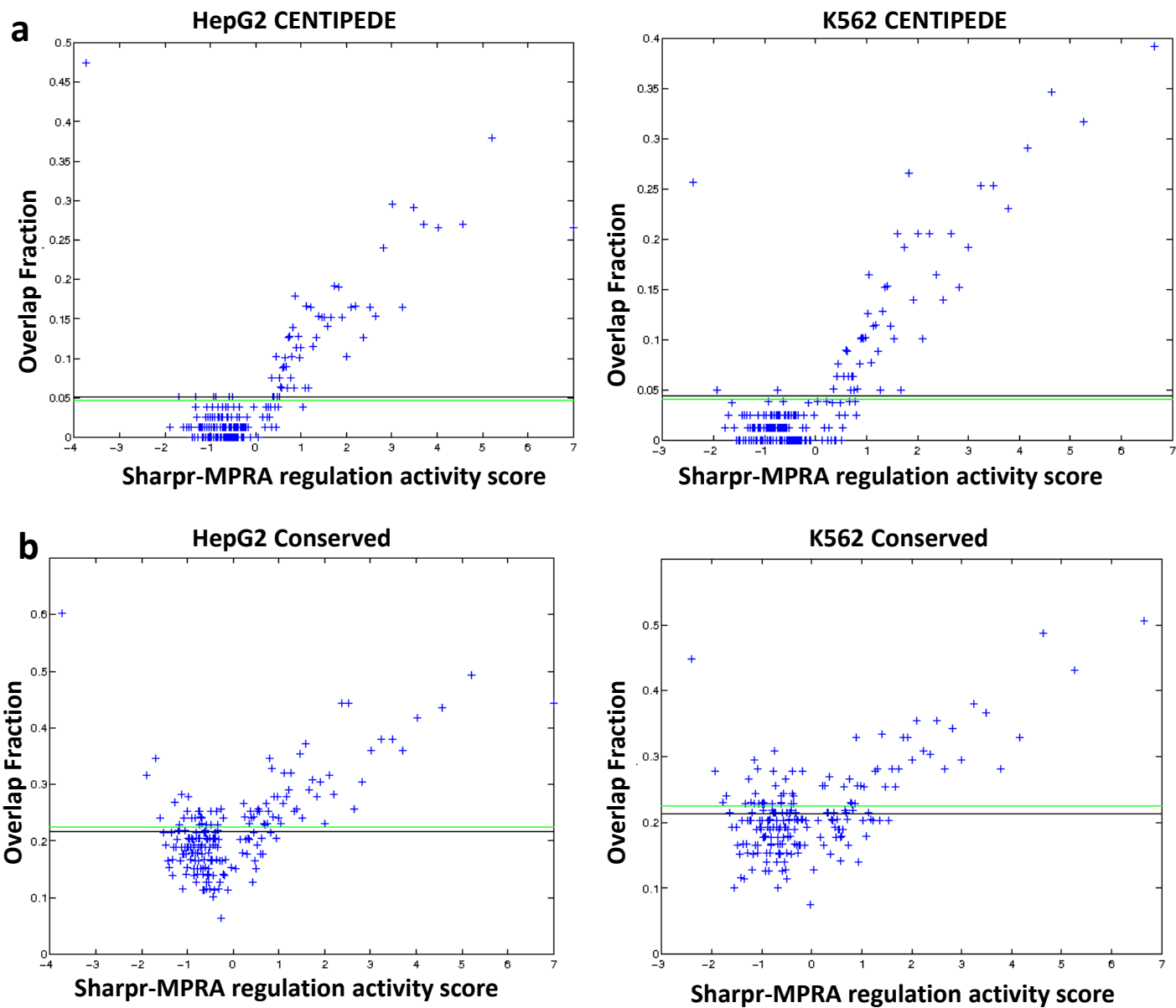
**Supplementary Figure 17 – Overlap with Conserved Elements as a function of Sharpr-MPRA regulatory activity score.** Extended version of **Fig. 3c** right and analogous to **Supplementary Fig. 13**, but showing the results for bases overlapping conserved elements from the SiPhy-PI method<sup>42,65</sup> for HepG2 cells (left) and K562 cells (right), based on: **(a)** the combined minP and SV40P data; **(b)** minP data only; and **(c)** SV40P data only.

**a**

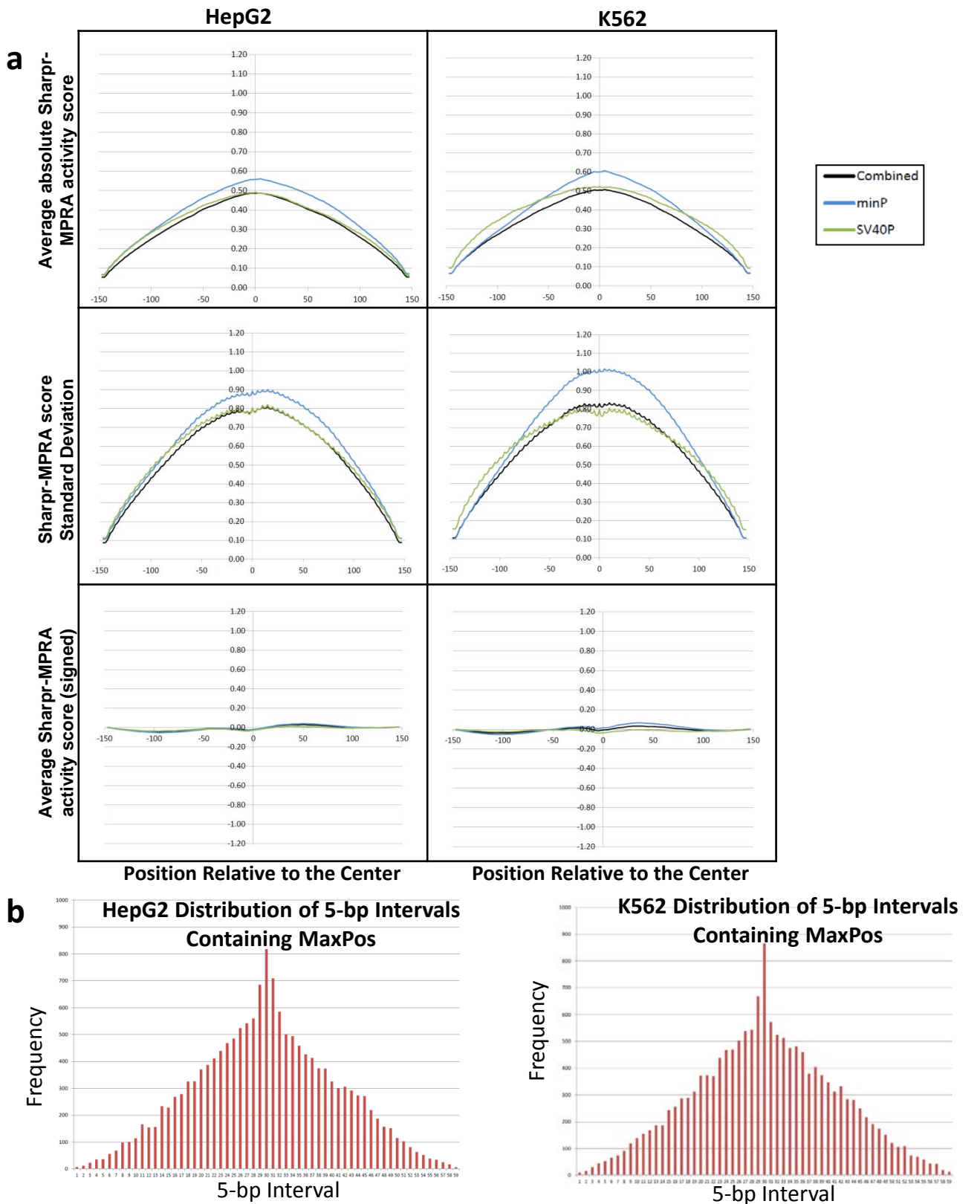
**Supplementary Figure 18(a). Comparison with DeltaSVM predictions. (a) Overlap with Nucleotide Positions predicted by DeltaSVM to contain Top 1% Negative Mutations to Reference Sequence.** Analogous plots to **Fig. 3c, Supplementary Figs. 13, 16-17**, but the enrichment is based on the 1% nucleotide positions tested that are predicted by DeltaSVM<sup>14</sup> to have a possible mutation to the reference sequence that would cause the greatest decrease in being regulatory (see Methods). Overlaps are shown for HepG2 (left) and K562 (right) for the combined SV40P and minP data (top row), minP only (middle row), and SV40P only (bottom row). An enrichment for DeltaSVM predictions is seen for nucleotides among our most activating predicted bases, but not our most repressive. **(b) Next Page.**

**b**

**Supplementary Figure 18(b). Comparison with DeltaSVM predictions. (a)** Previous page. **(b) Overlap with Nucleotides Predicted by DeltaSVM to contain Top 1% Positive Mutations to the Reference Sequence.** Analogous to (a), but based on the top 1% nucleotides having a possible mutation to the reference sequence leading to the greatest predicted increase in the sequence being regulatory (see Methods). Overlap is shown for HepG2 (left) and K562 (right), for the combined minP and SV40P data (top row), the minP data only (middle row), and the SV40P data only (bottom row). A depletion is seen in our most activating predicted bases.



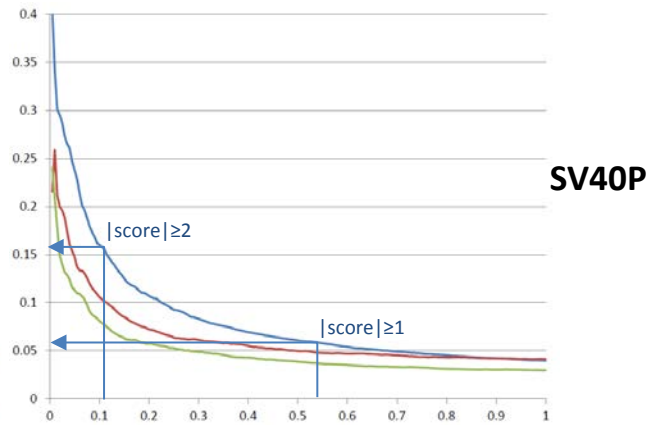
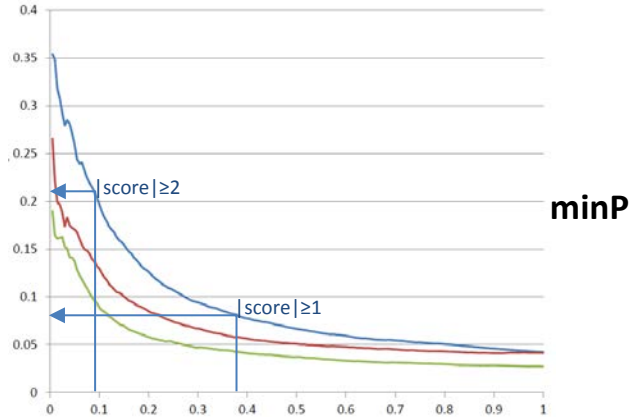
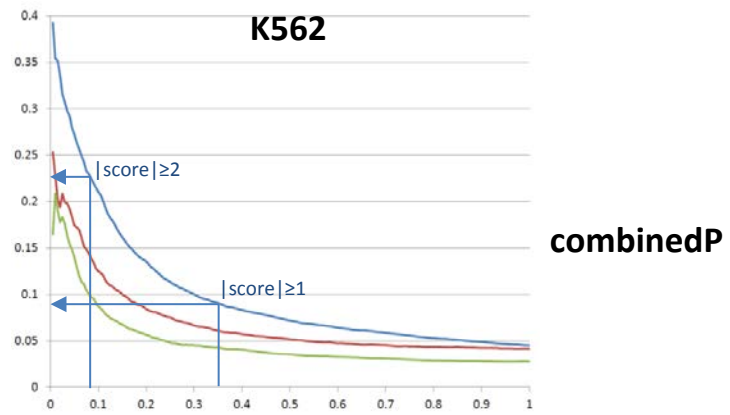
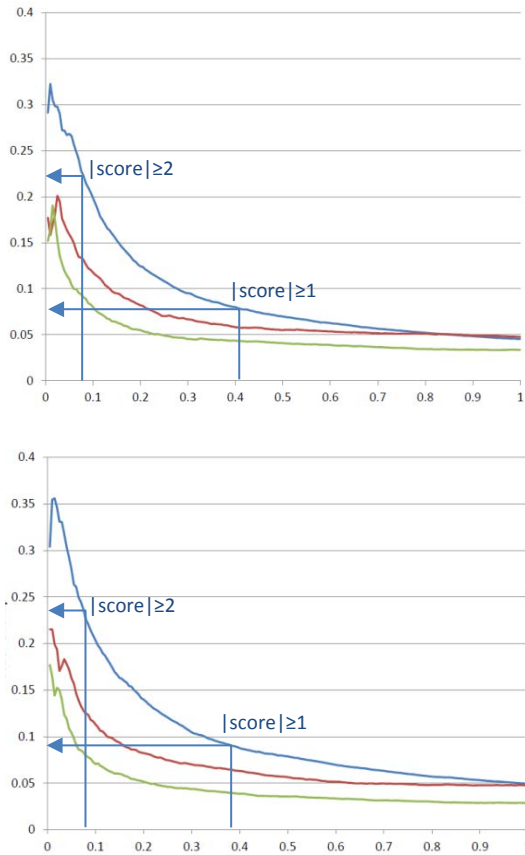
**Supplementary Figure 19 – Enrichment Based on Maximum Absolute Activity Position.** These are similar plots to those shown in **Fig. 3c** and **Supplementary Figs. 13 and 17** except the scatter plots are now based on only the single position which had the highest absolute Sharpr-MPRA regulatory activity score in each region tested (MaxPos nucleotides). The sign of the score is preserved for the analysis. In these plots there are 200 quantiles, so each point corresponds to around 79 unique MaxPos nucleotides. The green horizontal line represents the overlap fraction based on the center position and the black horizontal line overall overlap fraction of MaxPos nucleotides. The plots are **(a)** for transcription factor binding sites predicted by the CENTIPEDE method<sup>8</sup> for HepG2 (left) and K562 (right) and **(b)** for conserved elements detected by the SiPhy-PI method<sup>42,65</sup> for HepG2 (left) and K562 (right).



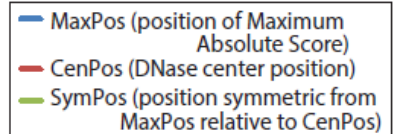
**Supplementary Figure 20 – Sharpr-MPRA activity score and standard deviation by position. (a)** Top row: Central positions show higher average absolute activity score for HepG2 (left) and K562 (right) for the minP data (blue), SV40P data (green), and combinedP data (black). Middle row: Central positions also show higher standard deviation of activity score. Bottom row: Average signed activity does not show a bias for central positions, as expected when averaging both positive and negative values. **(b)** Distribution of the location of maximum absolute regulatory activity (MaxPos) in each region among the 59 five-nucleotide intervals, for HepG2 (left) and in K562 (right).

**a****HepG2**Fig. 3d,  
left panel**K562**

Overlap Fraction with CENTIPEDE motif instances



Cumulative Fraction of Regulatory Regions



**Supplementary Figure 21 – Ranked MaxPos overlap with CENTIPEDE motifs and evolutionarily-conserved nucleotides. (a)** The three-way comparison of the cumulative overlap (y-axis) with CENTIPEDE<sup>8</sup> predicted transcription factor binding sites analogous to **Fig. 3d** (left panel) for HepG2, and shown here also for K562 for center nucleotide positions (CenPos, red), maximum-absolutely-score nucleotide positions (MaxPos, blue), and their symmetric nucleotide positions (SymPos, green), each ranked from highest to lowest based on the absolute Sharpr-MPRA scores (x-axis) for the (top) combinedP, (middle) minP, and (bottom) SV40P data indicating that our inference strategy captures regulatory nucleotides at high resolution. MaxPos, CenPos, and SymPos nucleotide positions are illustrated in the example of **Fig. 3b**. **(b)** next page.

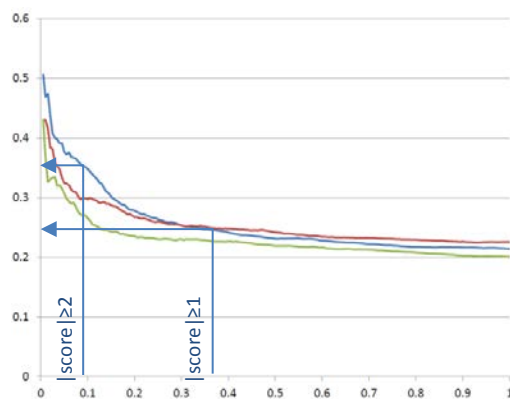
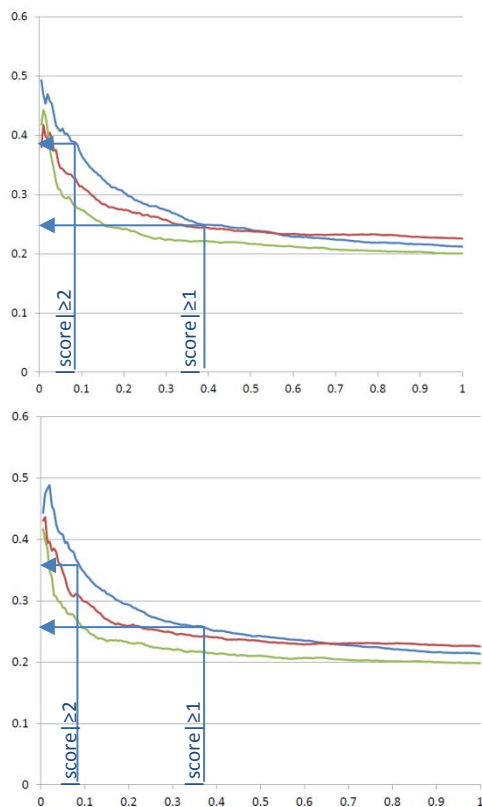
HepG2

K562

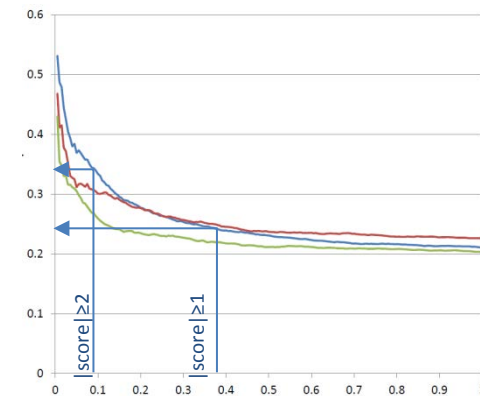
b

Fig. 3d,  
right panel

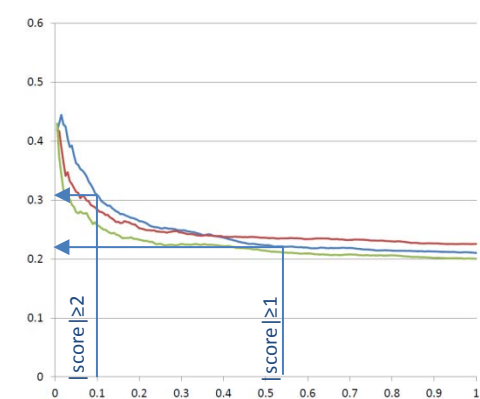
Overlap Fraction with evolutionarily-conserved elements



combinedP

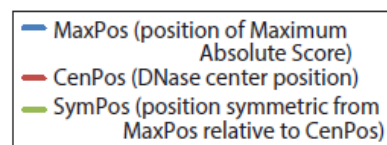


minP



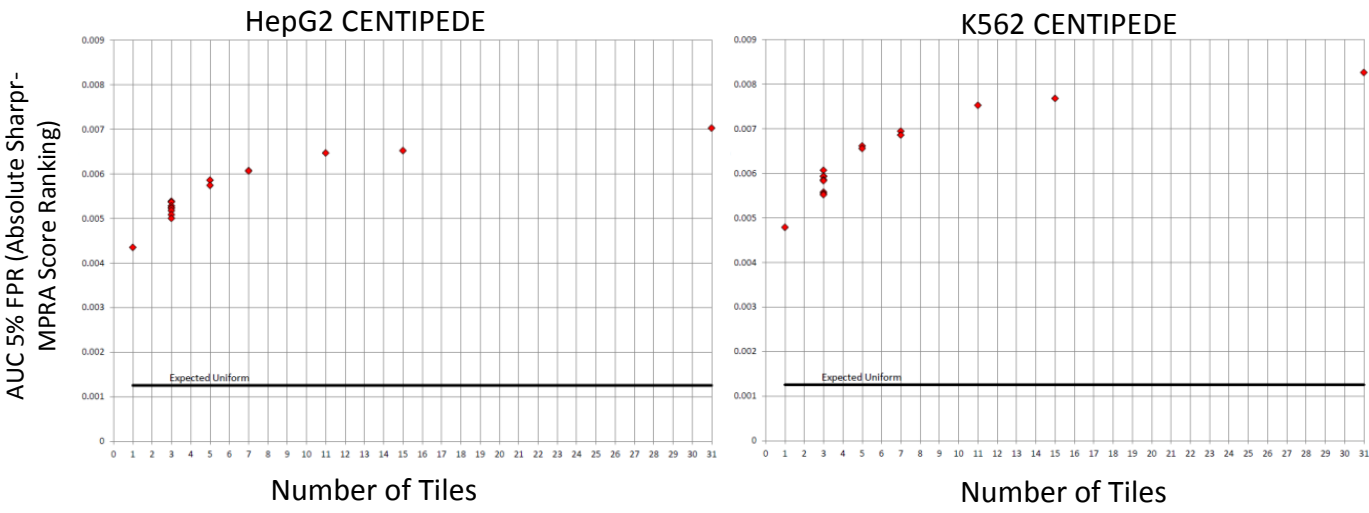
SV40P

Cumulative Fraction of Regulatory Regions

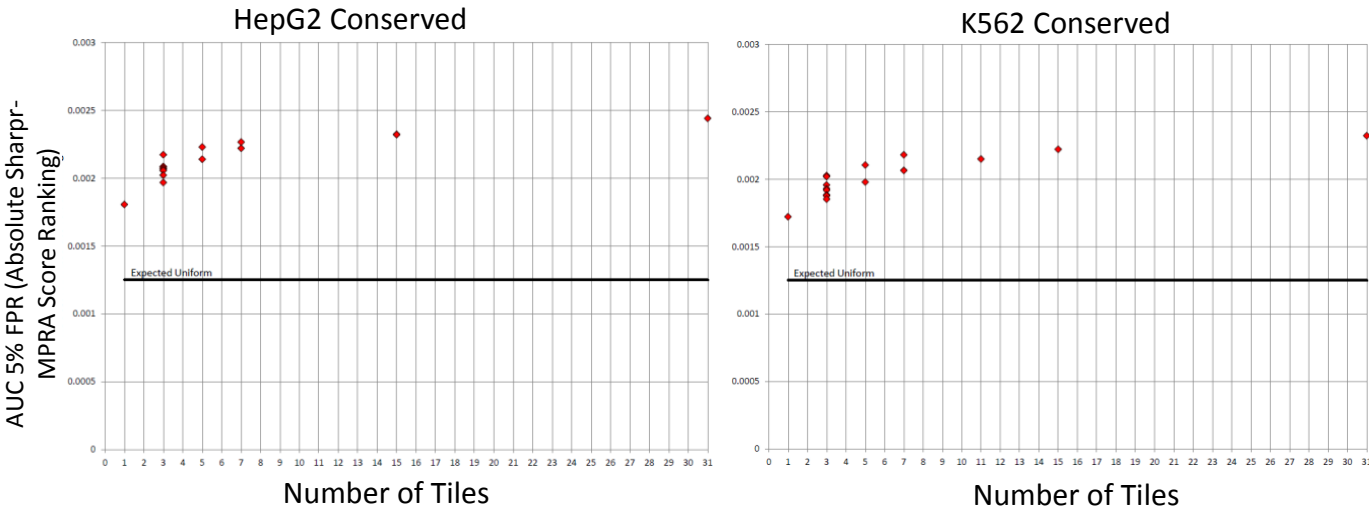


**Supplementary Figure 21 – Ranked MaxPos overlap with CENTIPEDE motifs and evolutionarily-conserved nucleotides.** (a) Previous page. (b) Overlap with conserved elements called by the SiPhy-PI method<sup>42,65</sup> analogous to Fig. 3d (right panel) for HepG2, and shown here also for K562 for center nucleotide positions (CenPos, red), maximum-absolutely-score nucleotide positions (MaxPos, blue), and their symmetric nucleotide positions (SymPos, green), each ranked from highest to lowest based on the absolute Sharpr-MPRA scores (x-axis) for the (top) combinedP, (middle) minP, and (bottom) SV40P data. MaxPos, CenPos, and SymPos nucleotide positions are illustrated in the example of Fig. 3b. The plots indicate that our inference strategy captures functional nucleotides at high resolution. The MaxPos nucleotides have higher overlap with conserved elements than CenPos nucleotides except at low absolute activity scores (Supplementary Fig. 7c).

**a. Motif recovery with varying numbers of tiles**

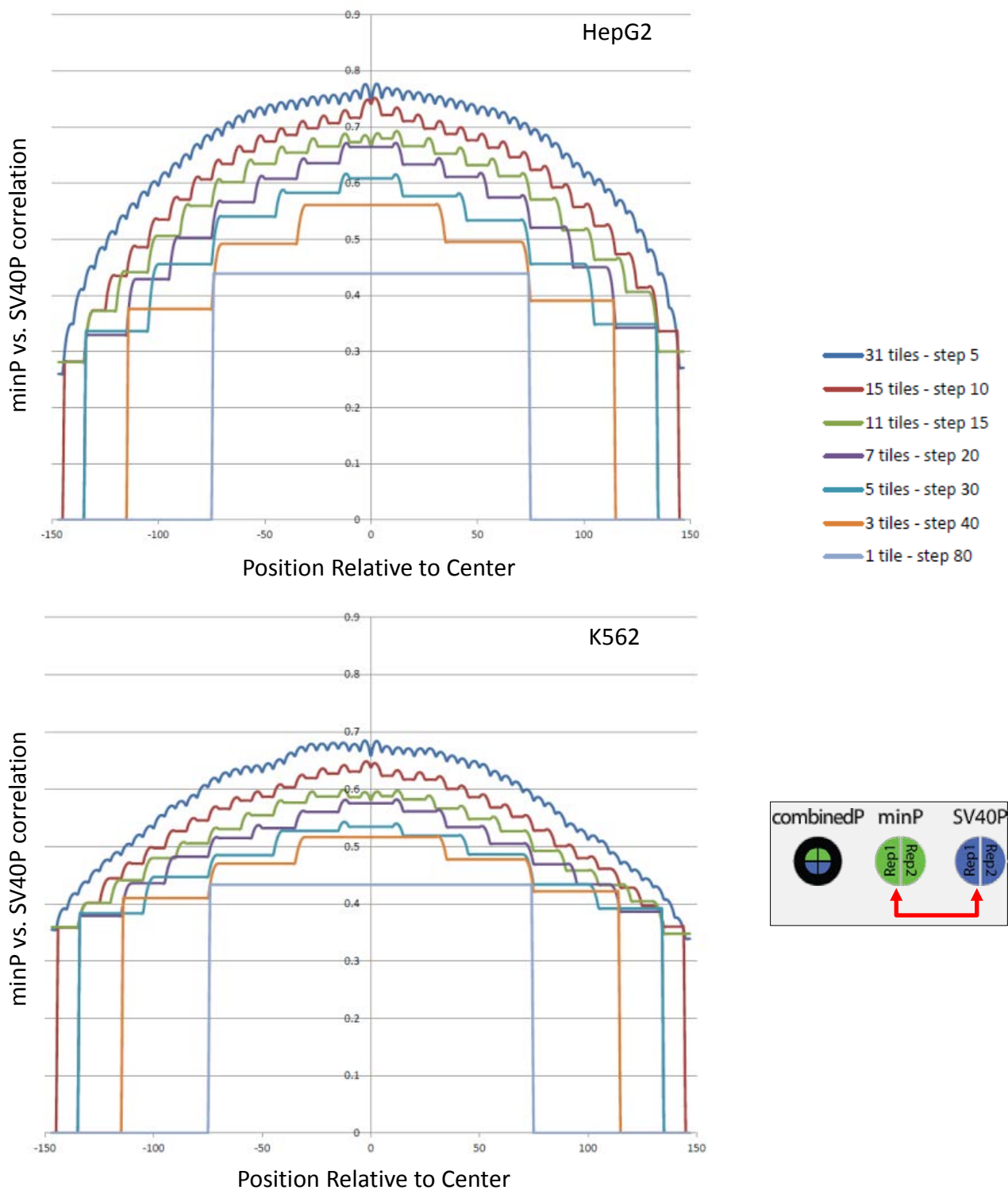


**b. Evolutionarily-conserved element recovery with varying numbers of tiles**

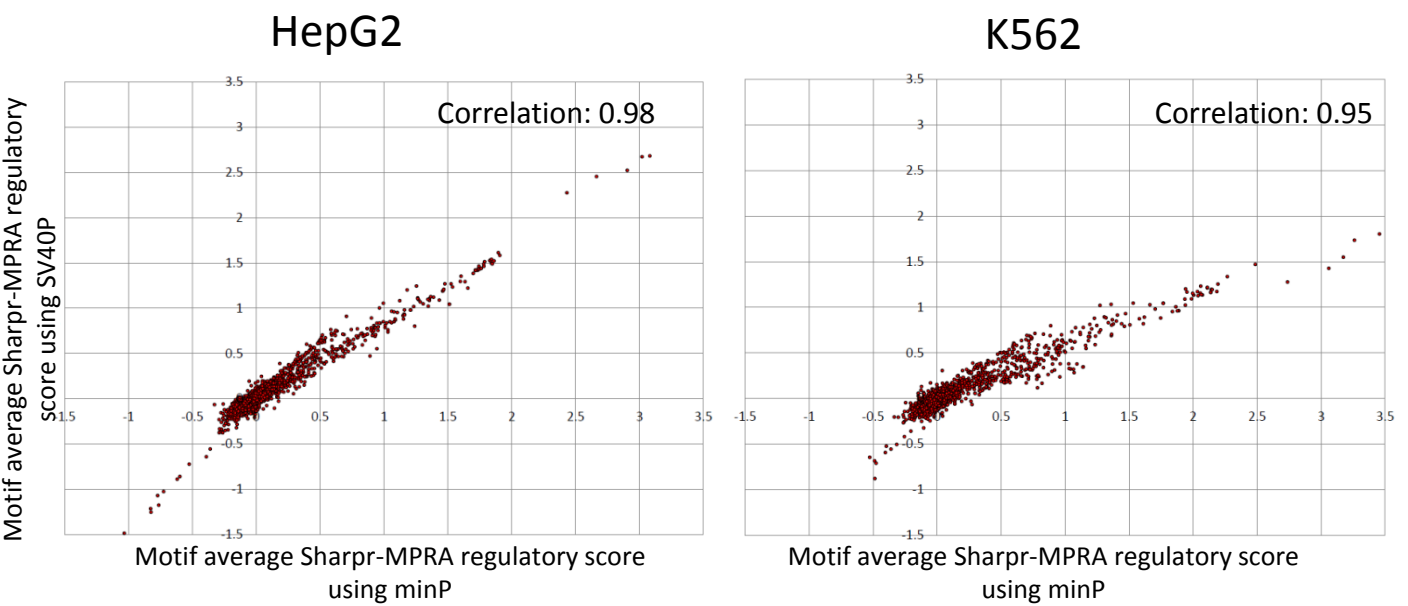


**Supplementary Figure 22 – Impact of Number of Tiles and Tiling Interval on Motif and Conserved Element Recovery.** Recovery of **(a)** CENTIPEDE<sup>8</sup> motif instances and **(b)** evolutionarily-conserved elements from the SiPhy-PI method<sup>42,65</sup> based on the Area under the ROC curve (AUC) up to a false positive rate (FPR) of 5% (y-axis), as a function of the number of tiles (x-axis), when ranking nucleotides by absolute Sharpr-MPRA regulatory activity score in HepG2 (left) and K562 (right). Multiple points for the same vertical line correspond to different step sizes leading to the same number of tiles within the 295bp region tested.





**Supplementary Figure 23 – Impact of Number of Tiles and Tiling Interval on Correlation.** Correlation between the minP and SV40P experiments (y-axis) at different positions relative to the center (x-axis) for regulatory activity inferred using only a subset of tiles selected by varying the step size (colors) for HepG2 (top) and K562 (bottom). In all cases the center tile is retained. If two or more step sizes led to the same number of tiles within the 295bp region tested, then only the correlations based on the smallest step size is plotted. We find that increasing the number of tiles (i.e. decreasing the step size) leads to increased correlation levels supporting the high-density tiling approach.



**Supplementary Figure 24 – Motif average Sharpr-MPRA regulatory score concordance between minP and SV40P data.** Scatter plot of the motif average Sharpr-MPRA regulatory obtained by averaging the activity at the central motif position for all motif instances using the minP scores (x-axis) or the SV40P scores (y-axis), for HepG2 (left) and K562 (right). Correlation between minP and SV40P motif scores is 0.98 in HepG2 and 0.95 in K562.

**a. Motifs with strongest difference in activation or in repression for HepG2 vs. K562**

b	c	d	e	f	g	h	i	k	m	-log <sub>10</sub> (j)	n	o	p	q	r	s	t	u	v	w																	
																					Summary			Average Score Analysis					Activating/Repressive Enrichment Analysis					Avg Activity in Matched			
																					AvgScore	Enriched in [score] ≥1		Avg Score		Pdiff	-log <sub>10</sub> Pdiff		Act Enr (-logP)		Repr Enr (-logP)		Enr HepG2-K562		HepG2match		K562match
SigHigh	HepG2	K562	N	HepG2	K562	Diff	corrP	corr	uncorr	HepG2	K562	HepG2	K562	ActDiff	ReprDiff	N	Avg	N	Avg																		
	TP53_7	HepG2	Activating	Neither	83	1.26	-0.10	1.36	4.7E-04	3.3	6.6	17.7	0.4	0.5	1.2	17.3	-0.7	21	2.113	17	-0.17																
	HNF4_known9	HepG2	Activating	Neither	393	0.60	-0.07	0.67	1.8E-25	24.7	28.0	41.3	0.0	0.0	1.0	41.3	-1.0	245	0.754	44	0.12																
	RXRA_disc1	HepG2	Activating	Neither	281	0.38	-0.03	0.41	2.8E-08	7.5	10.8	14.5	0.2	0.4	0.4	14.2	0.0	148	0.559	44	0.17																
	HNF1B_3	HepG2	Activating	Neither	172	0.27	-0.11	0.38	2.0E-06	5.7	9.0	4.9	0.0	0.0	0.1	4.9	-0.1	97	0.357	20	-0.05																
	PPARA_4	HepG2	Activating	Neither	323	0.37	0.00	0.37	4.1E-07	6.4	9.7	12.7	0.2	0.1	0.4	12.5	-0.4	157	0.589	46	0.15																
	NFkB_known2	HepG2	Activating	Neither	177	0.33	-0.03	0.36	4.2E-04	3.4	6.7	6.8	0.5	0.1	0.1	6.4	0.1	46	0.757	24	0.13																
	HNF1A_4	HepG2	Activating	Neither	204	0.26	-0.09	0.35	7.6E-07	6.1	9.4	5.1	0.0	0.0	0.1	5.1	0.0	111	0.363	22	-0.05																
	RXRG_3	HepG2	Activating	Neither	259	0.41	0.08	0.33	1.1E-04	4.0	7.3	12.6	1.2	0.1	0.5	11.5	-0.3	134	0.558	34	0.27																
	NFAT5_1	HepG2	Activating	Neither	236	0.27	-0.05	0.31	4.5E-05	4.3	7.6	5.9	0.1	0.3	0.2	5.8	0.1	57	0.755	47	0.04																
	RXRB_1	HepG2	Activating	Neither	252	0.40	0.09	0.31	2.3E-03	2.6	5.9	11.2	1.0	0.3	0.8	10.1	-0.5	133	0.559	37	0.25																
	HNF1_3	HepG2	Activating	Neither	271	0.16	-0.14	0.30	1.8E-05	4.7	8.0	2.5	0.0	0.3	0.1	2.5	0.1	126	0.361	41	-0.15																
	NR2C2_known1	HepG2	Activating	Neither	263	0.34	0.07	0.27	9.2E-03	2.0	5.3	8.3	0.9	0.7	1.0	7.4	-0.2	133	0.493	41	0.27																
	NR2F6_3	HepG2	Activating	Neither	253	0.33	0.07	0.26	6.2E-03	2.2	5.5	8.2	1.0	0.3	0.3	7.2	0.0	135	0.459	39	0.27																
	TCF7L1_2	HepG2	Activating	Neither	330	0.10	-0.11	0.20	8.6E-04	3.1	6.4	3.6	0.0	0.1	0.6	3.5	-0.5	152	0.319	51	0.01																
	TCF7_2	HepG2	Activating	Neither	315	0.11	-0.09	0.19	1.2E-03	2.9	6.2	4.3	0.0	0.1	0.5	4.2	-0.4	142	0.322	58	0.14																
	EP300_disc3	HepG2	Activating	Neither	541	0.06	-0.10	0.16	4.8E-05	4.3	7.6	3.2	0.0	0.3	0.3	3.2	0.0	259	0.162	66	-0.02																
	RFX5_known6	HepG2	Repressive	Neither	147	-0.34	-0.03	0.31	6.0E-05	4.2	7.5	0.3	0.8	18.6	0.9	-0.5	17.7	42	-0.32	18	0.04																
	RFX2_1	HepG2	Repressive	Neither	140	-0.33	-0.03	0.30	2.1E-04	3.7	7.0	0.3	0.9	16.8	1.0	-0.5	15.9	37	-0.341	17	0.1																
	RFX3_2	HepG2	Repressive	Neither	158	-0.31	-0.02	0.29	6.2E-05	4.2	7.5	0.5	1.1	17.8	0.8	-0.6	17.0	43	-0.288	20	0.12																
	BCL6B_1	HepG2	Repressive	Neither	131	-0.17	0.09	0.26	1.7E-03	2.8	6.1	0.0	0.6	2.5	0.4	-0.5	2.1	22	-0.095	22	0.18																
	MX1_disc1	HepG2	Repressive	Neither	280	-0.28	-0.02	0.26	3.1E-07	6.5	9.8	0.4	0.8	30.5	1.8	-0.3	28.7	67	-0.252	38	0.27																
	BCL_disc3	HepG2	Repressive	Repressive	136	-1.02	-0.58	0.43	3.9E-04	3.4	6.7	0.2	0.1	43.8	39.9	0.2	3.9	42	-0.618	28	-0.78																
	REST_disc1	HepG2	Repressive	Repressive	216	-0.92	-0.49	0.43	1.2E-08	7.9	11.2	0.0	0.0	59.7	47.6	0.0	12.1	45	-0.892	49	-0.48																
	SIN3A_disc1	HepG2	Repressive	Repressive	217	-0.88	-0.46	0.42	1.0E-08	8.0	11.3	0.2	0.2	59.6	44.9	0.0	14.7	46	-0.772	54	-0.41																
	SREBP_disc1	HepG2	Repressive	Repressive	254	-0.22	0.00	0.22	1.6E-02	1.8	5.1	1.1	1.7	27.9	2.1	-0.6	25.7	64	-0.203	40	0.13																
	MYC_disc2	HepG2	Repressive	Repressive	234	-0.31	-0.12	0.19	7.5E-03	2.1	5.4	0.0	0.1	22.6	2.9	0.0	19.8	56	-0.196	35	-0.01																
	FOXA_known2	HepG2	Neither	Neither	603	0.03	-0.12	0.15	1.2E-05	4.9	8.2	1.4	0.0	0.1	0.3	1.4	-0.2	260	0.128	79	-0.06																
	GATA_known14	K562	Neither	Activating	284	-0.12	0.36	0.48	3.6E-11	10.4	13.7	0.1	13.6	1.2	0.2	-13.5	1.0	39	0.041	152	0.39																
	TAL1_disc1	K562	Neither	Activating	293	-0.14	0.34	0.47	7.7E-12	11.1	14.4	0.2	12.5	0.9	0.2	-12.3	0.6	41	0.003	163	0.37																
	CCNT2_disc1	K562	Neither	Activating	313	-0.11	0.33	0.43	3.0E-10	9.5	12.8	0.3	9.8	1.0	0.5	-9.6	0.6	53	0.103	158	0.36																
	STAT_disc1	K562	Neither	Activating	279	0.01	0.36	0.35	4.9E-07	6.3	9.6	1.1	9.8	0.6	0.1	-8.7	0.6	57	0.143	96	0.69																
	KLF13_1	K562	Neither	Activating	142	0.18	0.43	0.26	2.6E-03	2.6	5.9	1.5	7.2	0.5	0.0	-5.7	0.4	36	0.239	33	0.48																
	KLF12_2	K562	Activating	Activating	230	0.19	0.52	0.32	3.3E-09	8.5	11.8	4.3	19.6	0.9	0.1	-15.4	0.8	52	0.035	65	0.64																
	SP4_2	K562	Activating	Activating	534	0.36	0.67	0.31	7.3E-21	20.1	23.4	24.0	68.7	1.0	0.0	-44.7	1.0	134	0.345	137	0.82																
	KLF14_1	K562	Activating	Activating	326	0.29	0.60	0.30	7.6E-12	11.1	14.4	10.3	37.7	1.3	0.2	-27.4	1.1	87	0.181	86	0.82																
	IRF_disc4	K562	Activating	Activating	879	0.35	0.64	0.30	1.6E-33	32.8	36.1	33.4	104.0	0.7	0.0	-70.6	0.7	228	0.237	231	0.71																
	KLF7_1	K562	Activating	Activating	534	0.24	0.52	0.28	2.6E-19	18.6	21.9	15.0	47.8	1.9	0.1	-32.9	1.7	131	0.162	142	0.68																
	KLF4_1	K562	Activating	Activating	638	0.22	0.48	0.26	5.6E-22	21.3	24.5	12.9	50.1	0.8	0.1	-37.1	0.7	171	0.122	170	0.63																
	SP1_known4	K562	Activating	Activating	904	0.32	0.56	0.24	5.0E-26	25.3	28.6	31.5	84.7	0.9	0.0	-53.2	0.9	244	0.193	231	0.65																
	TATA_disc4	K562	Activating	Activating	747	0.37	0.59	0.22	1.1E-15	15.0	18.3	30.9	73.6	0.9	0.1	-42.6	0.8	204	0.257	194	0.71																
	PAX5_disc5	K562	Activating	Activating	152	0.47	0.67	0.20	1.6E-02	1.8	5.1	16.4	22.0	0.7	0.3	-5.5	0.5	43	0.296	35	0.81																
	KLF16_1	K562	Activating	Activating	648	0.21	0.41	0.20	6.4E-14	13.2	16.5	9.6	39.3	0.6	0.3	-29.7	0.3	161	0.114	155	0.56																
	TRIM28_disc1	K562	Activating	Activating	652	0.51	0.71	0.20	3.4E-04	3.5	6.8	66.8	88.4	0.4	0.0	-21.6	0.4	167	0.9	152	0.95																
	BRACH2_1	K562	Activating	Activating	418	0.51	0.71	0.20	2.2E-02	1.6	4.9	35.0	51.5	0.0	0.0	-16.6	0.0	117	0.78	94	0.86																
	NRF1_disc3	K562	Activating	Activating	453	0.18	0.37	0.19	8.9E-10	9.0	12.3	5.6	26.0	1.1	0.2	-20.4	0.9	137	0.069	104	0.41																
	ATF3_disc4	K562	Activating	Activating	433	0.27	0.45	0.18	1.8E-07	6.8	10.0	12.5	29.5	1.2	0.0	-17.0	1.2	117	0.096	102	0.54																
	AP1_disc3	K562	Activating	Activating	773	0.42	0.60	0.17	2.1E-03	2.7	6.0	56.0	83.2	0.0	0.0	-27.2	0.0	194	0.787	150	0.84																
	TCF7L2_disc1	K562	Activating	Activating	819	0.58	0.75	0.17	1.7E-03	2.8	6.0	90.8	111.1	0.0	0.0	-20.3	0.0	213	0.944	169	0.99																
	PRDM1_disc2	K562	Activating	Activating	816	0.44	0.61	0.17	1.1E-03	3.0	6.2	65.2	98.1	0.0	0.0	-32.9	0.0	204	0.776	161	0.92																
	SP8_1	K562	Activating	Activating	461	0.22	0.38	0.16	8.5E-05	4.1	7.4	9.0	26.1	0.4	0.2	-17.0	0.2	120	0.12	108	0.54																
	TFAP2_disc1	K562	Activating	Activating	830	0.69	0.85	0.16	1.1E-02	1.9	5.2	113.5	128.2	0.0	0.0	-14.8	0.0	208	1.108	199	0.99																
	E2F_disc7	K562	Activating	Activating	504	0.27	0.42	0.15	3.9E-04	3.4	6.7	13.1	35.4	1.4	0.1	-22.2	1.3	132	0.112	124	0.55																
	CHD2_disc3	K562	Activating	Activating</																																	

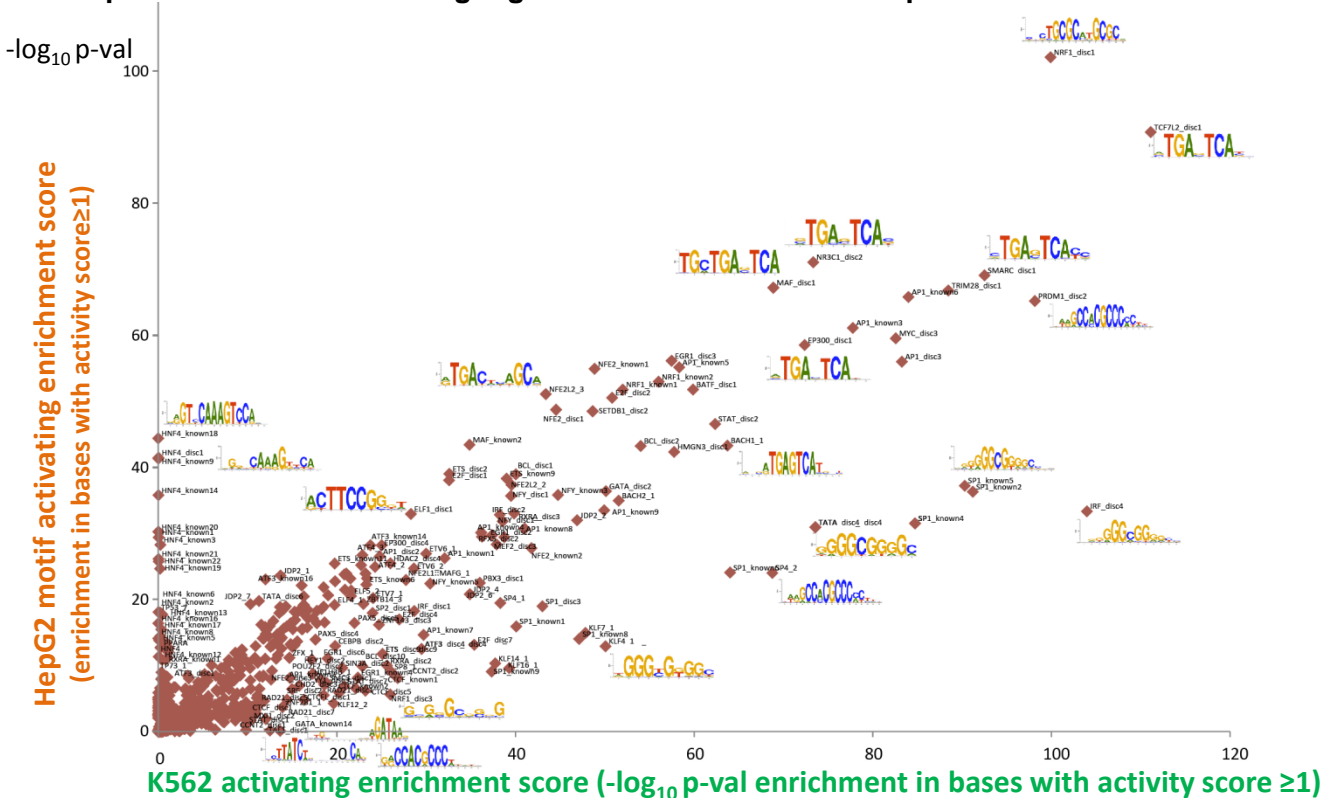
**Motifs with strongest activation or repression enrichments in either HepG2 or in K562 (page 1 of 3)**

b	af	d	e	n	o	p	q	r	s	f	g	h	i	m	c	t	u	v	w	For		Summary								Activating/Repressive Enrichment Analysis								Average Score Analysis				Avg Activity in Matched			
																				Sorting		Enrichment in  score ≥1		Act Enr (-logP)		Repr Enr (-logP)		Enr HepG2-K562		Avg Score		Pdiff		HepG2match		K562match									
																				MaxAct	HepG2	K562	HepG2	K562	HepG2	K562	ActDiff	ReprDiff	N	HepG2	K562	Diff	-lg10Cor	Higher	N	Avg	N	Avg							
TGA_TCA	TFAP2_disc1	128.24	Activating	Activating	113.5	128.2	0.0	0.0	-14.8	0.0	830	0.69	0.85	0.16	1.9	K562	208	1.108	199	0.99																									
TGA_TCA	TCF7L2_disc1	111.11	Activating	Activating	90.8	111.1	0.0	0.0	-20.3	0.0	819	0.58	0.75	0.17	2.8	K562	213	0.944	169	0.99																									
TGA_TCA	IRF_disc4	103.96	Activating	Activating	33.4	104.0	0.7	0.0	-70.6	0.7	879	0.35	0.64	0.30	32.8	K562	228	0.237	231	0.71																									
TGA_TCA	NRF1_disc1	102.12	Activating	Activating	102.1	99.9	0.2	0.0	2.2	0.1	249	1.35	1.44	0.09	0.0	N/S	46	1.196	53	1.46																									
TGA_TCA	PRDM1_disc2	98.14	Activating	Activating	65.2	98.1	0.0	0.0	-32.9	0.0	816	0.44	0.61	0.17	3.0	K562	204	0.776	161	0.92																									
TGA_TCA	SMARCD1_disc1	92.50	Activating	Activating	69.1	92.5	0.0	0.0	-23.4	0.0	729	0.58	0.71	0.14	0.6	N/S	175	0.88	164	0.89																									
TGA_TCA	SP1_known2	91.19	Activating	Activating	36.3	91.2	0.8	0.0	-54.9	0.8	1070	0.32	0.54	0.22	23.8	K562	273	0.233	270	0.65																									
TGA_TCA	TRIM28_disc1	88.44	Activating	Activating	66.8	88.4	0.4	0.0	-21.6	0.4	652	0.51	0.71	0.20	3.5	K562	167	0.9	152	0.95																									
TGA_TCA	AP1_known6	83.98	Activating	Activating	65.8	84.0	0.1	0.0	-18.2	0.1	637	0.61	0.68	0.06	0.0	N/S	150	1.063	140	0.82																									
TGA_TCA	MYC_disc3	82.56	Activating	Activating	59.6	82.6	0.3	0.0	-23.0	0.3	748	0.45	0.58	0.13	0.4	N/S	199	0.878	162	0.76																									
TGA_TCA	TATA_disc4	73.56	Activating	Activating	30.9	73.6	0.9	0.1	-42.6	0.8	747	0.37	0.59	0.22	15.0	K562	204	0.257	194	0.71																									
TGA_TCA	NR3C1_disc2	73.32	Activating	Activating	71.1	73.3	0.0	0.0	-2.3	0.0	717	0.56	0.66	0.10	0.0	N/S	145	0.748	174	0.77																									
TGA_TCA	EP300_disc1	72.36	Activating	Activating	58.5	72.4	0.3	0.0	-13.8	0.3	751	0.45	0.56	0.11	0.0	N/S	205	0.888	152	0.83																									
TGA_TCA	MAF_disc1	68.86	Activating	Activating	67.2	68.9	0.1	0.0	-1.6	0.0	447	0.68	0.78	0.10	0.0	N/S	112	1.094	119	0.96																									
TGA_TCA	SP4_2	68.75	Activating	Activating	24.0	68.7	1.0	0.0	-44.7	1.0	534	0.36	0.67	0.31	20.1	K562	134	0.345	137	0.82																									
TGA_TCA	BACH1_1	63.72	Activating	Activating	43.3	63.7	0.1	0.0	-20.4	0.1	618	0.47	0.57	0.10	0.0	N/S	161	0.846	131	0.75																									
TGA_TCA	REST_known2	63.63	Repressive	Repressive	0.3	0.2	63.6	52.4	0.0	11.2	132	-1.26	-0.68	0.58	6.2	HepG2	36	-0.894	32	-0.73																									
TGA_TCA	STAT_disc2	62.34	Activating	Activating	46.6	62.3	0.2	0.0	-15.8	0.2	698	0.36	0.50	0.15	1.2	N/S	189	0.727	146	0.84																									
TGA_TCA	BATF_disc1	59.87	Activating	Activating	51.8	59.9	0.1	0.0	-8.1	0.1	624	0.50	0.61	0.11	0.0	N/S	135	0.954	122	0.87																									
TGA_TCA	SIN3A_disc1	59.59	Repressive	Repressive	0.2	0.2	59.6	44.9	0.0	14.7	217	-0.88	-0.46	0.42	8.0	HepG2	46	-0.772	54	-0.41																									
TGA_TCA	HMGN3_disc1	57.75	Activating	Activating	42.4	57.8	0.1	0.3	-15.4	-0.1	493	0.59	0.70	0.11	0.0	N/S	124	0.911	114	0.85																									
TGA_TCA	EGR1_disc3	57.47	Activating	Activating	56.2	57.0	0.2	0.0	-1.3	0.2	153	1.26	1.40	0.14	0.0	N/S	30	0.914	40	1.45																									
TGA_TCA	NFE2_known1	54.94	Activating	Activating	54.9	48.8	0.6	0.0	6.1	0.6	458	0.56	0.62	0.06	0.0	N/S	103	1.01	118	0.76																									
TGA_TCA	BCL_disc2	54.00	Activating	Activating	43.2	54.0	0.5	0.0	-10.8	0.5	715	0.34	0.43	0.09	0.0	N/S	189	0.728	147	0.73																									
TGA_TCA	BACH2_1	51.54	Activating	Activating	35.0	51.5	0.0	0.0	-16.6	0.0	418	0.51	0.71	0.20	1.6	K562	117	0.78	94	0.86																									
TGA_TCA	NFE2L2_3	51.11	Activating	Activating	51.1	43.4	0.1	0.1	7.7	-0.1	391	0.65	0.64	0.01	0.0	N/S	83	1.187	125	0.61																									
TGA_TCA	E2F_disc2	50.81	Activating	Activating	50.6	50.8	0.3	0.1	-0.3	0.2	121	1.33	1.43	0.10	0.0	N/S	23	1.137	27	1.55																									
TGA_TCA	GATA_disc2	50.12	Activating	Activating	36.4	50.1	0.1	0.0	-13.7	0.0	510	0.38	0.49	0.10	0.0	N/S	106	0.782	109	0.73																									
TGA_TCA	KLF4_1	50.07	Activating	Activating	12.9	50.1	0.8	0.1	-37.1	0.7	638	0.22	0.48	0.26	21.3	K562	171	0.122	170	0.63																									
TGA_TCA	HDAC2_disc3	49.50	Repressive	Repressive	0.5	0.6	49.5	40.9	-0.1	8.6	243	-0.63	-0.33	0.29	2.5	HepG2	68	-0.31	60	-0.43																									
TGA_TCA	SETDB1_disc2	48.63	Activating	Activating	48.5	48.6	0.2	0.1	-0.1	0.1	143	1.19	1.32	0.13	0.0	N/S	32	1.018	30	1.27																									
TGA_TCA	KLF7_1	47.82	Activating	Activating	15.0	47.8	1.9	0.1	-32.9	1.7	534	0.24	0.52	0.28	18.6	K562	131	0.162	142	0.68																									
TGA_TCA	JDP2_2	46.89	Activating	Activating	32.0	46.9	0.1	0.0	-14.9	0.1	502	0.41	0.50	0.10	0.0	N/S	126	0.78	102	0.83																									
TGA_TCA	NFY_known3	44.75	Activating	Activating	35.8	44.7	0.7	0.0	-8.9	0.6	218	0.61	0.75	0.14	0.0	N/S	44	0.592	45	0.66																									
TGA_TCA	HNF4_known18	44.41	Activating	Neither	44.4	0.0	0.1	0.8	44.4	-0.8	241	0.82	-0.14	0.97	20.6	HepG2	166	0.971	25	-0.05																									
TGA_TCA	RXRα_disc3	39.83	Activating	Activating	33.0	39.8	0.5	0.0	-6.8	0.5	566	0.32	0.41	0.09	0.0	N/S	153	0.708	112	0.72																									
TGA_TCA	KLF16_1	39.29	Activating	Activating	9.6	39.3	0.6	0.3	-29.7	0.3	648	0.21	0.41	0.20	13.2	K562	161	0.114	155	0.56																									
TGA_TCA	ETS_disc2	39.03	Activating	Activating	39.0	32.6	0.2	0.1	6.5	0.0	74	1.70	1.81	0.11	0.0	N/S	15	2.229	11	2.02																									
TGA_TCA	MZF2_disc3	37.88	Activating	Activating	28.3	37.9	0.1	0.0	-9.5	0.1	620	0.38	0.42	0.04	0.0	N/S	154	0.789	125	0.45																									
TGA_TCA	KLF14_1	37.72	Activating	Activating	10.3	37.7	1.3	0.2	-27.4	1.1	326	0.29	0.60	0.30	11.1	K562	87	0.181	86	0.82																									
TGA_TCA	RFX5_disc2	36.24	Activating	Activating	29.5	36.2	0.5	0.2	-6.7	0.3	201	0.54	0.68	0.15	0.0	N/S	45	0.342	46	0.69																									
TGA_TCA	PBX3_disc1	35.99	Activating	Activating	22.5	36.0	0.4	0.0	-13.5	0.4	216	0.47	0.62	0.15	0.5	N/S	45	0.475	50	0.48																									
TGA_TCA	ELF1_disc1	32.96	Activating	Activating	33.0	28.3	0.2	0.1	4.7	0.0	75	1.48	1.57	0.10	0.0	N/S	15	1.859	13	1.34																									
TGA_TCA	MXI1_disc1	30.52	Repressive	Neither	0.4	0.8	30.5	1.8	-0.3	28.7	280	-0.28	-0.02	0.26	6.5	HepG2	67	-0.252	38	0.17																									
TGA_TCA	ETV6_1	29.98	Activating	Activating	26.9	30.0	0.1	0.2	-3.0	-0.1	170	0.78	0.97	0.19	0.0	N/S	26	1.192	51	0.91																									
TGA_TCA	ATF3_disc4	29.47	Activating	Activating	12.5	29.5	1.2	0.0	-17.0	1.2	433	0.27	0.45	0.18	6.8	K562	117	0.906	102	0.54																									
TGA_TCA	SREBP_disc1	27.87	Repressive	Repressive	1.1	1.7	27.9	2.1	-0.6	25.7	254	-0.22	0.00	0.22	1.8	HepG2	64	-0.203	40	0.13																									
TGA_TCA	NFE2L1::MAFG_1	27.73	Activating	Activating	23.0	27.7	0.1	0.1	-4.8	0.0	254	0.50	0.58	0.08	0.0	N/S	59	0.784	93	0.76																									
TGA_TCA	CONT2_disc2	26.81	Activating	Activating	8.1	26.8	1.3	0.0	-18.7	1.3	1248	0.13	0.23	0.10	5.9	K562	320	0.116	284	0.42																									
TGA_TCA	ATF4_3	26.73	Activating	Activating	26.7	22.9	0.3	0.2	3.9	0.1	112	1.10	1.16	0.06	0.0	N/S	37	1.319	32	1.22																									
TGA_TCA	SP8_1	26.09	Activating	Activating	9.0	26.1	0.4	0.2	-17.0	0.2	461	0.22	0.38	0.16	4.1	K562	120	0.12	108	0.54																									
TGA_TCA	ZNF143_disc3	24.70	Activating	Activating	16.2	24.7	0.9	1.3	-8.5	-0.4	554	0.30	0.34	0.04	0.0	N/S	146	0.381	105	0.42																									
TGA_TCA	ETV7_1	24.04	Activating	Activating	20.2	24.0	0.0	0.0	-3.8	0.0	250	0.59	0.57	0.03	0.0	N/S	46	0.629	61	0.73																									
TGA_TCA	SP2_disc1	23.99	Activating	Activating	18.0	24.0	0.5	0.1	-6.0	0.5	187	0.46	0.56	0.10	0.0	N/S	45	0.507	42	0.57																									
TGA_TCA	NR2C2_disc1	23.79	Activating	Activating	23.8	20.8	0.2	0.4	3.0	-0.2	70	1.25	1.34	0.10	0.0	N/S	19	1.209	13	1.25																									
TGA_TCA	ELF5_3	23.36	Activating	Activating	23.2	23.4	0.0	0.1	-0.2	-0.1	170	0.76	0.81	0.05	0.0	N/S	37	0.848	35	0.7																									
TGA_TCA	CTCF_disc5	23.19	Activating	Activating	6.2	23.2	0.6	0.4	-17.0	0.2	1084	0.16	0.22	0.06	0.0	N/S	304	0.097	213	0.33																									
TGA_TCA	ZBTB14_3	22.98	Activating	Activating	19.3	23.0	0.																																						

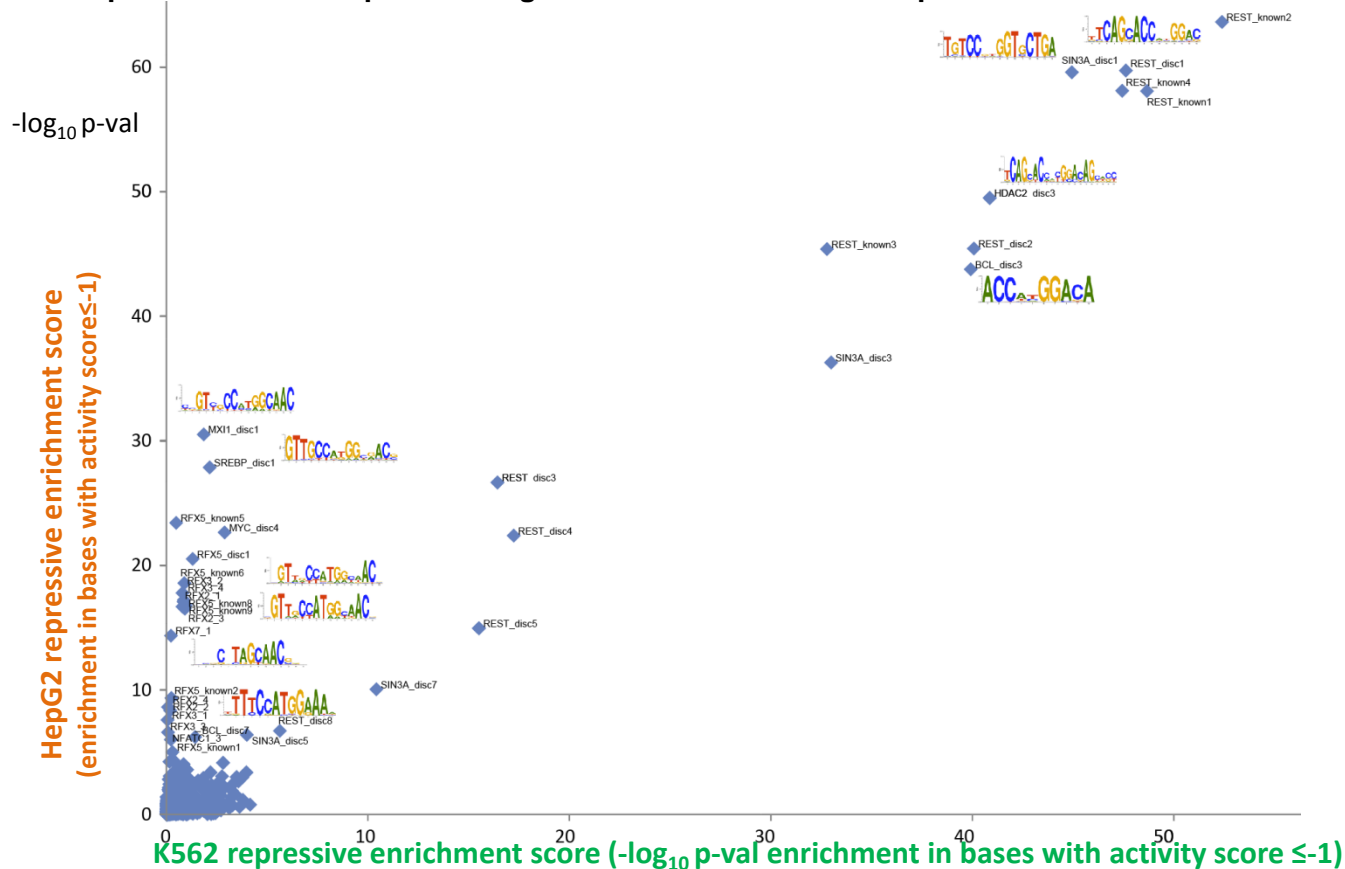
b		af d e n o p q r s f g h i m c t u v w																		
Motif	For Sorting	Summary		Activating/Repressive Enrichment Analysis						Average Score Analysis				Avg Activity in Matched						
		Enrichment in  score ≥1		Act Enr (-logP)		Repr Enr (-logP)		Enr HepG2-K562		Avg Score		Pdiff		HepG2match		K562match				
		HepG2	K562	HepG2	K562	HepG2	K562	ActDiff	ReprDiff	N	HepG2	K562	Diff	-lg10Cor	Higher	N	Avg	N	Avg	
CGTCA	XBP1_3	14.81	Activating	Activating	14.8	12.0	0.4	0.3	2.8	0.1	40	1.24	1.12	0.12	0.0	N/S	6	0.922	11	0.56
CCGCG	TFE3_1	14.75	Activating	Activating	14.8	13.3	0.2	0.1	1.4	0.1	69	0.87	0.78	0.09	0.0	N/S	22	0.943	6	1.73
CGCGCG	POU2F2_disc2	14.73	Activating	Activating	8.1	14.7	0.5	0.3	-6.6	0.2	962	0.13	0.18	0.06	0.0	N/S	239	0.119	195	0.21
CGCGT	ZFX_1	14.67	Activating	Activating	11.2	14.7	0.0	0.5	-3.5	-0.5	416	0.31	0.28	0.03	0.0	N/S	114	0.286	82	0.46
CGCGT	ELF3_3	14.48	Activating	Activating	14.5	12.3	0.1	0.5	-2.2	-0.4	114	0.73	0.75	0.02	0.0	N/S	28	0.812	21	0.77
CGCGT	SRF_disc2	14.42	Activating	Activating	4.8	14.4	0.7	0.1	-9.6	0.6	502	0.15	0.26	0.11	1.9	K562	147	0.063	117	0.34
CGTACAC	RFX7_1	14.36	Repressive	Neither	1.2	1.1	14.4	0.2	0.1	14.1	111	-0.21	0.06	0.28	1.2	N/S	26	0.379	18	0.17
ATGACGTCA	CREB5_1	14.09	Activating	Activating	14.1	11.8	0.2	0.1	2.3	0.0	80	0.85	0.76	0.09	0.0	N/S	21	0.608	17	0.5
CGCGCG	ZNF281_1	13.92	Activating	Activating	3.7	13.9	0.0	0.0	-10.2	0.0	1198	0.12	0.19	0.06	0.9	N/S	317	0.172	242	0.31
CGCGA	ETV5_1	13.64	Activating	Activating	13.6	13.3	0.3	0.7	0.3	0.4	45	1.55	1.56	0.01	0.0	N/S	6	1.455	11	1.19
CGCGA	ETV2_1	12.96	Activating	Activating	13.0	10.9	0.3	0.2	2.1	0.1	112	0.78	0.66	0.11	0.0	N/S	15	0.732	16	1.08
CGCGT	PPARA_4	12.72	Activating	Neither	12.7	0.2	0.1	0.4	12.5	-0.4	323	0.37	0.00	0.37	6.4	HepG2	157	0.589	46	0.15
CGCGT	RXRG_3	12.62	Activating	Neither	12.6	1.2	0.1	0.5	11.5	-0.3	259	0.41	0.08	0.33	4.0	HepG2	134	0.558	34	0.27
CGCGT	TAL1_disc1	12.51	Neither	Activating	0.2	12.5	0.9	0.2	-12.3	0.6	293	-0.14	0.34	0.47	11.1	K562	41	0.003	163	0.37
CGCGT	SIRT6_disc1	11.93	Activating	Activating	11.9	8.6	0.3	0.6	3.3	-0.3	54	0.98	0.69	0.29	0.0	N/S	16	0.969	7	1.15
CGCGT	SPI1_disc1	11.67	Activating	Activating	8.3	11.7	0.1	0.1	-3.4	0.0	276	0.26	0.32	0.07	0.0	N/S	43	0.141	94	0.37
CGCGT	ARNTL_1	11.61	Activating	Activating	11.6	9.3	0.3	0.2	2.3	0.1	56	0.91	0.73	0.19	0.0	N/S	17	0.804	9	1.43
CGCGT	YY2_2	11.52	Activating	Activating	11.5	9.1	0.3	0.7	2.4	-0.4	44	0.87	1.06	0.19	0.0	N/S	10	0.696	9	0.96
CGCGT	RXR8_1	11.17	Activating	Neither	11.2	1.0	0.3	0.8	10.1	-0.5	252	0.40	0.09	0.31	2.6	HepG2	133	0.559	37	0.25
CGCGT	BHLHE40_disc1	11.14	Activating	Activating	11.1	10.8	0.2	0.2	0.3	0.1	66	0.78	0.73	0.05	0.0	N/S	20	0.65	8	1.41
CGCGT	TP73_1	10.87	Activating	Neither	10.9	0.4	0.6	0.9	10.5	-0.3	68	1.03	-0.08	1.11	1.2	N/S	18	2.158	14	-0.18
CGCGT	TFEC_1	10.78	Activating	Activating	10.8	8.8	0.1	1.2	2.0	-1.0	84	0.69	0.53	0.16	0.0	N/S	26	0.833	10	0.67
CGCGT	ATF2_2	10.40	Activating	Activating	10.4	9.0	0.4	0.3	1.4	0.1	38	0.75	1.03	0.28	0.0	N/S	6	0.458	13	0.77
CGCGT	ZBTB7A_disc2	9.82	Activating	Activating	4.4	9.8	0.4	0.3	-5.4	0.1	98	0.29	0.43	0.14	0.0	N/S	16	0.24	25	0.57
CGCGT	BHLHE41_2	9.49	Activating	Activating	9.5	9.3	0.4	0.4	0.2	0.1	31	1.06	0.93	0.13	0.0	N/S	10	0.921	2	2.78
CGCGT	ZBTB33_disc1	9.31	Activating	Activating	9.3	9.2	0.8	1.7	0.2	-0.9	13	2.71	2.35	0.36	0.0	N/S	4	1.922	3	2.62
CGCGT	WT1_1	8.46	Activating	Activating	2.4	8.5	0.0	0.0	-6.1	0.0	513	0.13	0.23	0.09	1.0	N/S	137	0.188	108	0.26
CGCGT	ELF2_1	8.40	Activating	Activating	5.1	8.4	0.3	0.4	-3.3	-0.1	372	0.17	0.16	0.01	0.0	N/S	73	0.227	65	0.59
CGCGT	SPDEF_5	8.33	Activating	Activating	8.3	8.1	1.1	0.4	0.2	0.8	31	1.02	1.00	0.03	0.0	N/S	6	1.169	7	1.39
CGCGT	NR2F6_3	8.25	Activating	Neither	8.2	1.0	0.3	0.3	7.2	0.0	253	0.33	0.07	0.26	2.2	HepG2	135	0.459	39	0.27
CGCGT	NFKB_known9	8.13	Activating	Neither	8.1	0.9	0.1	0.1	7.2	0.1	175	0.36	0.02	0.35	2.6	HepG2	47	0.814	24	0.16
CGCGT	CREB3L1_4	8.05	Activating	Activating	8.0	3.4	0.2	0.4	4.6	-0.3	69	0.60	0.37	0.23	0.0	N/S	20	0.649	22	0.51
CGCGT	TP63_1	7.95	Activating	Neither	7.9	0.3	0.3	0.6	7.6	-0.3	54	0.99	-0.08	1.07	0.0	N/S	19	1.718	11	-0.12
CGCGT	ARNT_2	7.83	Activating	Activating	7.0	7.8	0.4	0.3	-0.8	0.1	39	0.76	0.60	0.17	0.0	N/S	14	1.234	4	0.37
CGCGT	MAFF_1	7.81	Dual	Activating	7.8	5.2	2.0	1.2	2.6	0.9	159	0.30	0.33	0.03	0.0	N/S	40	0.821	47	0.43
CGCGT	CREB3L2_2	7.41	Activating	Activating	5.0	7.4	0.3	0.2	-2.4	0.1	49	0.48	0.54	0.06	0.0	N/S	10	0.919	12	0.48
CGCGT	MYCN_2	7.37	Neither	Activating	1.4	7.4	0.4	0.3	-6.0	0.1	100	0.15	0.39	0.24	0.5	N/S	22	0.118	34	0.65
CGCGT	RARA_1	7.23	Activating	Neither	7.2	0.6	0.1	0.4	6.6	-0.3	240	0.27	0.03	0.24	0.9	N/S	107	0.43	53	0.21
CGCGT	KLF13_1	7.21	Neither	Activating	1.5	7.2	0.5	0.0	-5.7	0.4	142	0.18	0.43	0.26	2.6	K562	36	0.239	33	0.48
CGCGT	MLX_2	7.05	Activating	Activating	7.1	4.9	0.4	0.4	2.1	0.1	32	0.76	0.58	0.18	0.0	N/S	9	1.287	5	0.65
CGCGT	ATF6_1	6.73	Activating	Neither	6.7	1.8	0.2	0.2	5.0	0.1	57	0.43	0.29	0.14	0.0	N/S	17	0.551	14	0.42
CGCGT	HINFP_2	6.60	Activating	Activating	6.6	4.8	0.3	0.6	1.9	-0.4	50	0.48	0.48	0.00	0.0	N/S	9	0.874	12	0.17
CGCGT	ZNF263_disc1	6.40	Activating	Activating	2.6	6.4	0.9	0.4	-3.8	0.5	1013	0.06	0.10	0.04	0.0	N/S	221	0.085	221	0.26
CGCGT	PATZ1_1	6.37	Activating	Activating	4.1	6.4	0.1	0.0	-2.3	0.1	367	0.16	0.21	0.04	0.0	N/S	99	0.222	93	0.28
CGCGT	PAX4_1	6.33	Activating	Activating	6.3	4.7	0.6	0.9	1.7	-0.3	70	0.47	0.55	0.08	0.0	N/S	21	0.394	16	1.42
CGCGT	NFATC1_3	6.02	Repressive	Neither	0.2	0.0	6.0	0.2	0.2	5.8	113	-0.11	-0.09	0.01	0.0	N/S	29	0.113	22	0.18
CGCGT	HNF1B_4	5.88	Activating	Neither	5.9	0.0	0.1	0.1	5.9	0.0	198	0.28	-0.10	0.38	5.6	HepG2	111	0.388	20	-0.06
CGCGT	NFAT5_1	5.88	Activating	Neither	5.9	0.1	0.3	0.2	5.8	0.1	236	0.27	-0.05	0.31	4.3	HepG2	57	0.755	47	0.04
CGCGT	SP1B_2	5.77	Activating	Activating	2.6	5.8	0.1	0.3	-3.1	-0.2	360	0.09	0.24	0.14	0.4	N/S	54	0.258	136	0.23
CGCGT	MAFG_1	5.76	Dual	Activating	5.8	3.4	2.9	1.8	2.3	1.1	141	0.21	0.26	0.05	0.0	N/S	28	0.992	43	0.19
CGCGT	ZNF219_1	5.65	Neither	Activating	1.0	5.6	0.1	0.1	-4.6	0.0	534	0.11	0.14	0.03	0.0	N/S	140	0.148	108	0.21
CGCGT	CACD_1	5.58	Neither	Activating	1.9	5.6	0.4	1.7	-3.7	-1.3	222	0.09	0.13	0.04	0.0	N/S	51	0.22	48	0.35
CGCGT	E4F1_1	5.50	Activating	Activating	5.5	3.2	0.7	0.6	2.3	0.1	16	1.18	1.09	0.10	0.0	N/S	5	0.488	2	2.66
CGCGT	HNF1A_1	5.18	Activating	Neither	5.2	0.0	0.4	0.5	5.2	0.0	216	0.22	-0.14	0.36	5.2	HepG2	113	0.378	32	-0.03
CGCGT	HNF1_1	4.82	Activating	Neither	4.8	0.0	0.4	0.5	4.8	0.0	213	0.21	-0.15	0.35	4.5	HepG2	107	0.376	36	-0.12
CGCGT	NR4A_known4	4.75	Activating	Neither	4.7	0.8	0.4	0.1	3.9	0.3	215	0.20	0.02	0.19	0.0	N/S	93	0.457	41	0.18
CGCGT	CEBPD_1	4.75	Activating	Activating	2.6	4.7	0.1	0.1	-2.1	0.0	89	0.23	0.34	0.11	0.0	N/S	28	0.189	20	0.51
CGCGT	RARB_1	4.52	Activating	Neither	4.5	0.0	0.3	0.1	4.5	0.2	195	0.17	-0.02	0.19	0.0	N/S	84	0.362	40	0.03
CGCGT	EGR3_3	4.35	Activating	Activating	2.9	4.3	0.8	0.1	-1.5	0.7	96	0.24	0.40	0.16	0.0	N/S	20	0.271	27	0.64
CGCGT	MYB_1	4.33	Activating	Activating	4.3	3.0	0.5	0.1	1.4	0.4	78	0.25	0.34	0.09	0.0	N/S	16	0.		



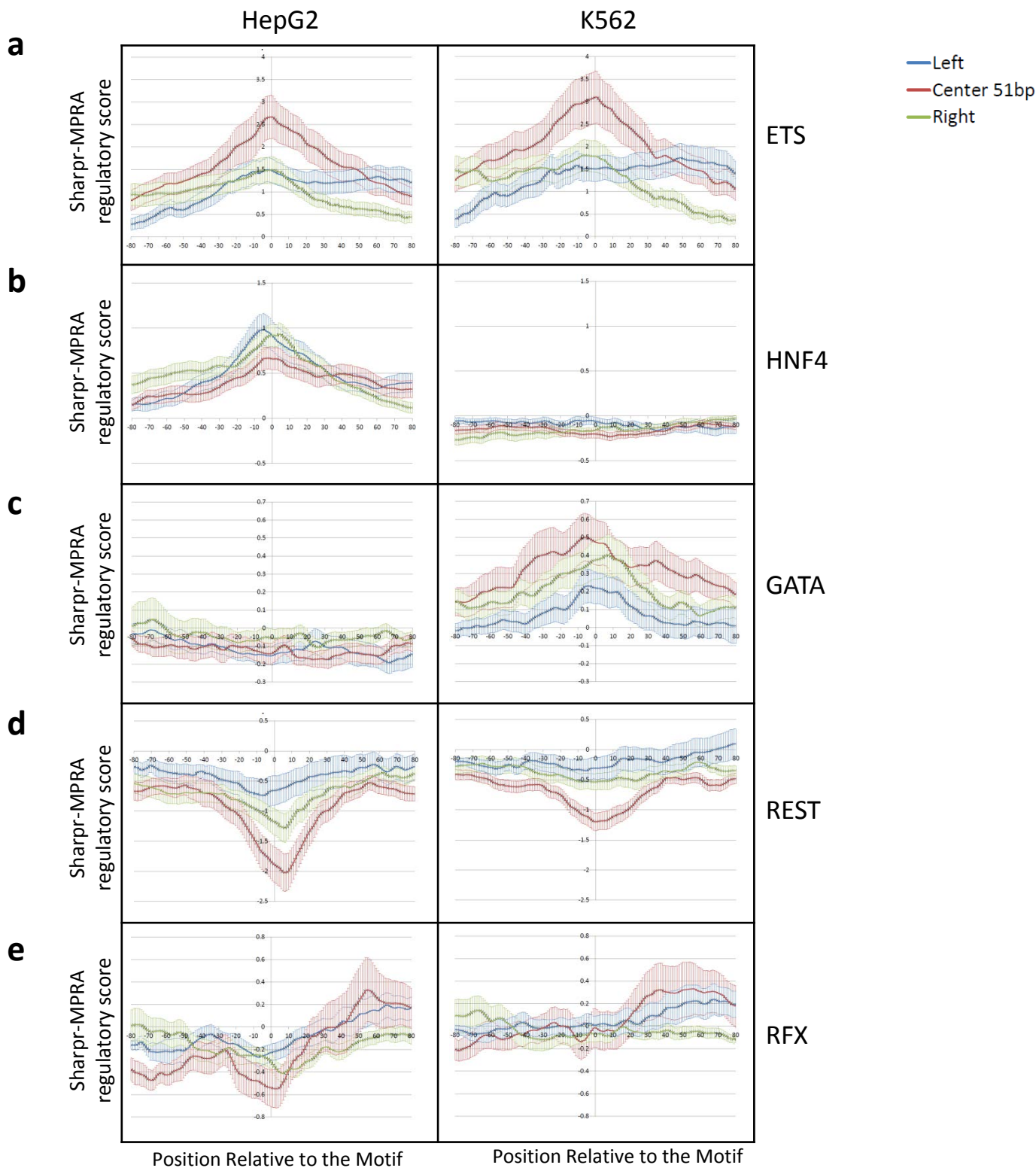
**a. Comparison of motif activating region enrichment between HepG2 and K562**



**b. Comparison of motif repressive region enrichment between HepG2 and K562**



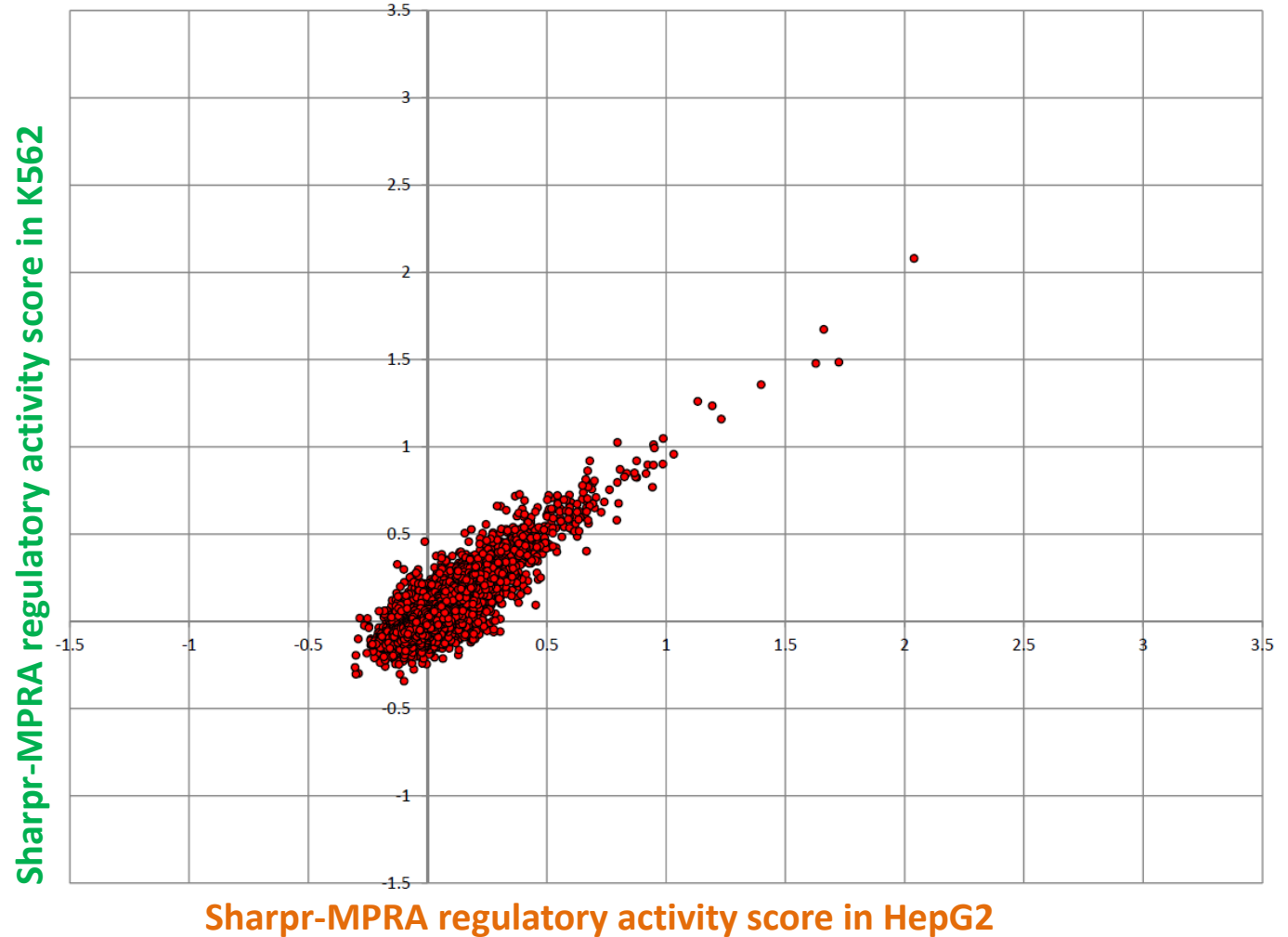
**Supplementary Figure 26 – Scatterplot of regulatory motif enrichments (a) comparing the  $-\log_{10}$  p-value of the enrichment for regulatory motif instances with activity  $\geq 1$  at the center position in HepG2 (y-axis) and in K562 (x-axis) for all regulatory motifs shown in **Supplementary Fig. 25** and in **Supplementary Table 2**. (b) Same as (a) except considering repression based on enrichment for regulatory motif instances with activity  $\leq -1$ .**



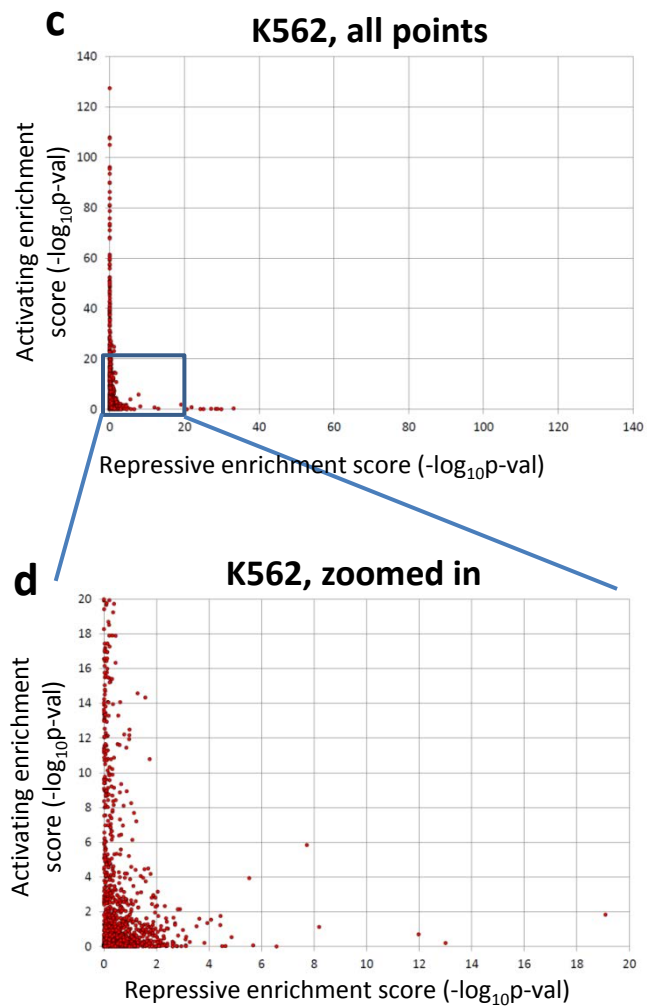
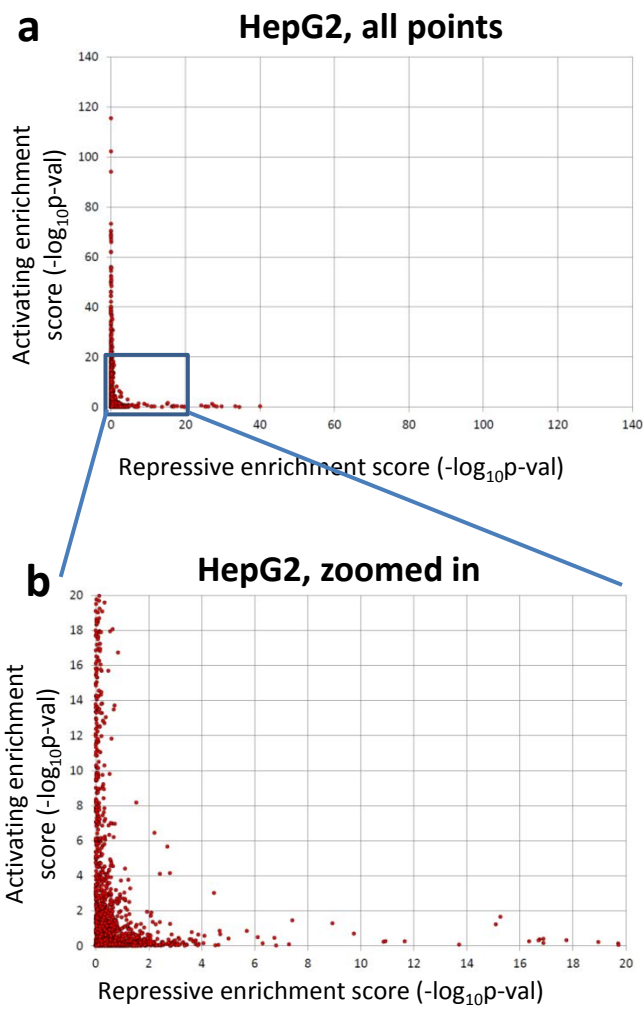
**Supplementary Figure 27 – Sharpr-MPRA Regulatory Score Aggregation Plots Relative to Motif Position.** Aggregation plots of Sharpr-MPRA regulatory scores relative to motif position from **Fig. 4b** shown separately for instances whose motif center fell within the central 51 bp (red), or to the left (blue) or right (green) of that in HepG2 (left) and K562 (right) cells, shown for: **(a)** a ETS motif, ETS\_known9; **(b)** a HNF4 motif, HNF4\_known18; **(c)** a GATA motif, GATA\_known14; **(d)** a REST motif, REST\_known2; **(e)** a RFX motif, RFX5\_known6, as specified in Ref. 13, which predicted motif instances independent of cell type. Vertical error bars indicate standard error.



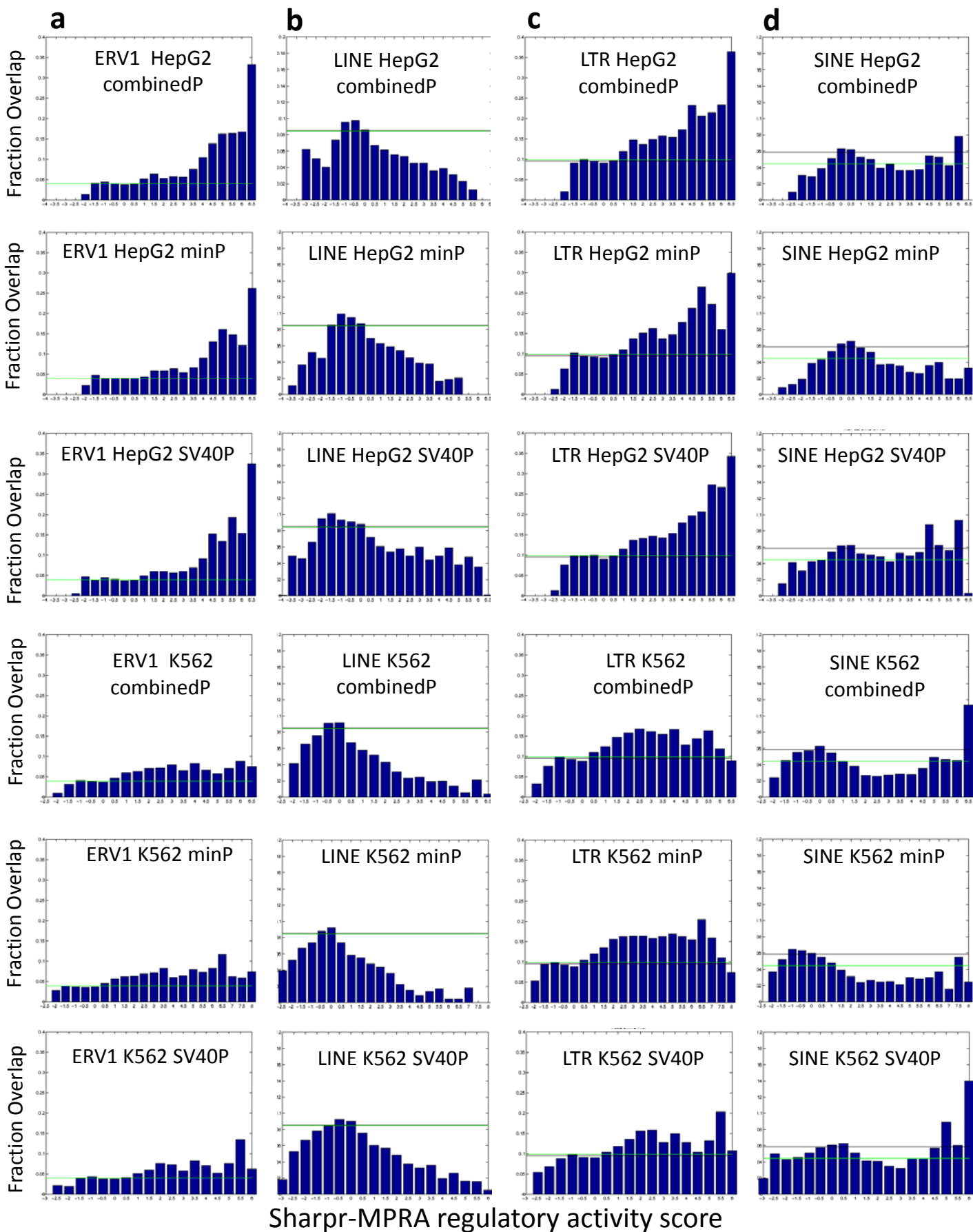
## Scatter plot of Sharpr-MPRA Regulatory Activity Score of 7-mers



**Supplementary Figure 28 – 7-mer Sharpr-MPRA regulatory activity score plot.** The plot has a point for each 7-mer appearing more than ten times based on the forward strand showing the average regulatory activity score in HepG2 cells (x-axis) and the average regulatory activity score in K562 cells (y-axis).

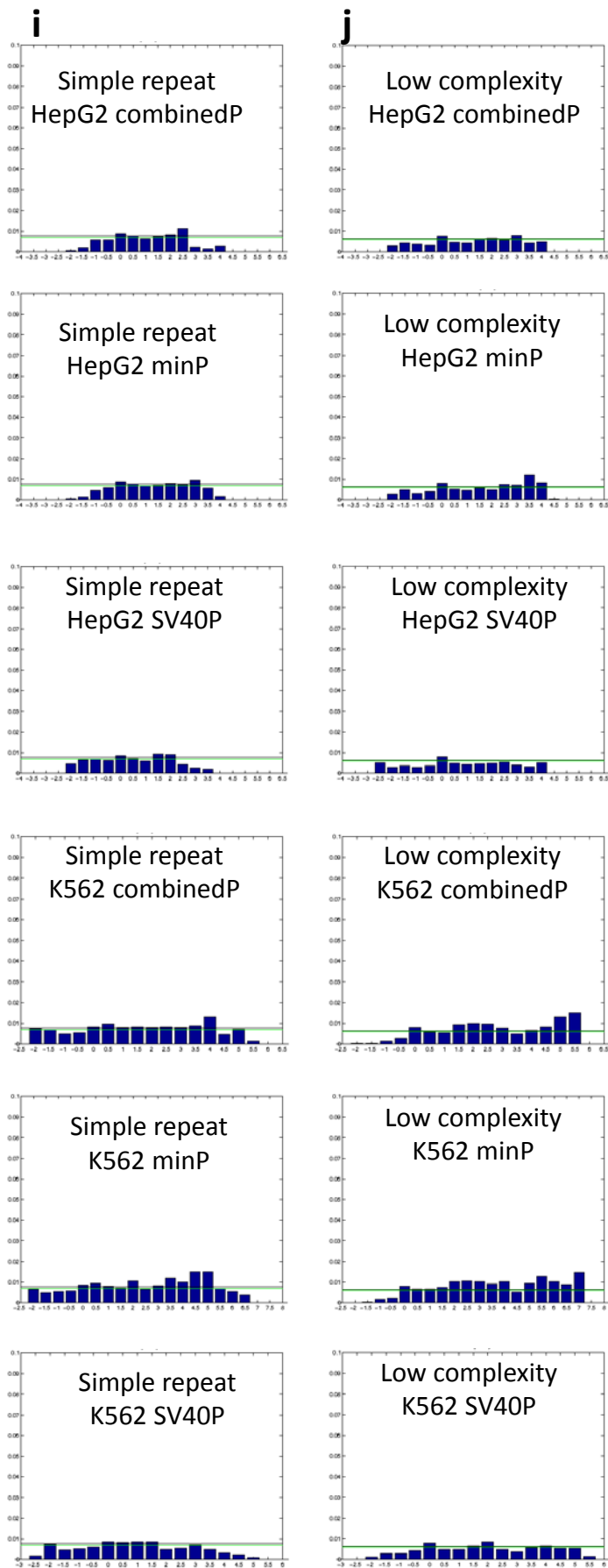


**Supplementary Figure 29 – Enrichment for regulatory motif instances in top 5% activated and top 5% repressed nucleotides.** These are similar to the plots shown in **Fig. 4c** except here the activating and repressive bases are defined as the 5% nucleotides that were given the most activating and repressive scores respectively instead of using the 1 and -1 thresholds. Activator score (y-axis) and repressor score (x-axis) for a compendium of regulatory motifs<sup>13</sup> (points) in HepG2 (left) and K562 (right), based on the statistical significance ( $-\log_{10}P$ ) for the enrichment of the center position of motif instances in positions that had a regulatory score in the lowest 5% (most repressive) or highest 5% (most activating), computed based on one-side binomial tests where the probability of success for the binomial distribution is the fraction of total nucleotides tested that met the threshold. Motif instances that appeared on both strands at the same position are only counted once. **(a,b)** The scatter plots are for **(a)** all points in HepG2 and **(b)** a zoomed in view of the points in HepG2 that have values less than 20 on both axes. **(c,d)** are similar except for K562.



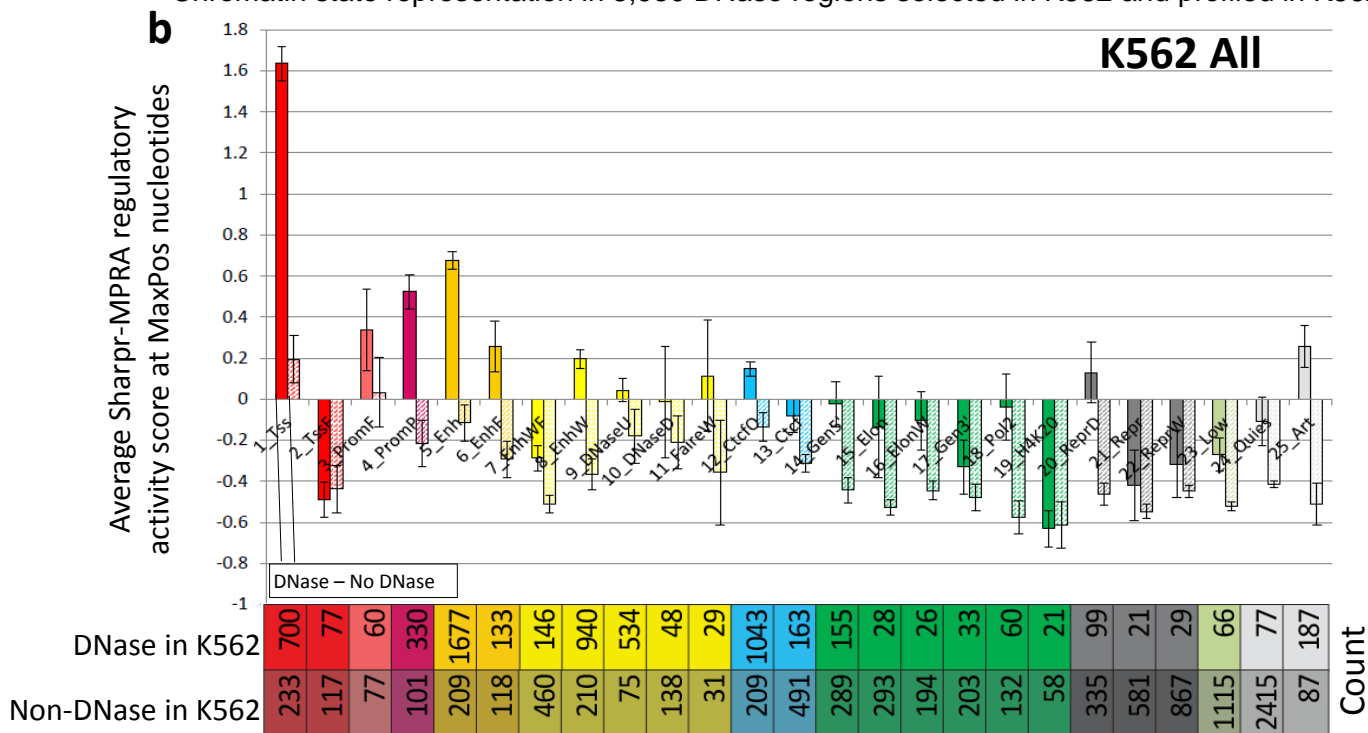
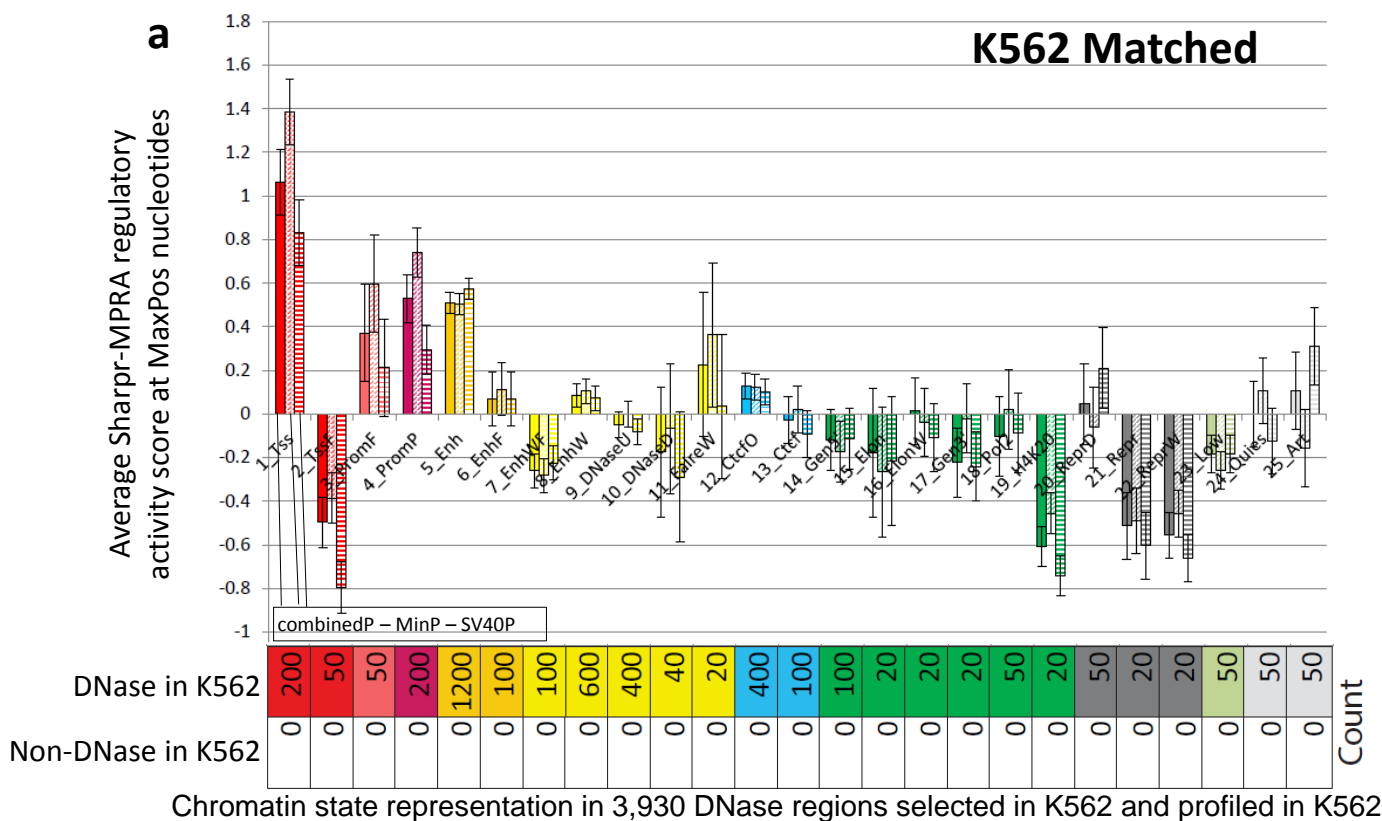
**Supplementary Figure 30a-d – Regulatory Activity in Repeat Classes and Families (a-d) This page. (e-h) Next page. (i,j) Following page. For legend, see i,j.**





**Supplementary Figure 30 – Regulatory Activity in Repeat Classes and Families.** Extended set of enrichments shown in **Fig. 5** now showing in each column for **(a)** ERV1, **(b)** LINE, **(c)** LTR, **(d)** SINE, **(e)** DNA, **(f)** ERVL, **(g)** ERVL-MaLR, **(h)** ERVK, **(i)** Simple Repeat, and **(j)** Low Complexity repeats as defined by RepeatMasker<sup>70</sup> for regulatory activity scores from top to bottom: HepG2 combinedP, HepG2 minP, HepG2 SV40P, K562 combinedP, K562 minP, and K562 SV40P data. Bins were formed by assigning each base to the nearest 0.5 value based on its regulatory score, and the two extreme bins contain all more extreme values. Green line denotes the expected fraction of overlap based on the center position (CenPos), and the black line the expected fraction based on all positions. Unlike in the corresponding main figure all bars are shown in the same blue color regardless if they are above or below expectation.

Sharpr-MPRA regulatory activity score



**Supplementary Figure 31 – Activity by Chromatin State – K562 cells.** Analogous figures as **Fig. 6 a** and **b** except for K562 cells. **(a)** For each chromatin state (x-axis) the average K562 regulatory score of all MaxPos nucleotides in K562, over all regions selected based on that state in that cell type. Results are shown for the combinedP, minP, and SV40P results (consecutive bars). The number of regions selected based on each state in the cell type is shown on the bottom. Since the regions were selected based on the indicated cell type they all have DNase sites. **(b)** For each chromatin state, the average score at MaxPos nucleotides for regions with the center position in the state based on all locations tested shown separately for locations overlapping DNase sites and those not overlapping DNase sites in K562. The number of regions in each set is indicated along the bottom. The vertical error bar indicates one standard error.

## a. DNase 'matched' (selected in same cell type as tested)

HepG2 state		HepG2 Combined			HepG2 minP			HepG2 SV40P			
id	state	All	Act	Rep	Non	Act	Rep	Non	Act	Rep	Non
1	Tss	200	36%	9%	55%	38%	15%	48%	37%	11%	52%
2	TssF	50	8%	10%	82%	12%	20%	68%	12%	18%	70%
3	PromF	50	16%	16%	68%	18%	16%	66%	14%	24%	62%
4	PromP	200	23%	15%	62%	25%	20%	56%	21%	18%	61%
5	Enh	1200	29%	7%	64%	30%	14%	56%	32%	14%	54%
6	EnhF	100	15%	10%	75%	18%	11%	71%	16%	18%	66%
7	EnhWF	100	9%	11%	80%	12%	23%	65%	10%	21%	69%
8	EnhW	600	22%	11%	67%	21%	14%	65%	27%	19%	55%
9	DNaseU	400	18%	13%	69%	19%	17%	64%	22%	17%	62%
10	DNaseD	40	18%	15%	68%	20%	25%	55%	23%	23%	55%
11	FaireW	20	10%	0%	90%	10%	5%	85%	15%	5%	80%
12	CtcfO	400	15%	17%	69%	17%	21%	62%	19%	22%	59%
13	Ctcf	100	2%	11%	87%	5%	17%	78%	4%	16%	80%
14	Gen5'	100	24%	13%	63%	22%	22%	56%	21%	20%	59%
15	Elon	20	15%	10%	75%	15%	15%	70%	15%	10%	75%
16	ElonW	20	15%	10%	75%	20%	20%	60%	15%	20%	65%
17	Gen3'	20	5%	25%	70%	5%	25%	70%	5%	40%	55%
18	Pol2	50	20%	14%	66%	22%	24%	54%	18%	16%	66%
19	H4K20	20	0%	20%	80%	5%	25%	70%	5%	20%	75%
20	ReprD	50	12%	26%	62%	18%	34%	48%	16%	28%	56%
21	Repr	20	0%	35%	65%	5%	35%	60%	0%	25%	75%
22	ReprW	20	20%	25%	55%	15%	35%	50%	20%	35%	45%
23	Low	50	16%	14%	70%	16%	20%	64%	14%	26%	60%
24	Quies	50	6%	16%	78%	8%	20%	72%	6%	24%	70%
25	Art	50	12%	16%	72%	14%	28%	58%	22%	22%	56%
All states		3930	22%	11%	67%	23%	17%	60%	24%	17%	58%

K562 state		K562 Combined			K562 minP			K562 SV40P			
id	state	All	Act	Rep	Non	Act	Rep	Non	Act	Rep	Non
1	Tss	200	41%	13%	46%	44%	16%	41%	44%	19%	37%
2	TssF	50	6%	20%	74%	10%	26%	64%	6%	42%	52%
3	PromF	50	24%	8%	68%	30%	12%	58%	30%	20%	50%
4	PromP	200	34%	18%	49%	42%	17%	42%	36%	27%	38%
5	Enh	1200	32%	13%	55%	31%	16%	53%	40%	21%	39%
6	EnhF	100	14%	8%	78%	13%	9%	78%	21%	21%	58%
7	EnhWF	100	8%	13%	79%	9%	18%	73%	11%	28%	61%
8	EnhW	600	22%	18%	61%	21%	20%	59%	30%	30%	41%
9	DNaseU	400	11%	11%	78%	13%	11%	77%	18%	24%	58%
10	DNaseD	40	18%	40%	43%	18%	38%	45%	25%	50%	25%
11	FaireW	20	15%	0%	85%	15%	5%	80%	10%	5%	85%
12	CtcfO	400	20%	12%	68%	20%	12%	68%	31%	25%	44%
13	Ctcf	100	12%	13%	75%	9%	12%	79%	18%	22%	60%
14	Gen5'	100	13%	25%	62%	10%	21%	69%	20%	34%	46%
15	Elon	20	10%	20%	70%	5%	20%	75%	10%	30%	60%
16	ElonW	20	10%	5%	85%	10%	5%	85%	15%	20%	65%
17	Gen3'	20	5%	5%	90%	5%	5%	90%	15%	25%	60%
18	Pol2	50	14%	12%	74%	12%	12%	76%	30%	24%	46%
19	H4K20	20	0%	15%	85%	0%	5%	95%	0%	40%	60%
20	ReprD	50	24%	20%	56%	24%	24%	52%	34%	26%	40%
21	Repr	20	0%	25%	75%	5%	20%	75%	5%	40%	55%
22	ReprW	20	0%	10%	90%	5%	10%	85%	10%	40%	50%
23	Low	50	2%	10%	88%	2%	6%	92%	18%	24%	58%
24	Quies	50	14%	10%	76%	14%	12%	74%	16%	18%	66%
25	Art	50	18%	14%	68%	16%	22%	62%	36%	16%	48%
All states		3930	23%	14%	63%	23%	16%	61%	30%	25%	45%

## b. DNase 'unmatched' (selected in different cell type from tested)

HepG2 state		HepG2 Combined			HepG2 minP			HepG2 SV40P			
id	state	All	Act	Rep	Non	Act	Rep	Non	Act	Rep	Non
1	Tss	599	53%	5%	42%	55%	9%	36%	53%	8%	39%
2	TssF	23	13%	22%	65%	13%	26%	61%	17%	35%	48%
3	PromF	4	25%	0%	75%	25%	25%	50%	25%	0%	75%
4	PromP	364	26%	10%	64%	31%	15%	54%	28%	14%	58%
5	Enh	191	43%	10%	47%	46%	17%	37%	47%	16%	37%
6	EnhF	13	23%	8%	69%	23%	15%	62%	15%	15%	69%
7	EnhWF	13	0%	15%	85%	8%	15%	77%	0%	23%	77%
8	EnhW	154	34%	12%	53%	37%	18%	45%	37%	16%	47%
9	DNaseU	45	38%	11%	51%	38%	27%	36%	33%	13%	53%
10	DNaseD	30	33%	20%	47%	37%	20%	43%	40%	23%	37%
11	FaireW	3	0%	33%	67%	33%	67%	0%	0%	67%	33%
12	CtcfO	319	14%	15%	70%	18%	24%	58%	18%	23%	60%
13	Ctcf	9	11%	0%	89%	11%	11%	78%	11%	11%	78%
14	Gen5'	27	26%	15%	59%	30%	15%	56%	30%	22%	48%
15	Elon	4	25%	25%	50%	25%	0%	75%	50%	25%	25%
16	ElonW	1	100%	0%	0%	100%	0%	0%	100%	0%	0%
17	Gen3'	11	45%	0%	55%	36%	0%	64%	45%	18%	36%
18	Pol2	8	50%	0%	50%	50%	13%	38%	75%	0%	25%
19	H4K20	0									
20	ReprD	25	28%	8%	64%	28%	16%	56%	28%	24%	48%
21	Repr	0									
22	ReprW	1	0%	0%	100%	0%	0%	100%	0%	0%	100%
23	Low	4	50%	0%	50%	25%	25%	50%	25%	0%	75%
24	Quies	1	0%	0%	100%	0%	0%	100%	0%	0%	100%
25	Art	61	20%	18%	62%	23%	20%	57%	21%	25%	54%
All states		1910	35%	10%	55%	38%	16%	46%	37%	15%	48%

K562 state		K562 Combined			K562 minP			K562 SV40P			
id	state	All	Act	Rep	Non	Act	Rep	Non	Act	Rep	Non
1	Tss	500	60%	9%	31%	64%	8%	28%	59%	15%	26%
2	TssF	27	0%	19%	81%	4%	11%	85%	11%	41%	48%
3	PromF	10	30%	0%	70%	20%	0%	80%	20%	30%	50%
4	PromP	130	38%	18%	45%	38%	17%	45%	41%	22%	38%
5	Enh	477	47%	11%	42%	47%	13%	40%	55%	18%	27%
6	EnhF	33	36%	6%	58%	33%	9%	58%	45%	18%	36%
7	EnhWF	46	2%	13%	85%	4%	17%	78%	11%	26%	63%
8	EnhW	340	32%	18%	50%	33%	20%	47%	41%	29%	30%
9	DNaseU	134	27%	12%	61%	26%	13%	61%	28%	22%	50%
10	DNaseD	8	38%	13%	50%	38%	13%	50%	63%	13%	25%
11	FaireW	9	22%	33%	44%	22%	33%	44%	11%	67%	22%
12	CtcfO	643	21%	13%	66%	21%	14%	65%	33%	23%	43%
13	Ctcf	63	6%	6%	87%	6%	5%	89%	10%	22%	68%
14	Gen5'	55	27%	13%	60%	25%	15%	60%	27%	35%	38%
15	Elon	8	38%	25%	38%	25%	25%	50%	25%	38%	38%
16	ElonW	6	0%	17%	83%	0%	17%	83%	17%	33%	50%
17	Gen3'	13	8%	23%	69%	8%	23%	69%	15%	31%	54%
18	Pol2	10	30%	10%	60%	20%	30%	50%	50%	10%	40%
19	H4K20	1	0%	100%	0%	0%	100%	0%	0%	100%	0%
20	ReprD	49	37%	27%	37%	33%	20%	47%	39%	41%	20%
21	Repr	1	100%	0%	0%	100%	0%	0%	100%	0%	0%
22	ReprW	9	44%	22%	33%	44%	33%	22%	44%	22%	33%
23	Low	16	0%	25%	75%	0%	31%	69%	19%	31%	50%
24	Quies	27	11%	19%	70%	11%	19%	70%	19%	37%	44%
25	Art	137	29%	18%	53%	26%	23%	51%	42%	24%	34%
All states		2752	35%	13%	52%	36%	14%	50%	42%	23%	35%

**Supplementary Figure 32a,b – Fraction of regions showing above activating or below repressive threshold at the maximum absolute score position (MaxPos). (a)** For each chromatin state (rows) in HepG2 (left panel) and K562 (right panel), the number of regions tested ('All') that were selected based on the chromatin state annotation in the same cell type (DNase Matched), the percentage of those regions with score at MaxPos  $\geq 1$  ('Act' for 'Activating'),  $\leq -1$  ('Rep' for 'Repressive'), or neither ('Non') using the combinedP score (first column group), minP score (second group), or SV40P score (third group). Each individual column is colored based on median (white), the 90<sup>th</sup> percentile (red), and the 10<sup>th</sup> percentile (blue), thus indicating chromatin states that are more often activating (e.g. Tss, Enh) or repressive (e.g. ReprD, Repr) than expected on average. Black boxes highlight the numbers discussed in the main text. **(b)** Similar to (a), but for DNase sites that were selected in a different cell type and are also DNase in the tested cell type. **(c-e)** Next pages.

### c. DNase All ('matched' and 'unmatched')

	HepG2 state				HepG2 Combined			HepG2 minP			HepG2 SV40P		
	id	state	All	Act	Rep	Non	Act	Rep	Non	Act	Rep	Non	
1 Tss		799	49%	6%	45%	51%	11%	39%	49%	9%	42%		
2 TssF		73	10%	14%	77%	12%	22%	66%	14%	23%	63%		
3 PromF		54	17%	15%	69%	19%	17%	65%	15%	22%	63%		
4 PromP		564	25%	12%	63%	29%	17%	54%	26%	15%	59%		
5 Enh		1391	31%	8%	61%	32%	14%	53%	34%	14%	52%		
6 EnhF		113	16%	10%	74%	19%	12%	70%	16%	18%	66%		
7 EnhWF		113	8%	12%	81%	12%	22%	66%	9%	21%	70%		
8 EnhW		754	25%	11%	64%	25%	14%	61%	29%	18%	53%		
9 DNaseU		445	20%	12%	67%	21%	18%	61%	23%	16%	61%		
10 DNaseD		70	24%	17%	59%	27%	23%	50%	30%	23%	47%		
11 FaireW		23	9%	4%	87%	13%	13%	74%	13%	13%	74%		
12 CctcO		719	14%	16%	70%	18%	22%	61%	18%	22%	59%		
13 CctcF		109	3%	10%	87%	6%	17%	78%	5%	16%	80%		
14 Gen5'		127	24%	13%	62%	24%	20%	56%	23%	20%	57%		
15 Elon		24	17%	13%	71%	17%	13%	71%	21%	13%	67%		
16 ElonW		21	19%	10%	71%	24%	19%	57%	19%	19%	62%		
17 Gen3'		31	19%	16%	65%	16%	16%	68%	19%	32%	48%		
18 Pol2		58	24%	12%	64%	26%	22%	52%	26%	14%	60%		
19 H4K20		20	0%	20%	80%	5%	25%	70%	5%	20%	75%		
20 ReprD		75	17%	20%	63%	21%	28%	51%	20%	27%	53%		
21 Repr		20	0%	35%	65%	5%	35%	60%	0%	25%	75%		
22 ReprW		21	19%	24%	57%	14%	33%	52%	19%	33%	48%		
23 Low		54	19%	13%	69%	17%	20%	63%	15%	24%	61%		
24 Quies		51	6%	16%	78%	8%	20%	73%	6%	24%	71%		
25 Art		111	16%	17%	67%	19%	23%	58%	22%	23%	55%		
All states		5840	26%	11%	63%	28%	16%	56%	28%	17%	55%		

	K562 state				K562 Combined			K562 minP			K562 SV40P		
	id	state	All	Act	Rep	Non	Act	Rep	Non	Act	Rep	Non	
1 Tss		700	55%	10%	35%	58%	10%	31%	55%	16%	29%		
2 TssF		77	4%	19%	77%	8%	21%	71%	8%	42%	51%		
3 PromF		60	25%	7%	68%	28%	10%	62%	28%	22%	50%		
4 PromP		330	35%	18%	47%	40%	17%	43%	38%	25%	38%		
5 Enh		1677	36%	12%	51%	36%	15%	49%	44%	20%	35%		
6 EnhF		133	20%	8%	73%	18%	9%	73%	27%	20%	53%		
7 EnhWF		146	6%	13%	81%	8%	18%	75%	11%	27%	62%		
8 EnhW		940	25%	18%	57%	25%	20%	55%	34%	30%	37%		
9 DNaseU		534	15%	11%	74%	16%	11%	73%	21%	24%	56%		
10 DNaseD		48	21%	35%	44%	21%	33%	46%	31%	44%	25%		
11 FaireW		29	17%	10%	72%	17%	14%	69%	10%	24%	66%		
12 CctcO		1043	21%	13%	67%	20%	13%	66%	33%	24%	43%		
13 CctcF		163	10%	10%	80%	8%	9%	83%	15%	22%	63%		
14 Gen5'		155	18%	21%	61%	15%	19%	66%	23%	34%	43%		
15 Elon		28	18%	21%	61%	11%	21%	68%	14%	32%	54%		
16 ElonW		26	8%	8%	85%	8%	8%	85%	15%	23%	62%		
17 Gen3'		33	6%	12%	82%	6%	12%	82%	15%	27%	58%		
18 Pol2		60	17%	12%	72%	13%	15%	72%	33%	22%	45%		
19 H4K20		21	0%	19%	81%	0%	10%	90%	0%	43%	57%		
20 ReprD		99	30%	23%	46%	28%	22%	49%	36%	33%	30%		
21 Repr		21	5%	24%	71%	10%	19%	71%	10%	38%	52%		
22 ReprW		29	14%	14%	72%	17%	17%	66%	21%	34%	45%		
23 Low		66	2%	14%	85%	2%	12%	86%	18%	26%	56%		
24 Quies		77	13%	13%	74%	13%	14%	73%	17%	25%	58%		
25 Art		187	26%	17%	57%	24%	22%	54%	41%	22%	37%		
All states		6682	28%	14%	58%	28%	15%	57%	35%	24%	41%		

### d. Non-DNase in tested cell type (selected as DNase in diff. cell type)

	HepG2 state				HepG2 Combined			HepG2 minP			HepG2 SV40P		
	id	state	All	Act	Rep	Non	Act	Rep	Non	Act	Rep	Non	
1 Tss		352	18%	17%	65%	20%	26%	54%	20%	24%	57%		
2 TssF		205	5%	17%	78%	5%	31%	63%	8%	27%	65%		
3 PromF		233	9%	14%	77%	10%	29%	61%	7%	21%	72%		
4 PromP		393	17%	18%	65%	20%	25%	55%	19%	20%	62%		
5 Enh		145	17%	12%	70%	20%	22%	58%	23%	18%	59%		
6 EnhF		142	14%	16%	70%	17%	27%	56%	13%	20%	66%		
7 EnhWF		272	8%	21%	71%	11%	29%	60%	9%	22%	69%		
8 EnhW		202	18%	13%	69%	20%	25%	55%	19%	19%	62%		
9 DNaseU		60	22%	18%	60%	20%	22%	58%	20%	25%	55%		
10 DNaseD		332	19%	23%	58%	22%	26%	52%	22%	30%	48%		
11 FaireW		20	10%	15%	75%	10%	40%	50%	10%	30%	60%		
12 CctcO		532	14%	15%	72%	17%	23%	60%	13%	23%	64%		
13 CctcF		528	8%	16%	76%	12%	23%	64%	9%	21%	70%		
14 Gen5'		253	9%	17%	73%	9%	26%	64%	9%	22%	69%		
15 Elon		588	6%	21%	73%	8%	30%	61%	6%	26%	68%		
16 ElonW		261	4%	15%	81%	6%	25%	69%	5%	24%	71%		
17 Gen3'		184	7%	15%	78%	10%	22%	68%	7%	21%	72%		
18 Pol2		124	15%	14%	72%	16%	23%	60%	14%	25%	61%		
19 H4K20		147	5%	22%	73%	5%	25%	69%	7%	29%	64%		
20 ReprD		329	13%	19%	67%	16%	25%	59%	13%	27%	60%		
21 Repr		469	7%	20%	74%	9%	30%	62%	7%	27%	66%		
22 ReprW		693	8%	19%	73%	11%	29%	60%	8%	26%	67%		
23 Low		1310	5%	19%	76%	8%	32%	60%	5%	26%	68%		
24 Quies		2043	6%	19%	75%	8%	29%	63%	7%	25%	69%		
25 Art		63	6%	19%	75%	8%	27%	65%	8%	25%	67%		
All states		9880	9%	18%	73%	11%	28%	61%	9%	25%	66%		

	K562 state				K562 Combined			K562 minP			K562 SV40P		
	id	state	All	Act	Rep	Non	Act	Rep	Non	Act	Rep	Non	
1 Tss		233	24%	27%	49%	27%	27%	46%	25%	36%	39%		
2 TssF		117	8%	27%	65%	8%	42%	50%	15%	36%	50%		
3 PromF		77	17%	22%	61%	19%	22%	58%	25%	25%	51%		
4 PromP		101	14%	23%	63%	13%	20%	67%	17%	38%	46%		
5 Enh		209	17%	26%	57%	17%	26%	57%	27%	33%	40%		
6 EnhF		118	11%	19%	70%	8%	35%	57%	20%	31%	48%		
7 EnhWF		460	8%	33%	59%	6%	39%	56%	17%	43%	40%		
8 EnhW		210	12%	32%	56%	11%	39%	50%	19%	42%	39%		
9 DNaseU		75	13%	20%	67%	15%	23%	63%	20%	29%	51%		
10 DNaseD		138	20%	29%	51%	19%	32%	49%	25%	39%	36%		
11 FaireW		31	3%	16%	81%	3%	23%	74%	6%	26%	68%		
12 CctcO		209	12%	20%	67%	12%	24%	64%	27%	29%	44%		
13 CctcF		491	10%	21%	69%	10%	24%	66%	18%	32%	50%		
14 Gen5'		289	8%	32%	60%	10%	40%	49%	20%	43%	37%		
15 Elon		293	1%	21%	77%	4%	22%	74%	9%	33%	58%		
16 ElonW		194	2%	14%	84%	2%	18%	80%	8%	27%	64%		
17 Gen3'		203	9%	23%	68%	10%	31%	59%	16%	38%	46%		
18 Pol2		132	6%	34%	60%	5%	32%	64%	19%	42%	39%		
19 H4K20		58	3%	28%	69%	5%	36%	59%	10%	47%	43%		
20 ReprD		335	10%	33%	57%	10%	42%	48%	21%	41%	38%		
21 Repr		581	6%	30%	64%	6%	41%	53%	18%	38%	44%		
22 ReprW		867	5%	23%	72%	6%	26%	68%	13%	36%	51%		
23 Low		1115	3%	23%	74%	3%	26%	71%	8%	35%	57%		
24 Quies		2415	5%	18%	77%	5%	20%	75%	10%	31%	58%		
25 Art		87	9%	37%	54%	7%	43%	51%	22%	44%	34%		
All states		9038	7%	24%	69%	7%	27%	65%	15%	35%	50%		

**Supplementary Figure 32c,d – Fraction of regions showing above activating or below repressive threshold at the maximum absolute score position (MaxPos). (a,b) Previous page. (c) Same as (a), but showing total counts and activity for all DNase in sites, regardless of the cell type in which they were selected. Total region counts ('All') are simply the sum of panels (a) and (b). Black boxes highlight the numbers discussed in the main text. (d) Same percentages but shown for the converse set regions, namely all the regions that do not show DNase activity in the tested cell type (and thus "Unmatched", as all regions selected were DNase in the cell type in which they were selected). Differences in the fraction of active regions between (a,b,c) vs. (d) reflect that among all regions in a given chromatin state, those that also overlap DNase regions are more likely to be activating (see also Fig. 6b and Supplementary Fig 31b). (e) Next Page.**



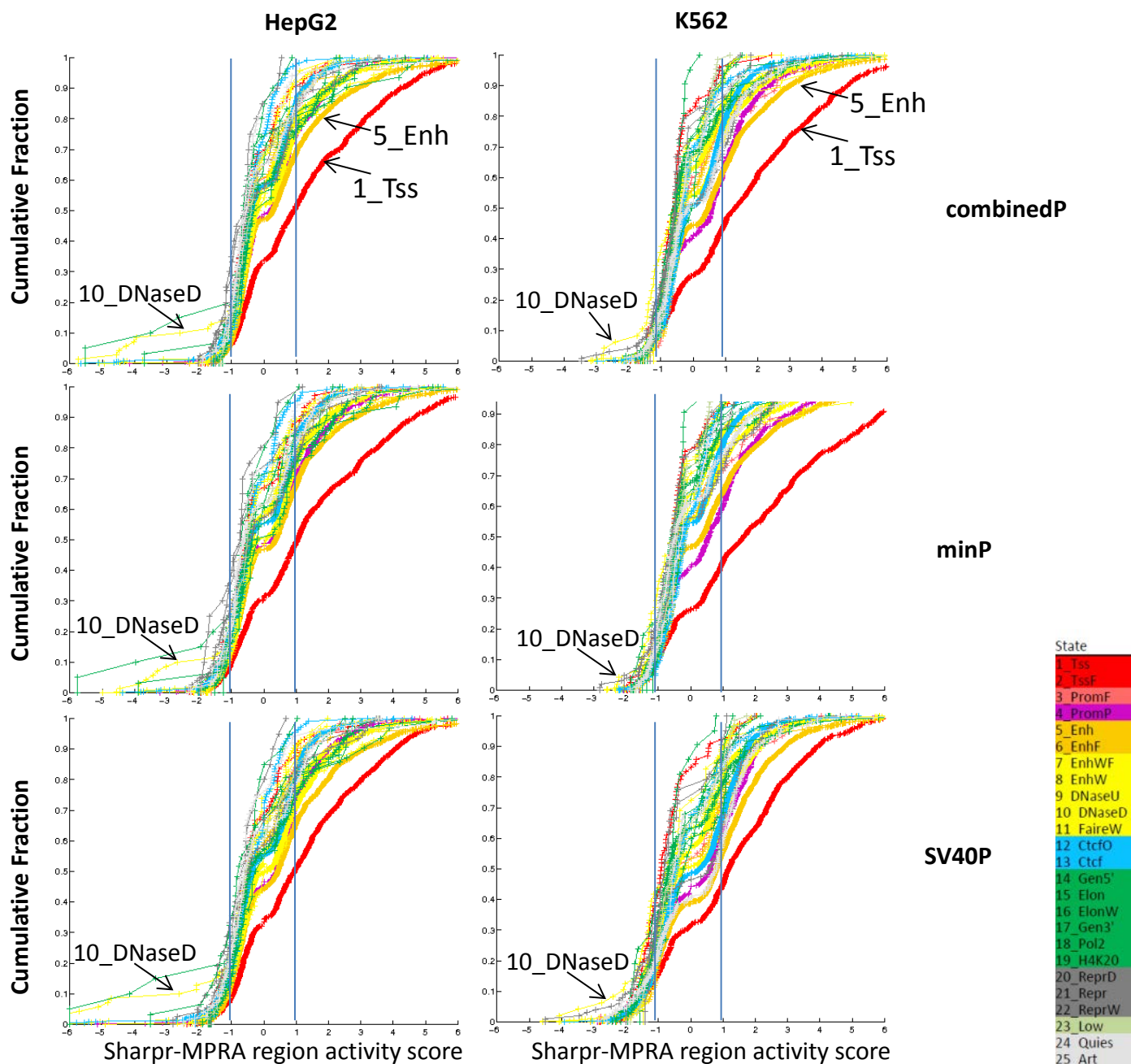
# e. All regions (DNase and non-DNase, matched and unmatched)

	HepG2 state			HepG2 Combined			HepG2 minP			HepG2 SV40P		
	id	state	All	Act	Rep	Non	Act	Rep	Non	Act	Rep	Non
	All 15720 region selections (all states, both DNase and non-DNase)	1	Tss	1151	39%	10%	51%	41%	15%	43%	40%	13%
	2	TssF	278	6%	16%	78%	7%	29%	64%	9%	26%	64%
	3	PromF	287	10%	14%	76%	12%	26%	62%	9%	21%	70%
	4	PromP	957	22%	14%	64%	25%	20%	55%	23%	17%	60%
	5	Enh	1536	30%	8%	62%	31%	15%	54%	33%	14%	53%
	6	EnhF	255	15%	13%	72%	18%	20%	62%	15%	19%	66%
	7	EnhWF	385	8%	18%	74%	11%	27%	62%	9%	22%	70%
	8	EnhW	956	23%	11%	65%	24%	17%	60%	27%	18%	55%
	9	DNaseU	505	20%	13%	67%	21%	19%	60%	23%	17%	60%
	10	DNaseD	402	20%	22%	58%	23%	26%	51%	23%	29%	48%
	11	FaireW	43	9%	9%	81%	12%	26%	63%	12%	21%	67%
	12	CtcfO	1251	14%	16%	70%	18%	22%	60%	16%	22%	62%
	13	Ctcf	637	7%	15%	78%	11%	22%	67%	8%	20%	71%
	14	Gen5'	380	14%	16%	69%	14%	24%	62%	14%	21%	65%
	15	Elon	612	6%	20%	73%	9%	30%	62%	6%	25%	68%
	16	ElonW	282	5%	15%	80%	7%	25%	68%	6%	23%	70%
	17	Gen3'	215	9%	15%	76%	11%	21%	68%	9%	23%	68%
	18	Pol2	182	18%	13%	69%	19%	23%	58%	18%	21%	61%
	19	H4K20	167	5%	22%	74%	5%	25%	69%	7%	28%	65%
	20	ReprD	404	14%	20%	66%	17%	25%	58%	14%	27%	59%
	21	Repr	489	6%	20%	73%	8%	30%	62%	7%	27%	66%
	22	ReprW	714	9%	19%	72%	11%	29%	59%	8%	26%	66%
	23	Low	1364	6%	19%	75%	8%	32%	60%	5%	26%	68%
	24	Quies	2094	6%	19%	75%	8%	29%	63%	7%	25%	69%
	25	Art	174	13%	18%	70%	15%	25%	60%	17%	24%	59%
	All states		15720	15%	15%	69%	17%	24%	59%	16%	22%	62%

	K562 state			K562 Combined			K562 minP			K562 SV40P		
	id	state	All	Act	Rep	Non	Act	Rep	Non	Act	Rep	Non
	All 15720 region selections (all states, both DNase and non-DNase)	1	Tss	933	47%	14%	39%	50%	14%	35%	47%	21%
	2	TssF	194	6%	24%	70%	8%	34%	59%	12%	38%	50%
	3	PromF	137	20%	15%	64%	23%	17%	60%	26%	23%	50%
	4	PromP	431	30%	19%	51%	34%	18%	48%	33%	28%	39%
	5	Enh	1886	34%	14%	52%	34%	16%	50%	42%	22%	36%
	6	EnhF	251	16%	13%	72%	14%	21%	65%	24%	25%	51%
	7	EnhWF	606	7%	28%	64%	6%	34%	60%	16%	39%	45%
	8	EnhW	1150	23%	20%	57%	23%	23%	54%	31%	32%	37%
	9	DNaseU	609	15%	12%	73%	16%	13%	71%	21%	24%	55%
	10	DNaseD	186	20%	31%	49%	19%	32%	48%	27%	40%	33%
	11	FaireW	60	10%	13%	77%	10%	18%	72%	8%	25%	67%
	12	CtcfO	1252	19%	14%	67%	19%	15%	66%	32%	25%	44%
	13	Ctcf	654	10%	18%	72%	9%	21%	70%	17%	30%	54%
	14	Gen5'	444	11%	28%	60%	12%	33%	55%	21%	40%	39%
	15	Elon	321	3%	21%	76%	4%	22%	74%	9%	33%	58%
	16	ElonW	220	3%	13%	84%	3%	16%	81%	9%	27%	64%
	17	Gen3'	236	8%	22%	70%	10%	28%	62%	16%	36%	47%
	18	Pol2	192	9%	27%	64%	7%	27%	66%	23%	36%	41%
	19	H4K20	79	3%	25%	72%	4%	29%	67%	8%	46%	47%
	20	ReprD	434	15%	31%	55%	14%	38%	48%	25%	39%	36%
	21	Repr	602	6%	30%	64%	6%	41%	53%	18%	38%	45%
	22	ReprW	896	5%	23%	72%	6%	26%	68%	14%	36%	50%
	23	Low	1181	3%	22%	75%	3%	25%	72%	9%	34%	57%
	24	Quies	2492	5%	18%	77%	6%	19%	75%	11%	31%	58%
	25	Art	274	21%	23%	56%	18%	29%	53%	35%	29%	36%
	All states		15720	16%	19%	65%	16%	22%	62%	23%	30%	46%

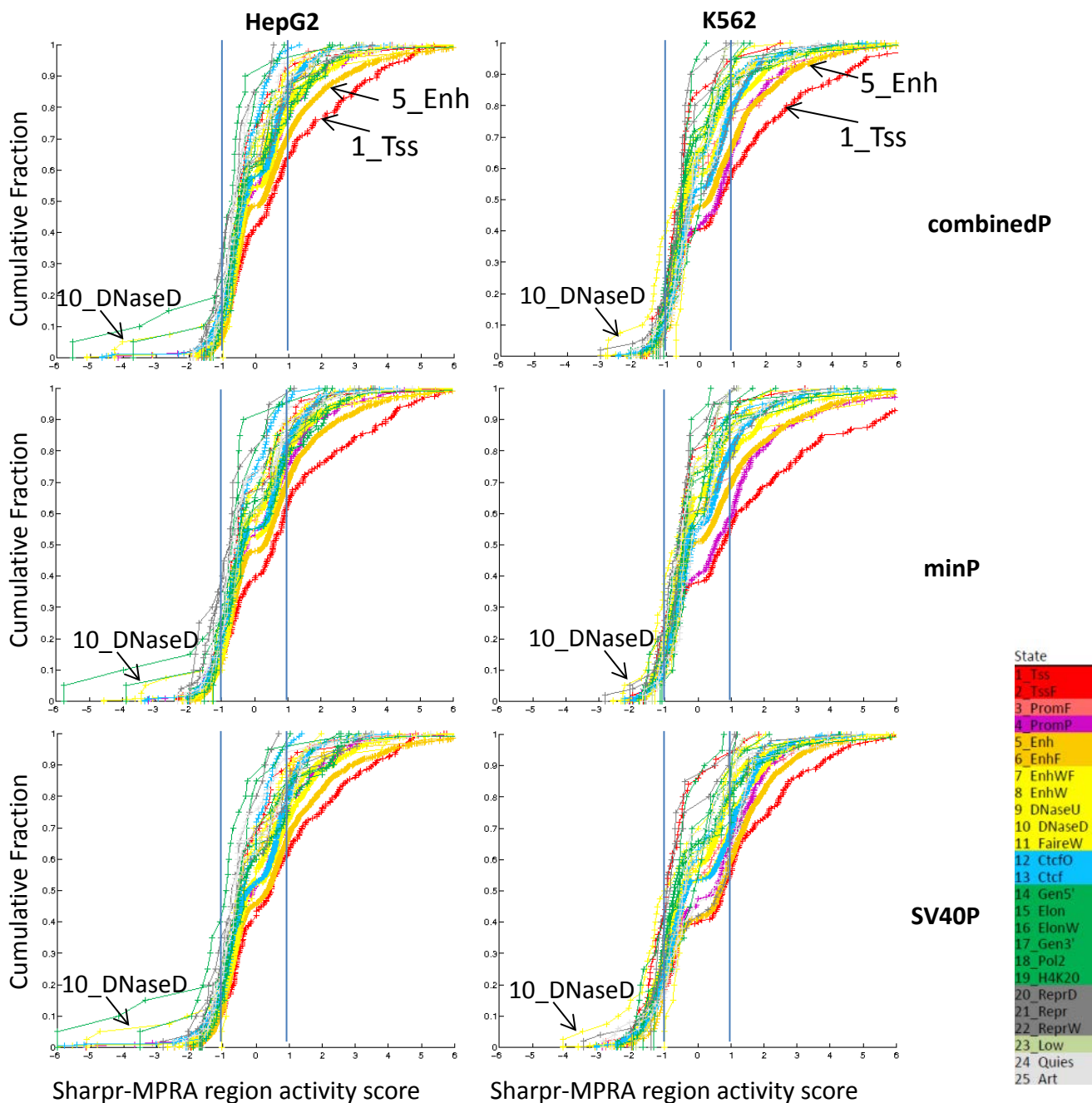
**Supplementary Figure 32e – Fraction of regions showing above activating or below repressive threshold at the maximum absolute score position (MaxPos). (a-d) Previous pages. (e) Same as (a-d), but showing all tested regions, regardless of what cell type they were selected in, and regardless of whether they are found in a DNase region or a non-DNase region. This shows that active chromatin states (e.g. Tss, Enh) are more likely to show activating regulatory scores ( $\geq 1$ ), that repressive chromatin states (e.g. Repr, ReprD, ReprW) are more likely to show repressive regulatory scores ( $\leq -1$ ).**

## a. DNase All (N=5840 in HepG2, N=6682 in K562)



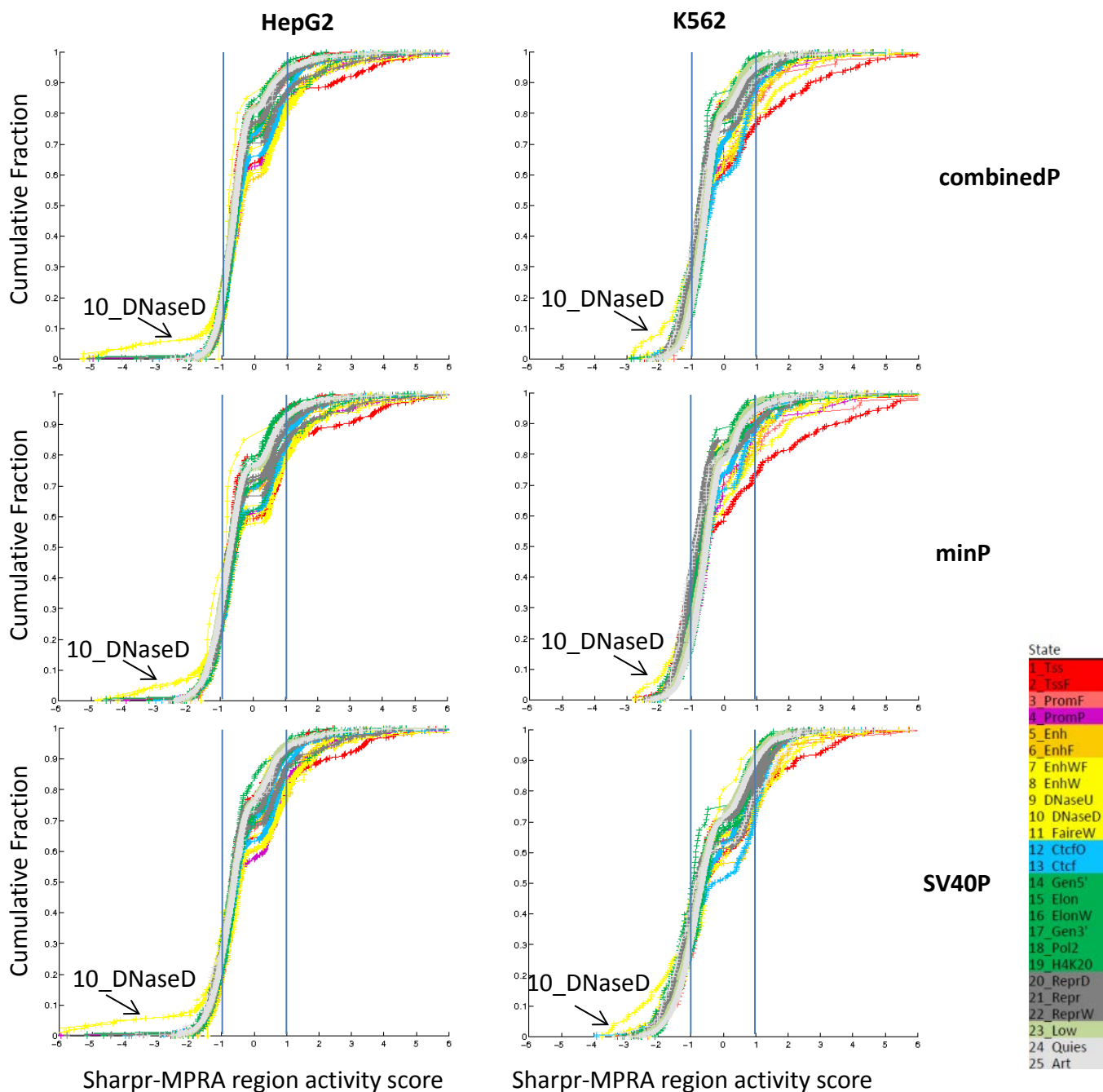
**Supplementary Figure 33 – Cumulative distribution of region activity scores for each chromatin state.** Cumulative fraction of regions (y-axis) showing a MaxPos activity score greater than indicated (on the x-axis) for each chromatin state (colors, see definitions in **Supplementary Fig. 5**) based on the combinedP score (top row), minP score (middle row) and SV40P score (bottom row) for HepG2 (left column) and K562 (right column), and shown for: **(a)** all selected regions from all four cell types that are DNase regions in the tested cell type (N=5840 in HepG2, N=6682 in K562); **(b)**, next page) the subset of DNase regions in the tested cell type that were selected in that cell type; **(c)**, following page) all DNase regions selected in another cell type, that were not in a DNase region in the tested cell type. Each region is shown as a separate point (plus signs). Arrows highlight the active promoter and enhancer states (1\_Tss, 5\_Enh) which had the greatest fraction of regions with strongly activating MaxPos scores, and the single-cut DNase state (10\_DNaseD) which had an increased presence among regions with the most repressive MaxPos scores.

## b. DNase Matched (N=3930 in HepG2, N=3930 in K562)

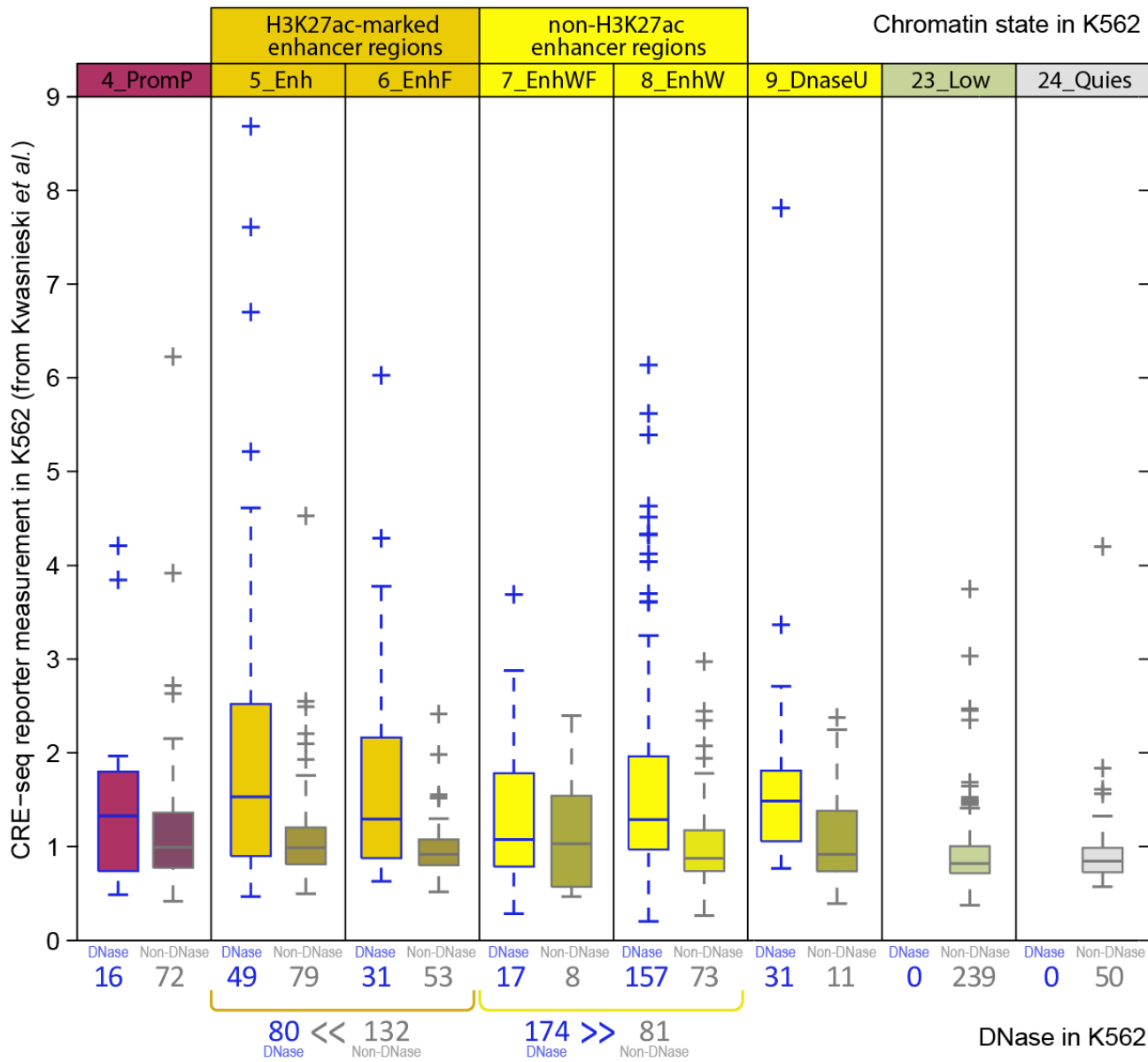


**Supplementary Figure 33b – Cumulative distribution of region activity scores for each chromatin state.** (a) Previous page. (b) Cumulative fraction of regions (y-axis) showing a MaxPos activity score greater than indicated (on the x-axis) for each chromatin state (colors, see definitions in **Supplementary Fig. 5**) based on the combinedP score (top row), minP score (middle row) and SV40P score (bottom row) for HepG2 (left column) and K562 (right column), and shown for the subset of DNase regions in the tested cell type that were selected in that cell type (N=3930 in each of HepG2 and K562); (c) Next page.

### c. Non-DNase (N=9880 in HepG2, N=9038 in K562)

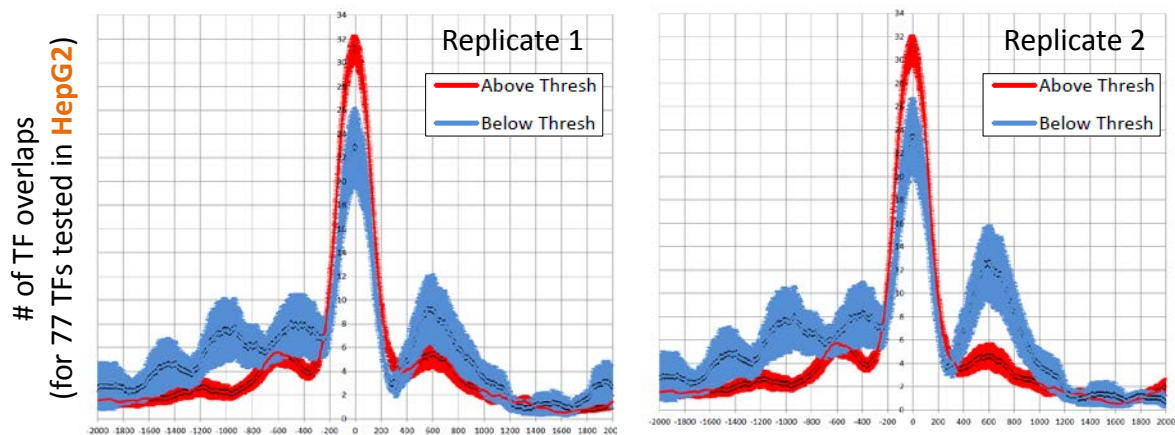


**Supplementary Figure 33c – Cumulative distribution of region activity scores for each chromatin state.** (a,b) Previous pages. (c) Cumulative fraction of regions (y-axis) showing a MaxPos activity score greater than indicated (on the x-axis) for each chromatin state (colors, see definitions in **Supplementary Fig. 5**) based on the combinedP score (top row), minP score (middle row) and SV40P score (bottom row) for HepG2 (left column) and K562 (right column), and shown for all DNase regions selected in another cell type, that were not in a DNase region in the tested cell type (N=9880 in HepG2, N=9038 in K562).

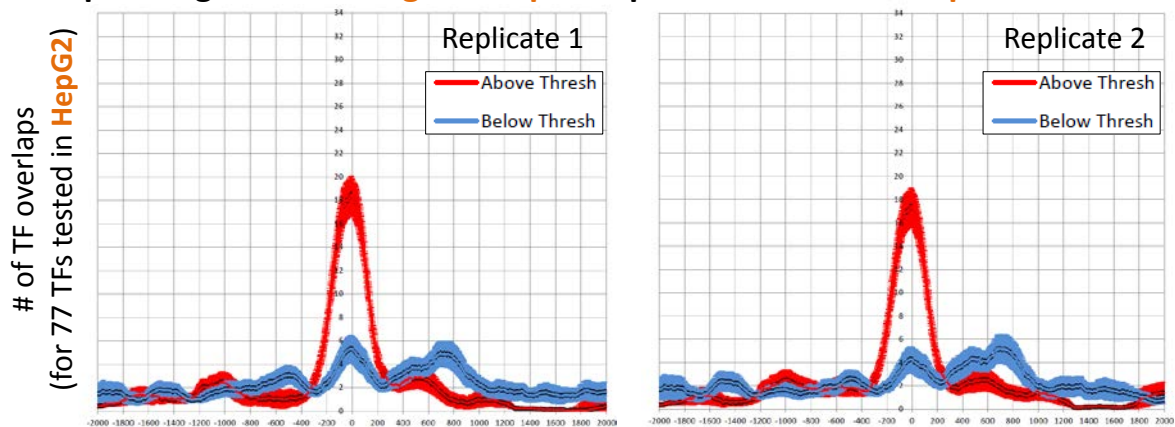


**Supplementary Figure 34 – CRE-seq activity distribution for ChromHMM chromatin states conditioned on presence or absence of DNase.** Average CRE-seq expression data<sup>31</sup> in K562 (y-axis) partitioned by ChromHMM state in K562 (fill color, see **Supplementary Fig. 5**) and by DNase (blue lines) vs. non-DNase (grey lines) based on the center position of tested segments in Ref. 31. The study reported lower CRE-seq activity for H3K27ac-marked enhancer regions, and higher activity for non-H3K27ac enhancer regions, contrary to the previous literature and to our results here and elsewhere<sup>2,4,30,59</sup>. As the figure shows, this reversal in effect is partly explained by the study mixing together two very different classes of regions (DNase vs. non-DNase), whose representation was highly imbalanced between the two groups compared, a statistical effect known as Simpson's Paradox. Only 80 of 212 H3K27ac-marked enhancer regions tested captured DNase sites, compared to 174 of 255 non-H3K27ac enhancer regions. When we separate the tested regions by both chromatin state and by DNase overlap, we find that regions overlapping DNase sites show higher median CRE-seq activity than non-DNase regions within the same chromatin state (as we saw in **Fig. 6b** and **Supplementary Fig. 31b**), and that H3K27ac-marked enhancer regions show higher median activity than non-H3K27ac enhancer regions within DNase regions (as we saw in **Fig. 6a** and **Supplementary Fig. 31a**). This emphasizes the importance of centering tested elements on candidate driver nucleotides (e.g. regulatory motifs, H3K27ac dips, DNase peaks), as we did here and elsewhere<sup>30</sup>. This is especially true when testing short elements without tiling (130-bp in the case of CRE-seq), as positioning the elements without such evidence may exclude driver regulatory nucleotides from tested regions, which can have a strong effect on reporter activity, as we show here. All states with >5 regions tested shown. Boxplots generated with the Matlab boxplot command.

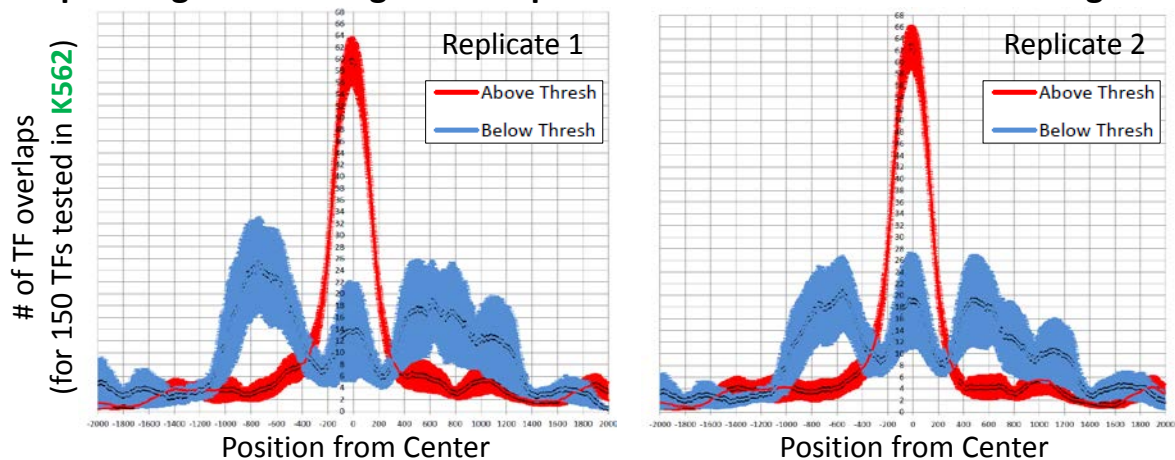
**a. 100 pilot regions with high HepG2 dip scores tested in HepG2 vs. TF binding in HepG2.**



**b. 100 pilot regions at a range of HepG2 dip scores tested in HepG2 vs. TF binding in HepG2.**

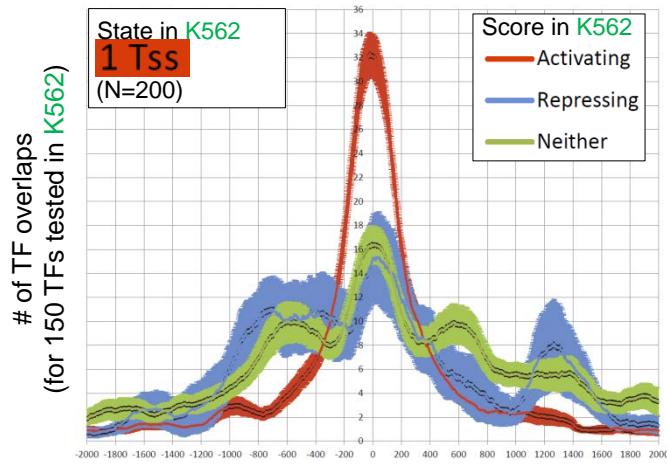
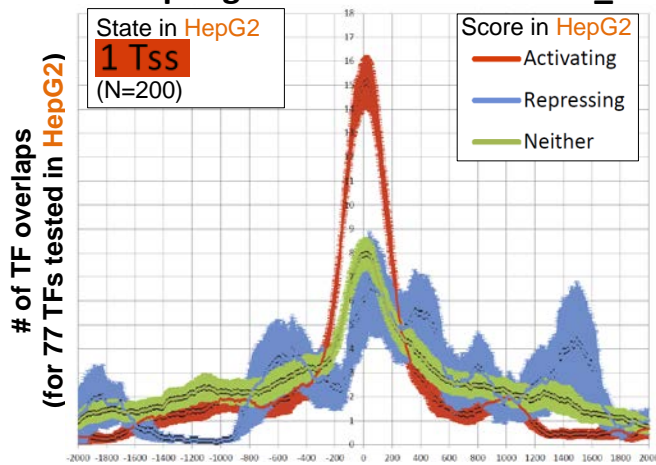


**c. 50 pilot regions with high K562 dip scores tested in K562 vs. TF binding in K562.**

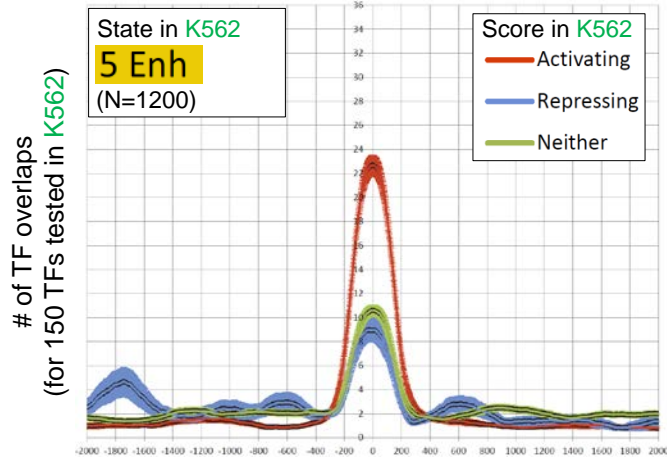
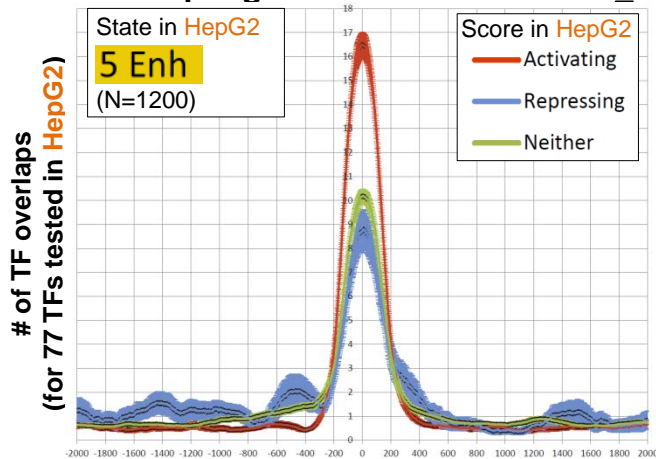


**Supplementary Figure 35 – *In vivo* TF binding is higher and more centrally positioned for higher-scoring regions in pilot experiments. (a,b)** Number of HepG2 ENCODE TF binding ChIP-Seq peak calls out of 77 sets (y-axis) overlapping regions with higher reporter scores (above threshold, red) and lower reporter scores (below-threshold, blue) for each nucleotide position relative to the selected regions center (x-axis) for: **(a)** the 100 regions tested in HepG2 and selected to have the highest *in vivo* dip scores in HepG2; **(b)** and the 100 regions selected to represent a range of dip scores in HepG2 cells. **(c)** The same plot as in (a,b) except for the 50 regions selected based on being in an active enhancer in K562 cells and for the 150 ENCODE peak call sets from K562 cells. Reporter activity threshold selected at the 95<sup>th</sup> percentile of outmost tile offsets. Vertical error bars indicate one standard error. These plots indicate that regions with higher reporter activity show higher TF binding within the tested region (higher red peak in the center), and more concentrated TF binding (lower red peaks in the surrounding) compared to lower reporter activity regions.

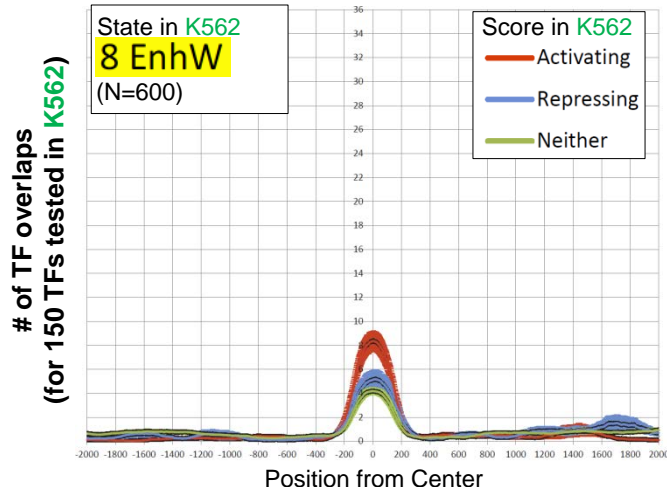
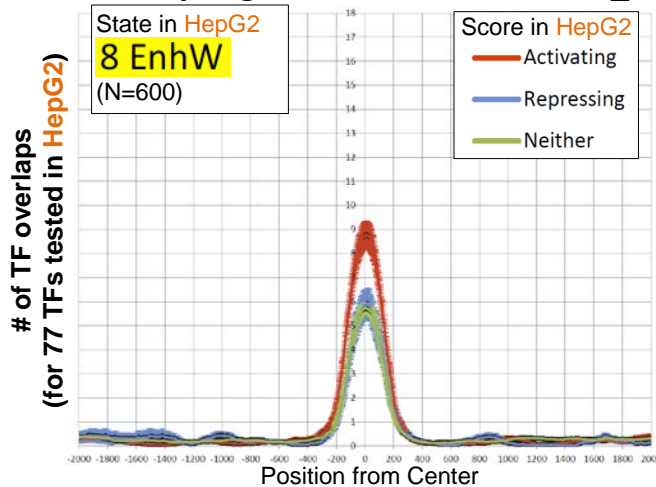
### a. 400 scale-up regions selected to be in 1\_Tss



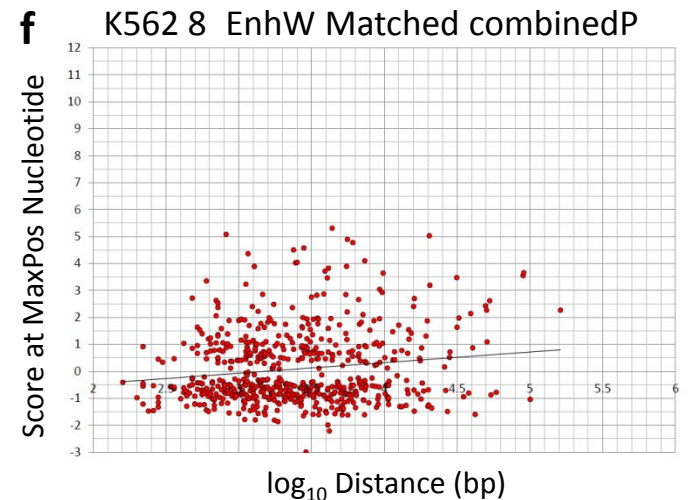
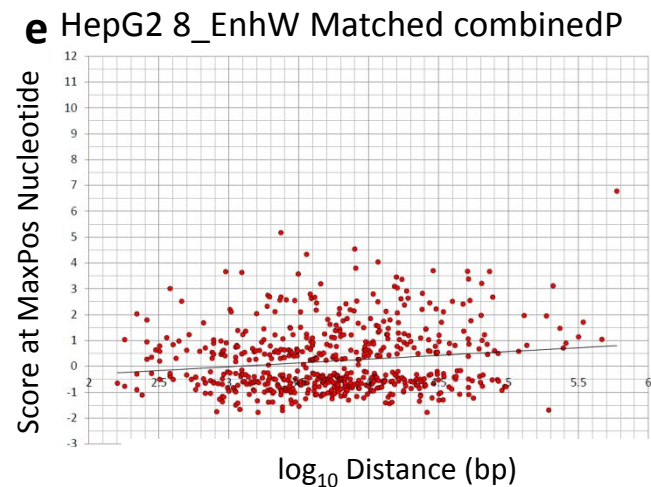
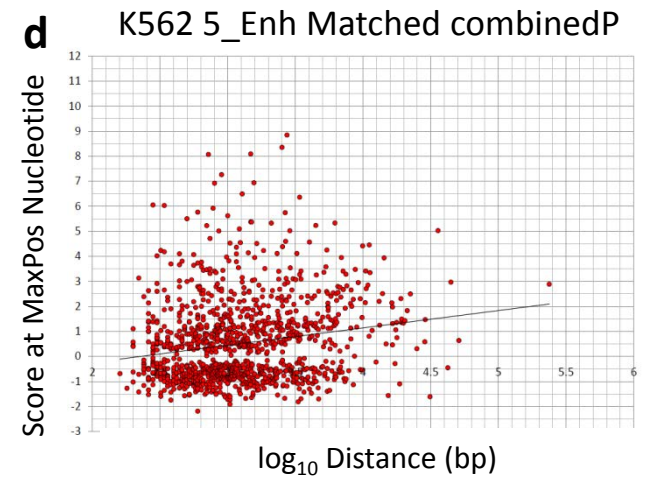
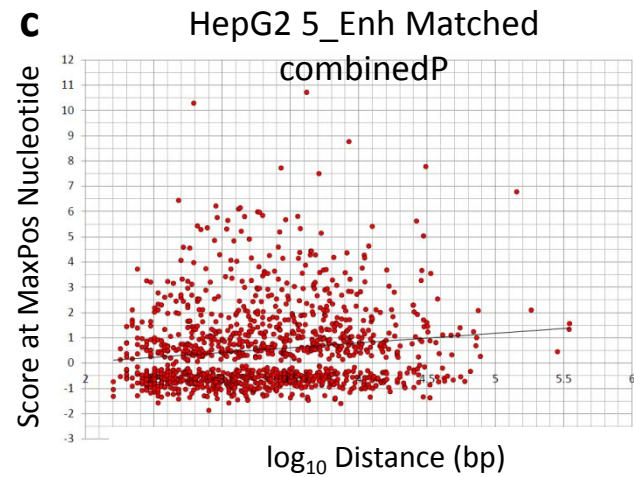
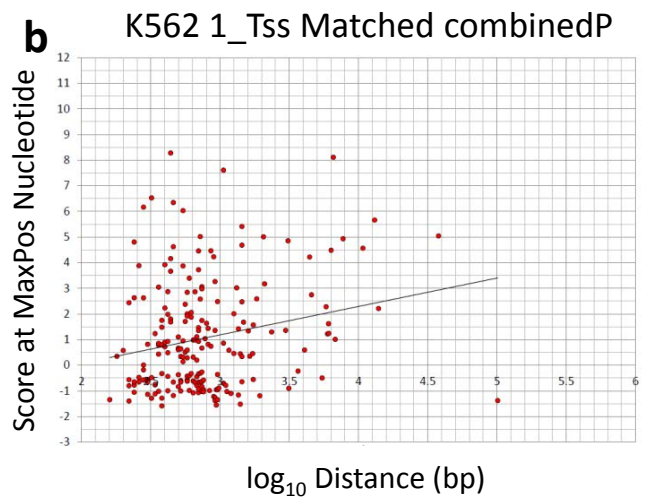
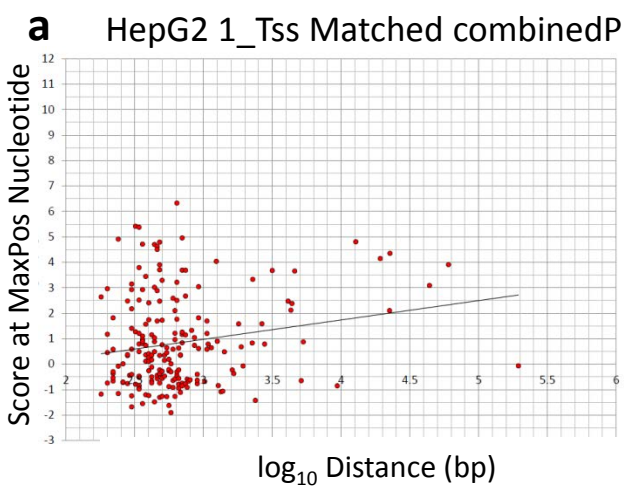
### b. 2400 scale-up regions selected to be in 5\_Enh



### c. 1200 scale-up regions selected to be in 8\_EnhW



**Supplementary Figure 36 – *In vivo* TF binding is higher and more centrally positioned for activating regions in active regulatory chromatin states for scale-up experiments.** Number of HepG2 (left) and K562 (right) ENCODE TF binding ChIP-Seq peak calls (y-axis) overlapping regions with activating regulatory scores (red, combinedP MaxPos score $\geq$ 1), repressive regulatory scores (blue, combinedP MaxPos score $\leq$ -1) or neither (green) for each position relative to the selected regions center (x-axis) for DNase regions selected in the matched cell type where they were tested (DNase matched) for: **(a)** 200 active promoter states (1\_Tss) per cell type; **(b)** 1200 strong enhancer states (5\_Enh); **(c)** 600 weak enhancer states (8\_EnhW). Vertical error bars indicate one standard error. These plots indicate that regions with higher regulatory activity show higher TF binding within the tested region (higher red peak in the center), and more concentrated TF binding (lower red peaks in the surrounding) compared to lower regulatory activity regions.



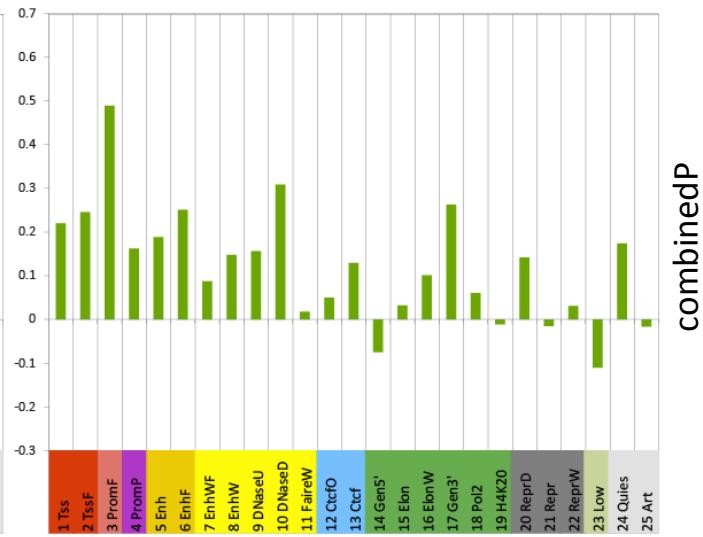
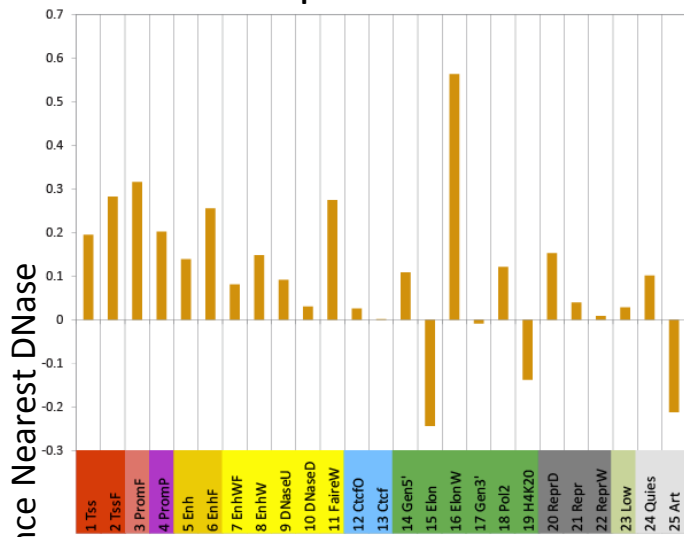
**Supplementary Figure 37 – Scatter Plot Relationship between Region Activity and Distance to Nearest DNase Site for Selected Chromatin States. (a,b)** (a) Scatter plot showing for each tiled region selected based on being in a DNase site in promoter state 1\_Tss in HepG2 cells on the x-axis the  $\log_{10}$  distance in bp to the nearest DNase site in HepG2 (excluding itself) and on the y-axis the Sharpr-MPRA regulatory activity score at the MaxPos nucleotide for the region based on the HepG2 combinedP data. In black is shown a line of best fit for the scatter plots. (b) The same as (a) except for K562 cells. (c,d) The same as (a,b) except for enhancer state 5\_Enh. (e,f) The same as (a,b) except for weak enhancer state 8\_EnhW. The plots show for these states that more isolated DNase sites tend to have greater positive activity. The correlation for all states are shown in **Supplementary Fig. 38**.



## HepG2 Matched

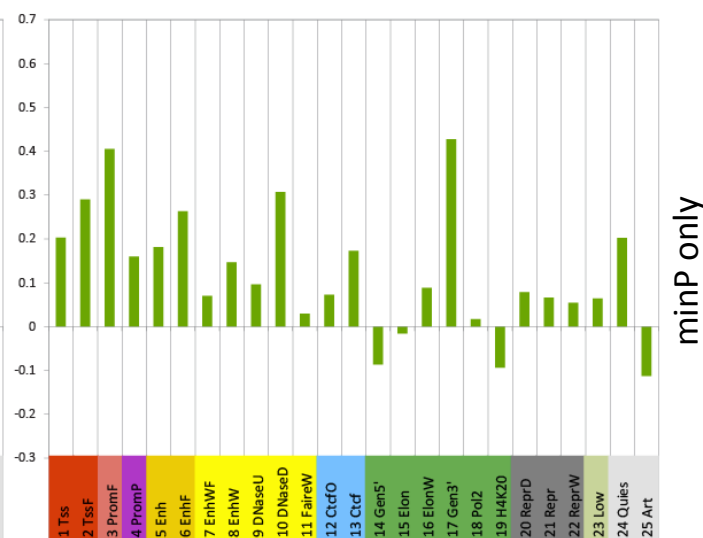
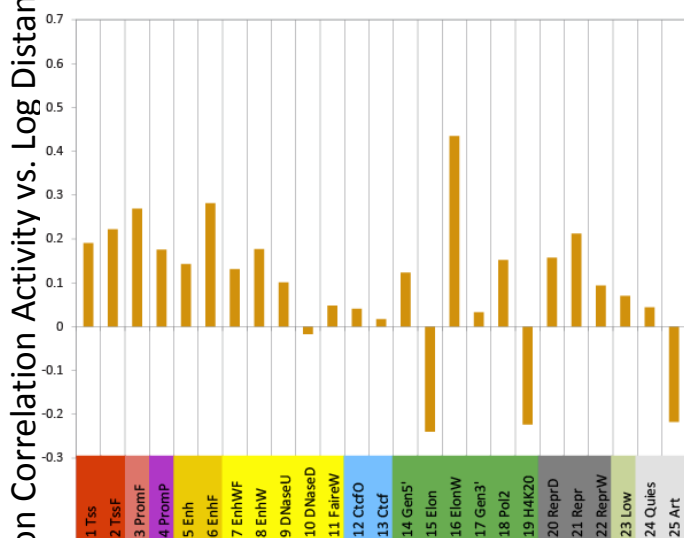
## K562 Matched

a



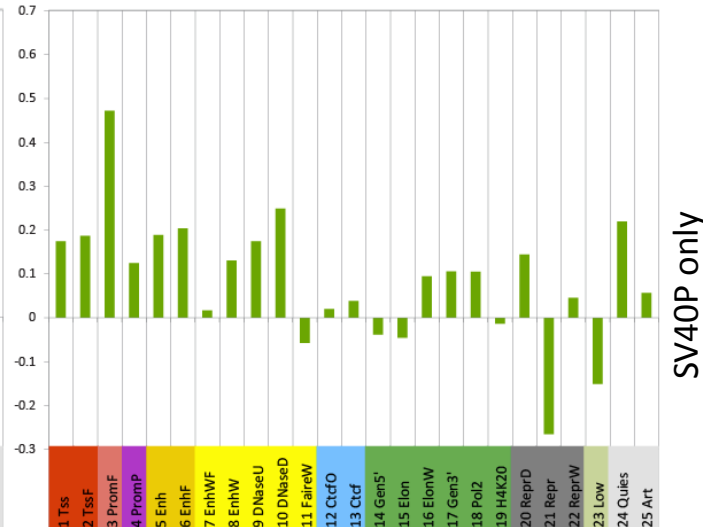
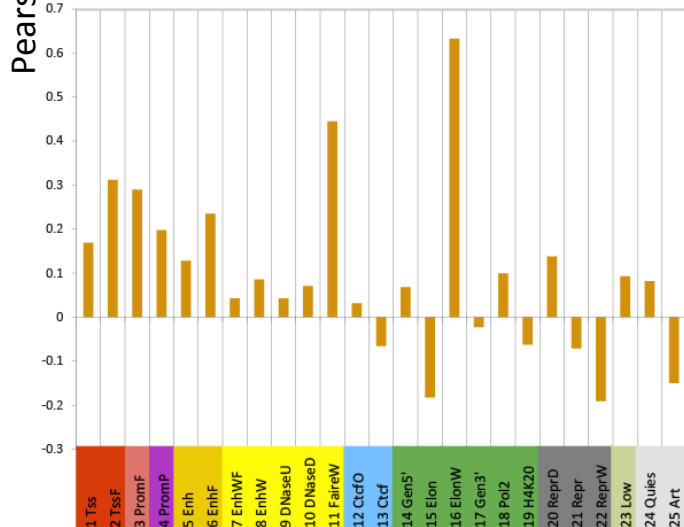
combinedP

b



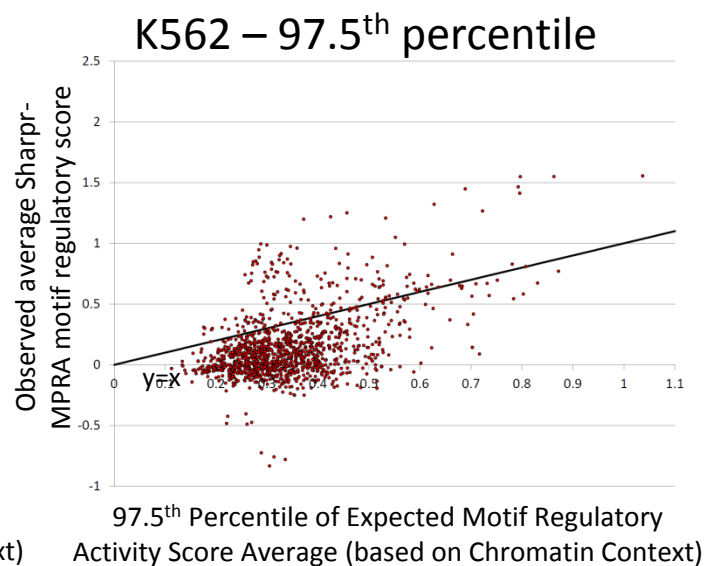
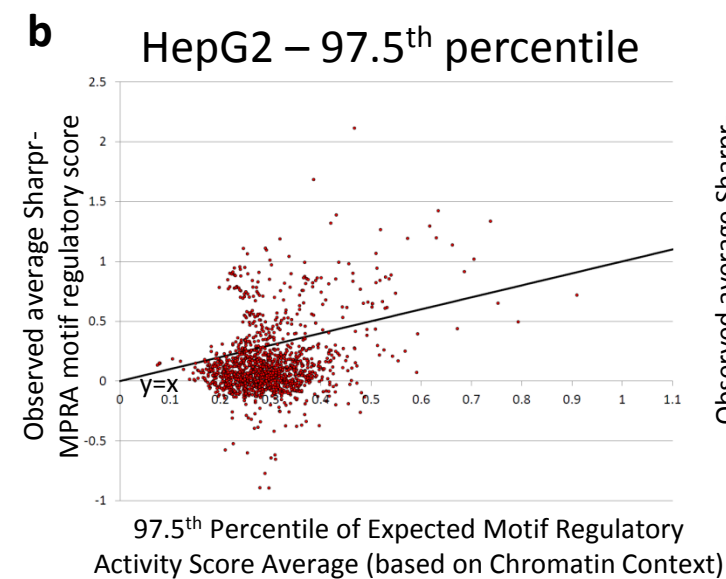
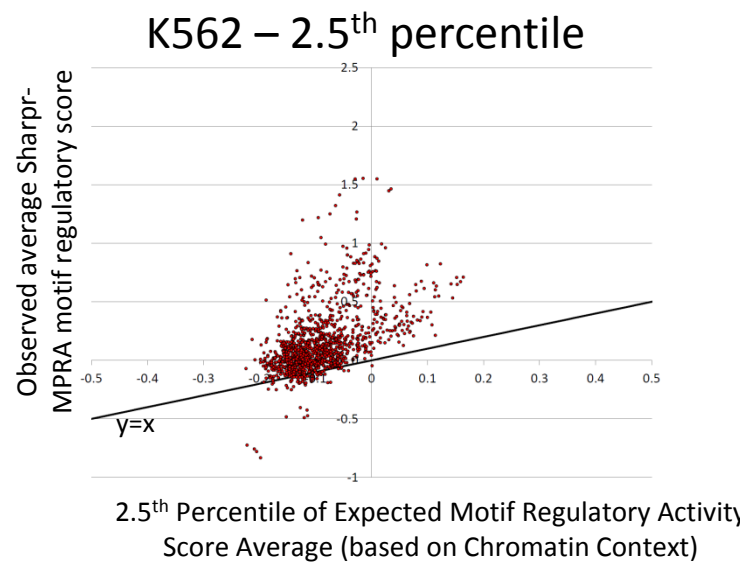
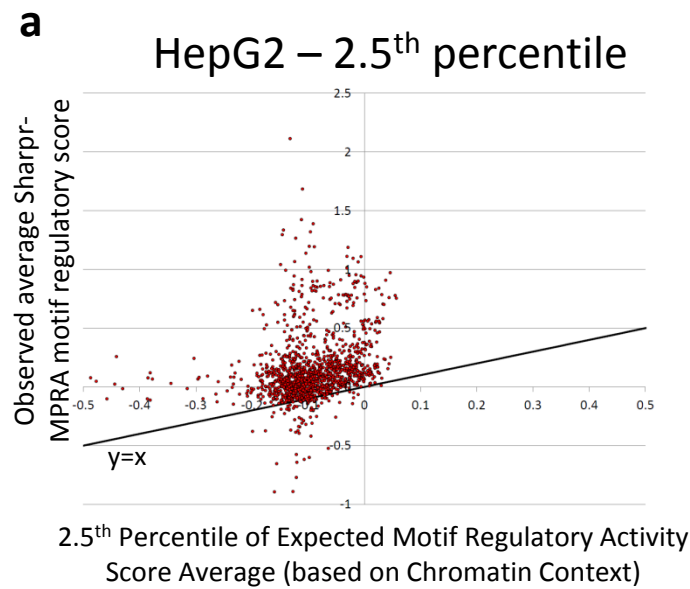
minP only

c



SV40P only

**Supplementary Figure 38 – Correlation of Each Chromatin State for Region Activity and Distance to Nearest DNase Site.** Bar graph showing for the regions selected based on each of the 25 chromatin states the correlation between the log of the distance to the nearest DNase site in that cell type (excluding itself) and the regulatory activity score at the MaxPos nucleotide for HepG2 chromatin states (left) and K562 chromatin states (right), using: **(a)** the combined minP and SV40P score; **(b)** the minP data only; **(c)** the SV40P data only. Full scatter plot for selected states are found in **Supplementary Fig. 37**.



**Supplementary Figure 39 – 95% Confidence Interval Bounds for Expected Motif Regulatory Activity Score Averages based on Chromatin Context. (a)** Scatter plots analogous to Fig. 6c showing the observed average regulatory activity (y-axis) compared with the *lower bound* of a 95% confidence interval (2.5<sup>th</sup> percentile) on the motif regulatory score expected based on the chromatin context (x-axis) for HepG2 (left) and K562 (right). **(b)** Same as (a) except the x-axis represents the *upper bound* on a 95% confidence interval (97.5<sup>th</sup> percentile).

HepG2 PI-conserved elements (AUC 1%/expected)

K562 PI-conserved elements (AUC 1%/expected)

Active Ranking

Active Ranking

	0.1	0.5	1	5	10	50	100
0.1	2.41	2.42	2.43	2.49	2.51	2.60	2.65
0.5		2.42	2.44	2.49	2.52	2.60	2.66
1			2.44	2.49	2.52	2.61	2.67
5				2.52	2.55	2.67	2.71
10					2.56	2.70	2.73
50						2.73	2.71
100							2.69

	0.1	0.5	1	5	10	50	100
0.1	2.41	2.42	2.42	2.45	2.47	2.53	2.54
0.5		2.42	2.42	2.45	2.48	2.54	2.54
1			2.43	2.45	2.48	2.55	2.55
5				2.47	2.50	2.57	2.57
10					2.51	2.59	2.58
50						2.62	2.60
100							2.59

Repressive Ranking

Repressive Ranking

	0.1	0.5	1	5	10	50	100
0.1	1.69	1.69	1.68	1.66	1.63	1.54	1.48
0.5		1.69	1.68	1.66	1.62	1.54	1.48
1			1.68	1.65	1.62	1.53	1.48
5				1.63	1.60	1.52	1.46
10					1.59	1.48	1.40
50						1.31	1.24
100							1.21

	0.1	0.5	1	5	10	50	100
0.1	1.43	1.43	1.42	1.41	1.38	1.24	1.17
0.5		1.42	1.42	1.41	1.38	1.24	1.16
1			1.42	1.40	1.37	1.23	1.16
5				1.37	1.33	1.21	1.13
10					1.31	1.17	1.10
50						1.10	1.04
100							0.99

Absolute Ranking

Absolute Ranking

	0.1	0.5	1	5	10	50	100
0.1	2.46	2.47	2.48	2.55	2.59	2.72	2.77
0.5		2.47	2.48	2.56	2.60	2.72	2.78
1			2.49	2.56	2.60	2.73	2.79
5				2.58	2.63	2.78	2.82
10					2.64	2.81	2.81
50						2.78	2.67
100							2.61

	0.1	0.5	1	5	10	50	100
0.1	2.42	2.43	2.43	2.46	2.48	2.54	2.55
0.5		2.43	2.43	2.46	2.48	2.55	2.56
1			2.43	2.46	2.49	2.55	2.56
5				2.48	2.50	2.58	2.57
10					2.52	2.59	2.59
50						2.60	2.53
100							2.49

HepG2 PI-conserved elements (AUC 5%/expected)

K562 PI-conserved elements (AUC 5%/expected)

Active Ranking

Active Ranking

	0.1	0.5	1	5	10	50	100
0.1	1.85	1.86	1.86	1.88	1.90	1.94	1.93
0.5		1.86	1.86	1.89	1.90	1.94	1.93
1			1.87	1.89	1.90	1.94	1.93
5				1.91	1.92	1.95	1.94
10					1.93	1.95	1.93
50						1.93	1.88
100							1.85

	0.1	0.5	1	5	10	50	100
0.1	1.82	1.82	1.83	1.84	1.86	1.88	1.88
0.5		1.83	1.83	1.85	1.86	1.88	1.88
1			1.83	1.85	1.86	1.89	1.88
5				1.86	1.87	1.89	1.88
10					1.88	1.89	1.87
50						1.86	1.81
100							1.77

Repressive Ranking

Repressive Ranking

	0.1	0.5	1	5	10	50	100
0.1	1.28	1.28	1.28	1.27	1.26	1.21	1.18
0.5		1.28	1.28	1.27	1.26	1.21	1.18
1			1.28	1.27	1.26	1.20	1.18
5				1.26	1.25	1.18	1.15
10					1.23	1.15	1.13
50						1.08	1.05
100							1.03

	0.1	0.5	1	5	10	50	100
0.1	1.18	1.18	1.18	1.16	1.15	1.09	1.06
0.5		1.18	1.18	1.16	1.15	1.09	1.06
1			1.18	1.16	1.15	1.08	1.06
5				1.14	1.12	1.06	1.03
10					1.11	1.04	1.02
50						0.99	0.97
100							0.96

Absolute Ranking

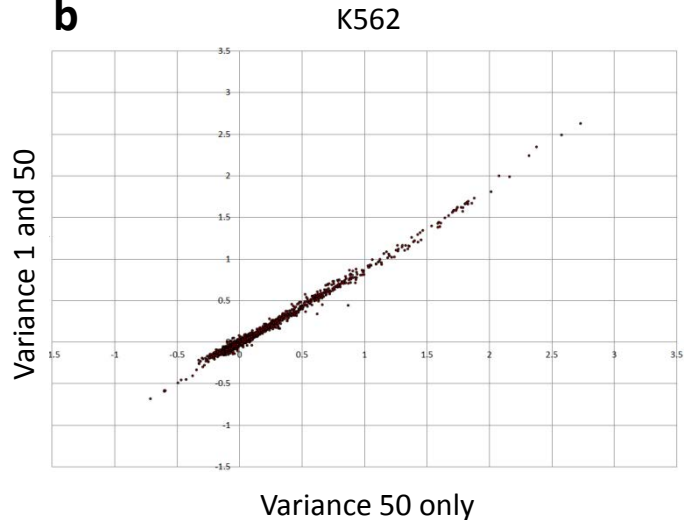
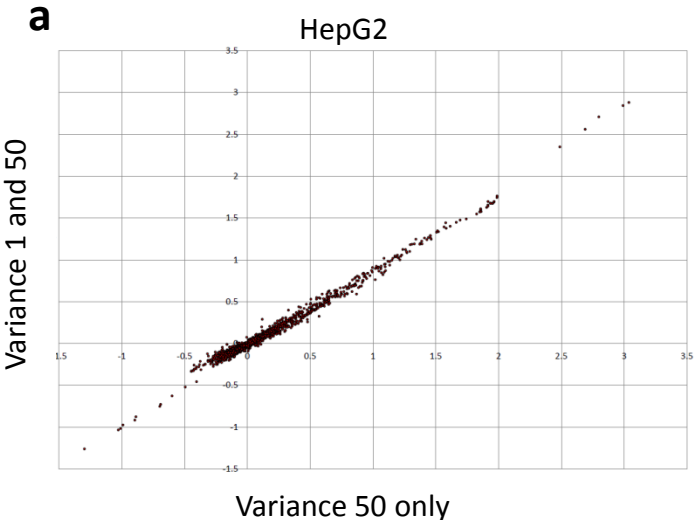
Absolute Ranking

	0.1	0.5	1	5	10	50	100
0.1	1.88	1.89	1.89	1.91	1.92	1.95	1.93
0.5		1.89	1.89	1.91	1.92	1.95	1.93
1			1.90	1.92	1.93	1.95	1.93
5				1.93	1.95	1.96	1.92
10					1.95	1.95	1.90
50						1.88	1.78
100							1.75

	0.1	0.5	1	5	10	50	100
0.1	1.82	1.82	1.83	1.84	1.85	1.86	1.83
0.5		1.83	1.83	1.85	1.85	1.86	1.83
1			1.83	1.85	1.86	1.86	1.83
5				1.86	1.86	1.86	1.82
10					1.87	1.86	1.81
50						1.80	1.72
100							1.67

### Supplementary Figure 40 – Conserved Element Recovery as a Function of Variance Prior

**(a)** The table shows for each indicated combination of the variance prior parameters for the regulatory activity variables (see Methods), when ranking the resulting nucleotides by how activating their regulatory activity score is, the ratio of the Area under the ROC curve (AUC) up to a 1% false positive rate for recovering SiPhy-PI conserved elements<sup>42,65</sup> to what would be expected by random guessing (0.00005) for HepG2 (left) and K562 (right). **(b)** The same as (a) except based on ranking nucleotides by how repressive their regulatory score is. **(c)** The same as (a) and (b) except ranking nucleotides based on the absolute value of their regulatory score. **(d-f)** The same as (a-c) except based on an AUC to a 5% false positive rate, where the expected AUC to a 5% false positive rate by random guessing is 0.00125. The color shading scale is specific to each heatmap.



**Supplementary Figure 41 – Effect on motif average Sharpr-MPRA scores when using just the higher-variance parameter.** Scatter plots for **(a)** HepG2 and **(b)** K562 cells comparing the average motif regulatory score using (x-axis) just the larger variance parameter (50) and (y-axis) using the combined results from two variance parameters (1,50) as reported in **Fig. 4a**. Each point corresponds to one motif. Only motifs with >10 instances are shown. In both plots the correlation of the values are >0.99.

## Supplementary Note

Here we derive the equations for the terms of  $\Sigma_{M_{r,t},M_{r,t}}$  and  $\Sigma_{A_{r,t},M_{r,t}}$  presented in the Methods.

For  $\Sigma_{M_{r,t},M_{r,t}}$  we have

$$\begin{aligned}\Sigma_{M_{r,t,u},M_{r,t,u}} &= \sigma_{m_r}^2 + \sigma_a^2/N \\ \Sigma_{M_{r,t,u},M_{r,t,v}} &= \frac{\sigma_a^2(N - |u - v|)}{N^2} \text{ if } (0 < |u - v| < N) \\ \Sigma_{M_{r,t,u},M_{r,t,v}} &= 0 \text{ if } (|u - v| \geq N)\end{aligned}$$

The expression for  $\Sigma_{M_{r,t,u},M_{r,t,u}}$  above follows from observing

$$\text{var}(M_{r,t,u}) = \text{var}\left(\sigma_{m_r}Z_{r,t,u} + \left(\frac{1}{N}\sum_{l=0}^{N-1}A_{r,t,u+l}\right)\right) = \sigma_{m_r}^2 + \frac{1}{N^2}\sum_{l=0}^{N-1}\text{var}(A_{r,t,u+l}) = \sigma_{m_r}^2 + \sigma_a^2/N$$

where  $Z_{r,t,u}$  represents a standard normal random variable.

The expression for  $\Sigma_{M_{r,t,u},M_{r,t,v}}$  follow since

$$\begin{aligned}\text{cov}(M_{r,t,u},M_{r,t,v}) &= \text{cov}\left(\sigma_{m_r}Z_{r,t,u} + \frac{1}{N}\sum_{l=0}^{N-1}A_{r,t,u+l}, \sigma_{m_r}Z_{r,t,v} + \frac{1}{N}\sum_{l=0}^{N-1}A_{r,t,v+l}\right) \\ &= \text{cov}\left(\frac{1}{N}\sum_{l=0}^{N-1}A_{r,t,u+l}, \frac{1}{N}\sum_{l=0}^{N-1}A_{r,t,v+l}\right)\end{aligned}$$

where  $Z_{r,t,u}$  and  $Z_{r,t,v}$  represent independent standard normal random variables. The expression

is 0 if  $(|u - v| \geq N)$  and otherwise is

$$\frac{1}{N^2}\sum_{l=|u-v|}^{N-1}\text{var}(A_{r,t,(\min(u,v)+l)}) = \frac{\sigma_a^2(N - |u - v|)}{N^2}$$

For  $\Sigma_{A_{r,t}, M_{r,t}}$  we have:

$$\Sigma_{A_{r,t,k}, M_{r,t,u}} = \frac{\sigma_a^2}{N} \text{ if } u \leq k < u + N$$

$$\Sigma_{A_{r,t,k}, M_{r,t,u}} = 0 \text{ otherwise}$$

This follows from observing

$$\begin{aligned} \text{cov}(A_{r,t,k}, M_{r,t,u}) &= \text{cov}\left(A_{r,t,k}, \frac{1}{N} \sum_{l=0}^{N-1} A_{r,t,u+l}\right) \\ &= \frac{1}{N} \sum_{l=0}^{N-1} \text{cov}(A_{r,t,k}, A_{r,t,u+l}) \end{aligned}$$

All terms are 0 except if  $u \leq k < u + N$  in which case the expression reduces to

$$= \frac{1}{N} \text{var}(A_{r,t,k})$$