

# Supplementary data

## Contents

<b>Supplementary methods</b>	<b>2</b>
Method Overview . . . . .	2
Motif models curation . . . . .	2
Detecting motif-associated footprint patterns from chromatin accessibility data . . . . .	3
Profiling motif-associated sequence features . . . . .	4
Sparse logistic regression training (MocapS) . . . . .	4
Implementation of cross-sample TFBS prediction (MocapX) . . . . .	5

## Supplemental Figures

Figure S1	Filtering for highly redundant motifs representing the same TF. . . . .	8
Figure S2	ChIP-Seq peak coverage by motifs. . . . .	10
Figure S3	Cut count density distribution across experiment types. . . . .	11
Figure S4	Features correlations with ChIP-Seq signals. . . . .	11
Figure S5	Training sparse logistic regression classifiers by bootstrap aggregation. . . . .	13
Figure S6	Footprint profiles vary across factors, cell types and experimental protocols. . . . .	14
Figure S7	Performance comparison between methods. . . . .	15

## Supplemental Tables

Table S1	Performance comparison for DGF samples . . . . .	16
Table S2	Performance comparison for ATAC-Seq samples . . . . .	17

## Appendices

Appendix A	Footprint profile plots . . . . .	18
Appendix B	Sparse logistic regression training plots . . . . .	38

# Supplementary methods

## Method Overview

Mocap combines information from motif, chromatin accessibility and a range of sequence-associated genomic features to predict TFBS. We start with a liberal set of candidate motif sites, and create a generic unsupervised method (MocapG) that ranks motif sites based on local accessibility. We then, building on this baseline method, adopt a supervised machine learning scheme (training sparse logistic regression classifiers MocapS) to investigate the contribution of different features to TFBS prediction and probe the diversity in TF binding behaviors. In training MocapS, we model both cell type-specific chromatin accessibility and footprint patterns around motif sites. We investigate (for different types of TFs) the contribution of sequence properties, such as GC/CpG content, surrounding motif sites to TFBS prediction. We use L1-regularized logistic regression to integrate motif-associated genomic features, which we find to be a good balance of performance and interpretability (in comparison to methods such as support vector machines, prior published methods, and methods adopting other regularization schemes). Lastly, we extend the trained classifiers of MocapS to TFCT conditions without ChIP-Seq label in a data-driven manner (MocapX). The goal is to create a TFBS prediction pipeline that balances method sensitivity with selectivity, and be able to generalize the specifically trained models to more factors and cell types via cross-sample TFBS predictions.

Below we detail the curation of motif models extracted from public databases, the computation of motif-associated genomic features, and the implementation of MocapS and MocapX.

## Motif models curation

We drew PWMs from two motif collections (Figure S1A). To reduce information redundancy for TFs with multiple motif representations, we filtered motifs based on pairwise PWM comparisons. Briefly, for a given motif pair, we vectorize the PWMs and calculate length-normalized Euclidean distance between the two vectors as follow

$$d_2 = \sqrt{\sum_{i=1}^v \frac{(p_{1i} - p_{2i})^2}{v}}$$

where  $p_1$  and  $p_2$  represent the length-normalized PWM vectors,  $v$  denotes the length of the vectorized PWMs. For motif pairs of unequal length, we let  $v = \min(p_1, p_2)$  where we used the alignment that results in the smallest Euclidean distance among all possible ungapped PWM alignments. For motif pairs that are highly similar (with normalized Euclidean distance  $< 0.14$ ), the PWM of lower entropy was retained (Figure S1B-D).

## Detecting motif-associated footprint patterns from chromatin accessibility data

To assess the probability that a footprint profile ( $F$ ) exists around a motif site, we used a pair of binomial tests to gauge whether the motif region is protected from enzyme digestion as compared to its left and right flanking regions (4).

For DNase-Seq, we scored for the depletion of DNase I cuts at motif sites in a strand-specific manner, as strand imbalance is often conspicuous because of the size selection steps in DNase-Seq protocol (Figure S6). Specifically, on the positive strand, we assessed if a motif region is depleted of cuts as compared to the left flanking region.

$$P(c_m^+; p, n^+) = \sum_{i=0}^{c_m^+} \binom{n^+}{i} p^i (1-p)^{n^+-i}$$

On the negative strand, we assessed if a motif region is depleted of cuts as compared to the right flanking region.

$$P(c_m^-; p, n^-) = \sum_{i=0}^{c_m^-} \binom{n^-}{i} p^i (1-p)^{n^- - i}$$

where  $p = \frac{l_m}{l_m + l_f}$ .  $p$  denotes the expected size-normalized cut probability, where  $l_m$  and  $l_f$  denotes the size of the motif and its flanking region respectively. We empirically set the size of flanking region as 1.75 the size of motif length to optimize ChIP-Seq peak recovery.  $n^+ = c_m^+ + c_f^+$  is the total number of cut counts within the motif site  $c_m^+$  and its left flanking region  $c_f^+$  on the positive strand, whereas  $n^- = c_m^- + c_f^-$  is the total number of cuts within the motif site  $c_m^-$  and its right flanking region  $c_f^-$  on the negative strand.

For ATAC-Seq, as the strand imbalance immediately flanking binding motif is less pronounced, we merged strand information to assess if a motif region is depleted of cuts as compared to the left and right flanking region respectively.

$$P(c_m; p, n^L) = \sum_{i=0}^{c_m} \binom{n^L}{i} p^i (1-p)^{n^L - i}$$

$$P(c_m; p, n^R) = \sum_{i=0}^{c_m} \binom{n^R}{i} p^i (1-p)^{n^R - i}$$

where  $n^L = c_m + c_f^L$  is the total number of cuts within the motif site  $c_m$  and its left flanking region  $c_f^L$  on both strands, whereas  $n^R = c_m + c_f^R$  is the total number of cuts within the motif site  $c_m$  and its right flanking region  $c_f^R$  on both strands.

The footprint score was calculated for the DNase I footprint as

$$F = P(c_m^+; p, n^+) \cdot P(c_m^-; p, n^-)$$

and for the ATAC-Seq footprint as

$$F = P(c_m; p, n^L) \cdot P(c_m; p, n^R)$$

We shuffled the per base-pair cut count data around each site and reran the same binomial tests 500 times. A  $p$ -value threshold corresponding to a false discovery rate of 0.01 was chosen as the threshold for determining whether a footprint profile exists for each motif site.

## Profiling motif-associated sequence features

To profile the sequence environment surrounding TF binding motifs, we calculated GC content and CpG frequency for a region 100bp upstream and downstream of each motif site. We also computed a binary CpG island feature to assess if a motif resides in a larger region of inflated CpG count. A CpG island was assessed with the following criteria (2, 5)

$$CGI(N_{CpG}, N_C, N_G, L) = \begin{cases} 2 & \text{if } \frac{N_{CpG}L}{((N_C + N_G)/2)^2} > 0.6 \ \& \ \frac{N_C + N_G}{L} > 0.5 \\ 1 & \text{otherwise} \end{cases}$$

where  $L$  is the size of the region encompassing 400bp upstream and 400bp downstream of the length of a motif,  $N_C$ ,  $N_G$  and  $N_{CpG}$  are the total count of C, G and CpG dinucleotide respectively.

## Sparse logistic regression training (MocapS)

We trained sparse logistic regression models to classify candidate motif sites as true (motif overlapping with ChIP-Seq) or false (motif not overlapping with ChIP-Seq peaks) binding sites using a range of motif-associated genomic features, including motif matching quality, chromatin accessibility, TF footprint, TSS proximity, evolutionary conservation, motif region mapability and sequence-associated GC/CpG features. The goal is to build an ensemble of TFCT-specific and yet generalizable predictive models that address the differences in TF binding preferences and facilitate cross-sample binding site prediction. The logistic regression model is specified as below:

$$p(y = 1|x; \theta) = \sigma(\theta^T x) = \frac{1}{1 + e^{-(\theta^T x + c)}}$$

where  $x$  represents motif samples,  $y \in [1, -1]$  corresponds to bound and unbound motif labels respectively;  $\theta$  specifies the 30 motif-associated genomic features that are mostly highly correlated with the ChIP-Seq

feature in the tested motif set and cell type environment (after two rounds of correlation filtering);  $c$  specifies the intercept. Features, such as SNP density, TES proximity and exon, have relatively low correlation with ChIP-Seq signals, and a weak contribution to the overall prediction of TFBS, so we excluded them from the model *a priori* for all TFCTs.

To create generalizable models and ensure only a most relevant subset of the genomic features are included in the logistic regression model for each TFCT, we specified a sparsity constraint to minimize the below objective function with L1 regularization

$$\min_{\theta} \sum_{i=1}^M -\log p(y^{(i)}|x^{(i)}; \theta) + \lambda \|\theta\|_1$$

where  $M$  is the number of candidate motif sites used for model training;  $\lambda$  is the shrinkage parameter corresponding to the sparsity of the trained logistic regression model where with L1-regularization, all but the most relevant features are shrink to zero. As the number of true binding sites takes up only 0.01%-5% of the total number of candidate sites (imbalanced class labels), we weighed the true and false motif classes to create a weight-balanced training dataset. For each TFCT condition, we tuned the shrinkage parameter  $\lambda$  to select a  $\lambda$  that optimizes the AUPR to achieve more sensible site ranking (also, AUPR scores are less biased than AUROC scores in assessing classification performance in imbalanced datasets). To avoid overfitting and promote model generalization,  $\lambda$  parameter was further tuned such that the resulting classifier is sparser than an optimally performing classifier with near-optimum (within one standard error of the optimum) cross-validation performance. We adopted a final bootstrap aggregation step to reduce estimator bias and help achieve more robust classifications.

## Implementation of cross-sample TFBS prediction (MocapX)

To extend the utility of the above sparse logistic regression classifiers ensemble to TFCT conditions where ChIP-Seq data is missing, we used robustly weighted least squares regression to derive a mapping vector  $\beta$  to assign new TFCT prediction problems to trained classifiers based on genomic feature distance. When minimizing the below error function

$$\min \sum_{i=1}^n \sum_{j=i+1}^n w_{ij} \left( \gamma \log \frac{Y_{ij} + Y_{ji}}{2} - \sum_{k=1}^p X_{ij,k} \beta_k \right)^2 = \min \sum_{i=1}^n \sum_{j=i+1}^n w_{ij} (e_{ij})^2$$

where  $n$  is the number of TFCT conditions;  $p$  is the number of features;  $X_{ij,k}$  is the binary probability that feature  $k$  from TFCT  $i$  and TFCT  $j$  are derived from the same distribution based on KS test for continuous features and Chi-square test for categorical features;  $Y_{ij}$  is the predictive performance (AUPR) of model

trained in TFCT  $i$  when applied to classify motifs in TFCT  $j$ . We used predictive performance of MocalpG as a baseline and used adaptive grid search to assign  $\gamma$  to separate sample pairs that cross-predict well from those that predicts poorly. For DNase-Seq sample mappings, for example, we set  $\gamma$  as follows

$$\gamma = \begin{cases} 1 & \frac{Y_{ij}+Y_{ji}}{2Y_{MocalpG}} > 1.1 \\ 1.31 & \frac{Y_{ij}+Y_{ji}}{2Y_{MocalpG}} \leq 1.1 \end{cases}$$

$\gamma$  was re-optimized with the exclusion of footprint features, when deriving mappings between ATAC-Seq samples and DNase-Seq data trained models.  $w$  was fitted via iteratively re-weighted least squares regression with the Tukey’s bisquare family psi functions as below (1, 3).

$$w(u) = \begin{cases} ((1 - (\frac{u}{c})^2))^2 & |u| \leq c \\ 0 & |u| > c \end{cases}$$

where  $u$  is the regression residual.

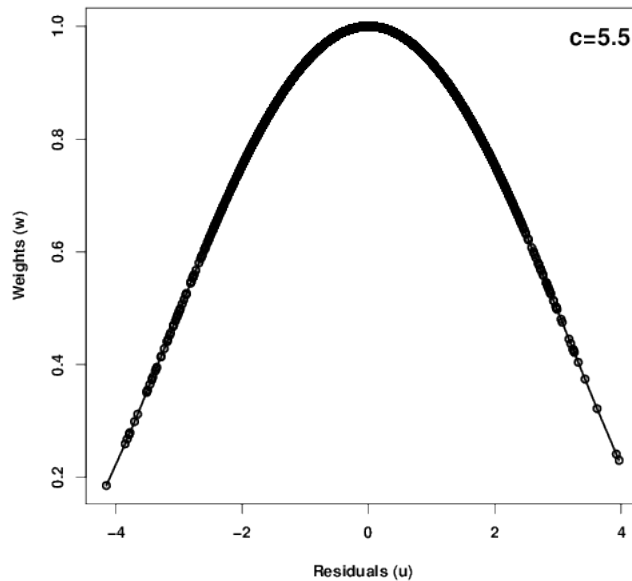


Figure 1: Tukey’s bisquare function specifies sample weights based on regression residuals.

We selected the bisquare psi function hyperparameter  $c$  based on cross-validation sample-mapping performance. Distances between feature vectors were then used to assign weight  $w$  to a hold-out TFCT condition based on the weights of its nearest neighbors. We used  $\hat{w}$  and  $\hat{\beta}$  to compute a sample mapping  $i \rightarrow j$  that max-

minizes  $\hat{Y} = \arg \max_{i \rightarrow j} \hat{w}_{ij} \sum_{k=1}^p X_{ij,k} \hat{\beta}_k$ , but to avoid mappings that do not significantly improve performance as compared to the generic method MocalpG, we assign the mapping if and only if the estimated cross-prediction performance  $\hat{Y}$  compares favorably with  $Y_{MocalpG}$  in TFCT  $i$  as below

$$Y_{MocalpX}^{(j)} = \begin{cases} Y_{MocalpG}^{(j)} & \hat{Y} < Y_{MocalpG}^{(i)} \\ Y_{MocalpS_i}^{(j)} & \hat{Y} \geq Y_{MocalpG}^{(i)} \end{cases}$$

## References

1. Robert Andersen. *Modern methods for robust regression*. Number 152. Sage, 2008.
2. M Gardiner-Garden and M Frommer. CpG islands in vertebrate genomes. *Journal of molecular biology*, 196(2):261–282, 1987.
3. Peter J Huber. Wiley series in probability and mathematics statistics. *Robust Statistics*, pages 309–312, 1981.
4. Jason Piper, Markus C Elze, Pierre Cauchy, Peter N Cockerill, Constanze Bonifer, and Sascha Ott. Wellington: a novel method for the accurate identification of digital genomic footprints from dnase-seq data. *Nucleic acids research*, page gkt850, 2013.
5. Serge Saxonov, Paul Berg, and Douglas L Brutlag. A genome-wide analysis of cpg dinucleotides in the human genome distinguishes two distinct classes of promoters. *Proceedings of the National Academy of Sciences of the United States of America*, 103(5):1412–1417, 2006.

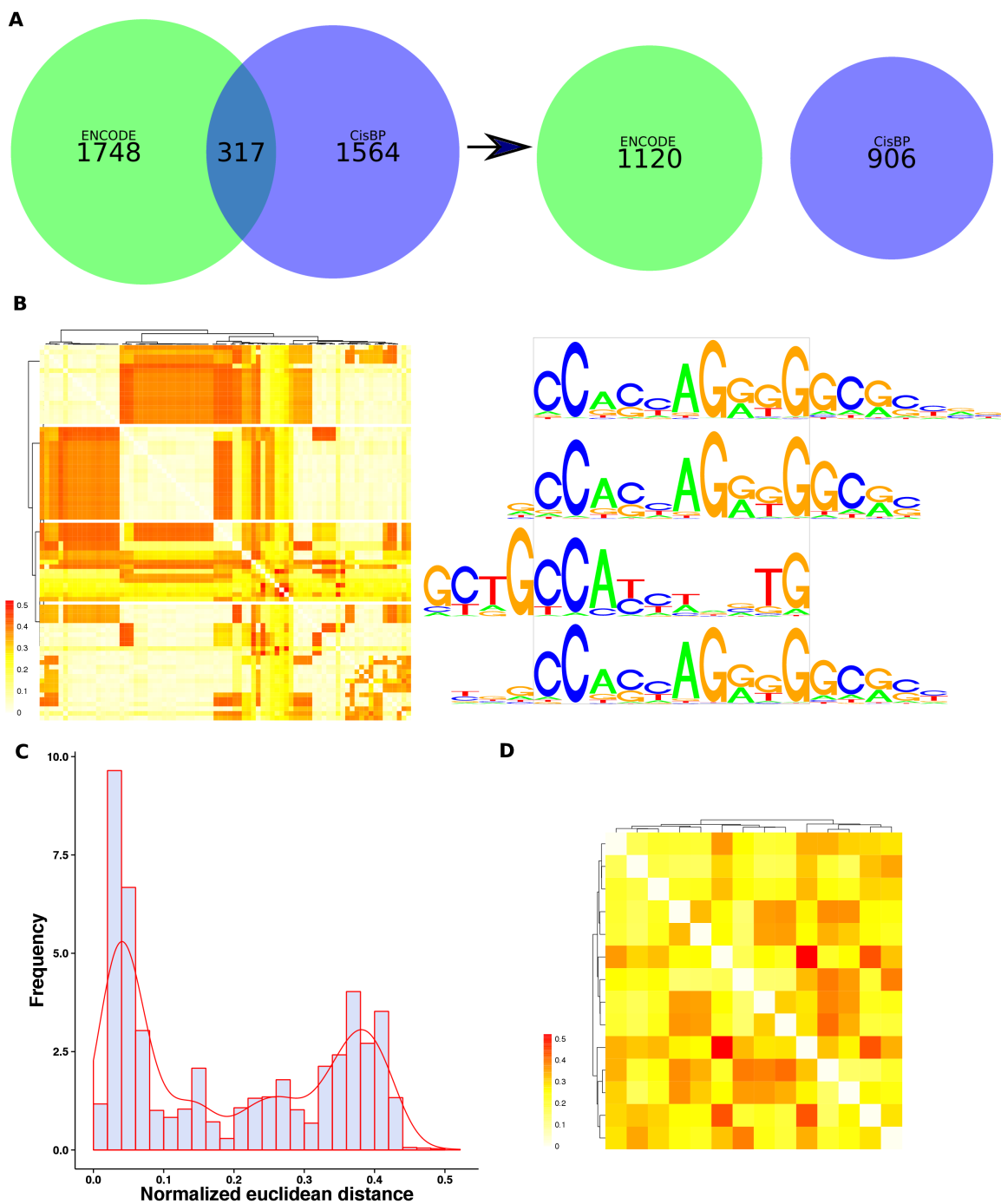
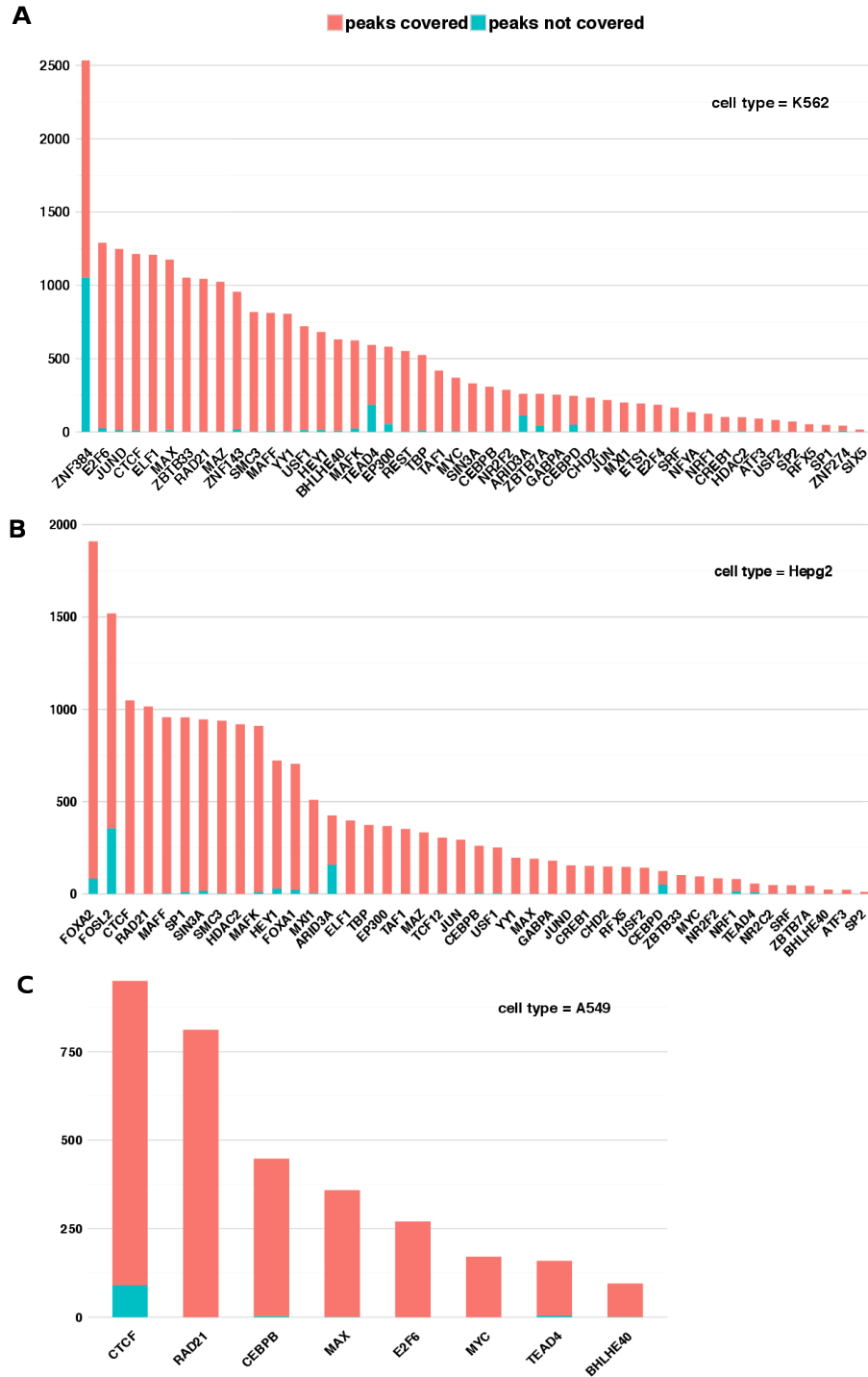


Figure S1: Filtering for highly redundant motifs representing the same TF. **(A)** Venn diagram showing size and overlap of the two motif collection, ENCODE (green) and CisBP (blue) motif set before and after redundancy filtering. ENCODE motif collection harbors motifs predominantly derived from ENCODE ChIP-Seq experiments, whereas CisBP motif collection includes motifs derived from large-scale PBM and SELEX studies. The overlap between the two motif collections are motifs drawn from databases such as TRANSFAC and JASPAR. **(B)** Similarity heatmap showing 79 motifs representing the TF, CTCF. White: highly similar motifs; Red: dissimilar motifs. Motifs are clustered based on normalized Euclidean distances between vectorized motifs and separated by row into four major groups. Representative motif logos for each of the four groups are drawn on the right. Cluster 1, 2, 4 represent motifs containing the canonical binding



site, with additional degenerate sites on the right, left and both side of the core motif respectively. Cluster 3 contains motifs that are more diverged from the canonical motif reflecting the multivalency in CTCF binding. (C) Histogram showing the distribution of between motif distances for CTCF. A conservative cutoff was selected at 0.14. (D) Similarity heatmap after redundancy filtering leaving 14 motifs that are more dissimilar from each other.



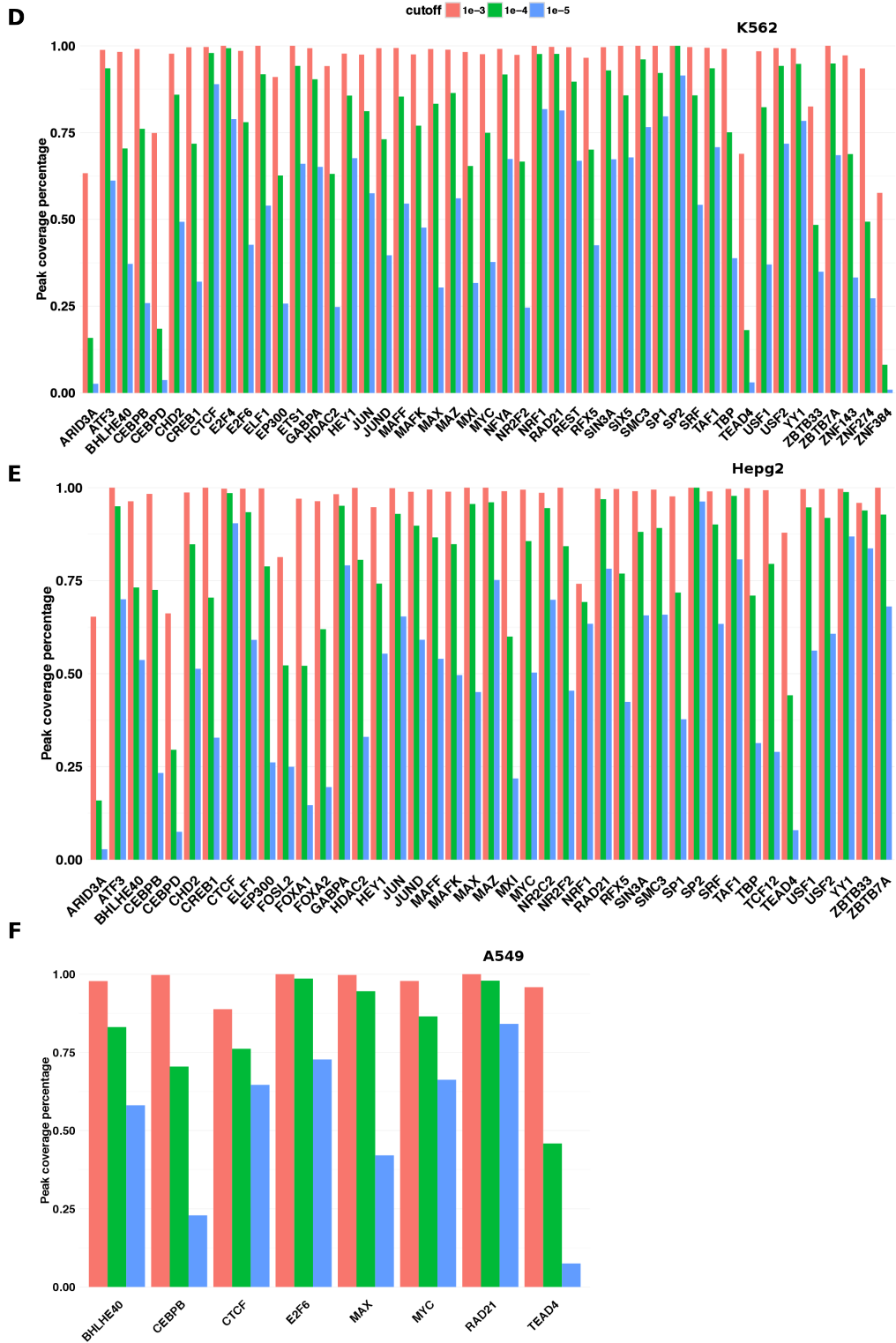


Figure S2: ChIP-Seq peak coverage by motifs. (A-C) Bar plot showing the number of ChIP-Seq peaks covered (red) and not covered (blue) by our analyses on testing chromosome 15 in cell type Hepg2, K562 and A549. (D-F) Grouped barplot showing the percentage of ChIP-Seq peak coverage using different FIMO motif scan cutoffs in the three cell types on chromosome 3.

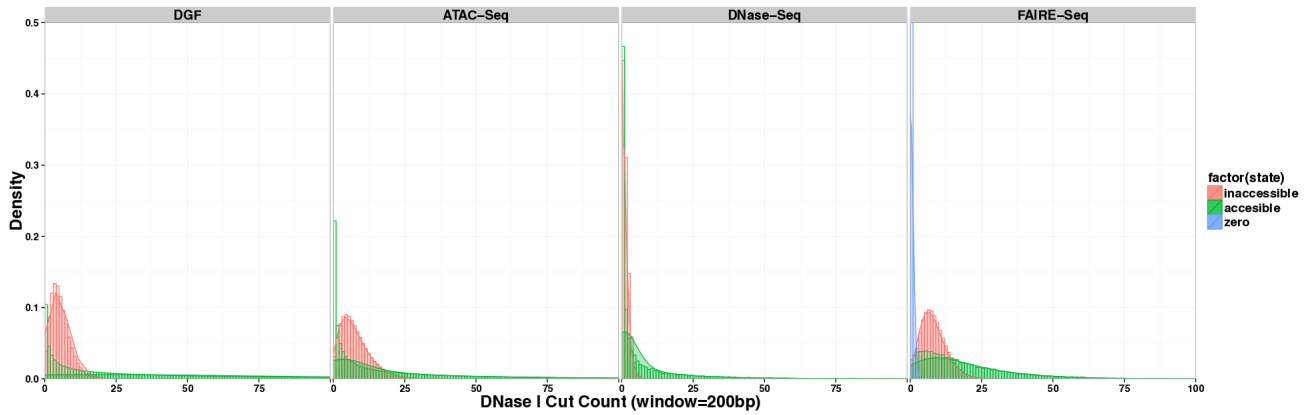


Figure S3: Cut count density distribution across experiment types. Data are simulated using zero-inflated negative binomial mixture model parameters derived from different experiment types. Zero component is insignificant in all but FAIRE-Seq. Red: inaccessible component; Green: accessible component; Blue; zero component.

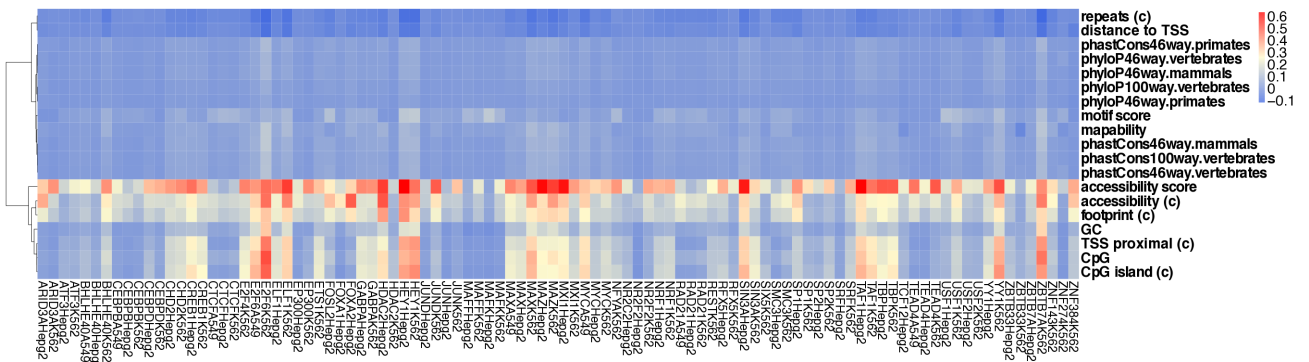
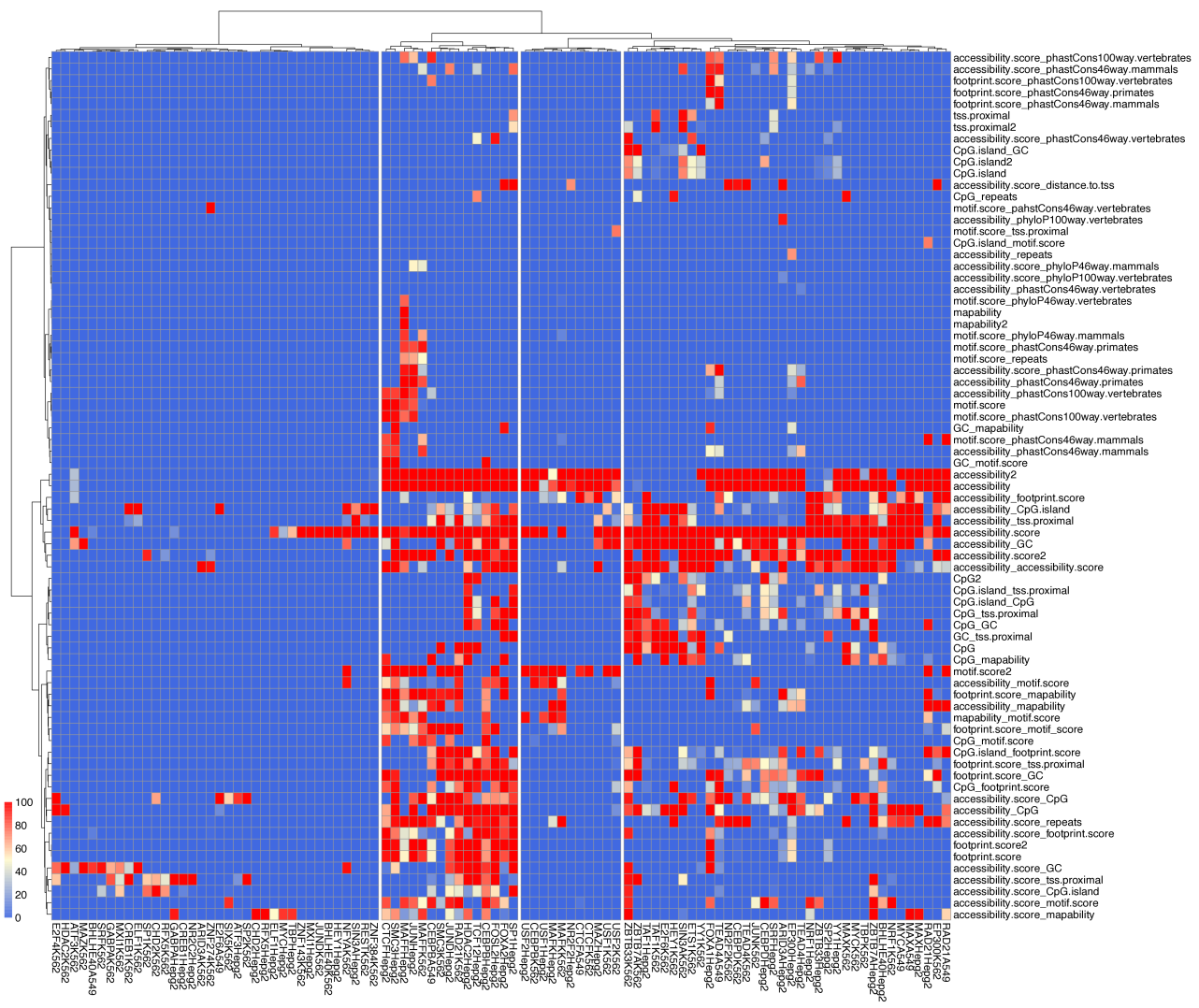


Figure S4: Feature correlations with ChIP-Seq signals. Heatmap showing the PCC between 19 genomic features and ChIP-Seq signals. Red indicates high positive correlation between a genomic feature and its associated ChIP-Seq signal; Blue indicates low positive or negative correlation. Genomic features are clustered based on their correlation value across 98 TF-cell type samples. Categorical features are marked with (c).

A



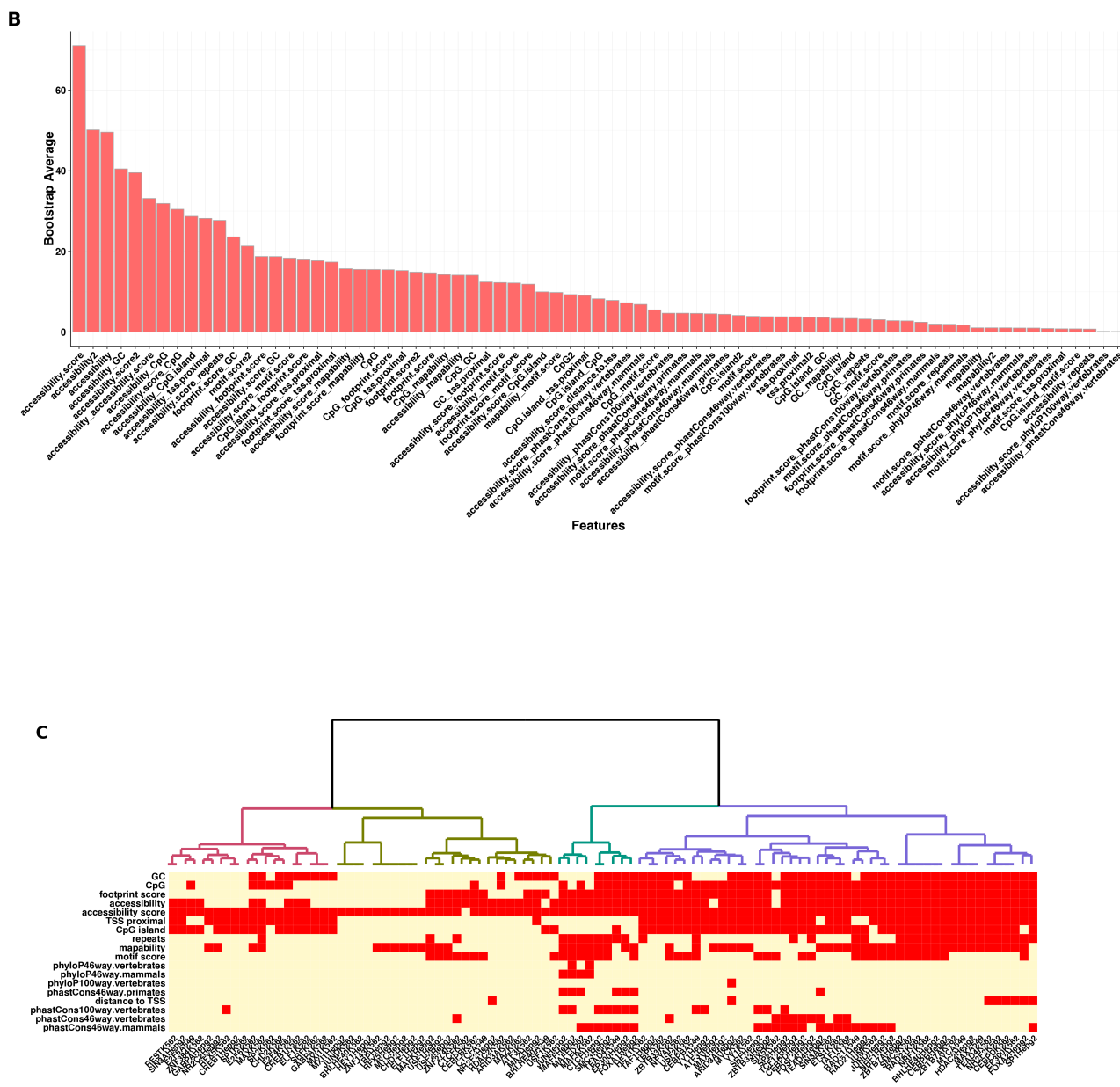


Figure S5: Training sparse logistic regression classifiers by bootstrap aggregation. (A) Heatmap showing the number of bootstraps each single or interaction feature was selected during 100 bootstrap aggregation. Red indicates more stable features that are likely to be selected by the majority of bootstraps, whereas blue indicates features less likely to be selected. Features are clustered based on the majority of bootstraps that select them. (B) A barplot summary of the selected features ranked based on their bootstrap stability across samples. (C) Heatmap showing feature usage by different TFCT conditions. TFCT conditions are clustered into four major clusters based on feature usage preference. Group 1 (red) prefers models with local accessibility, TSS proximity and GC/CpG sequence features. Group 2 (khaki) predominantly selected a combination of accessibility and motif scores. Group 3 (green) and Group 4 (purple) both make use of a fuller range of features, with Group 4, additionally, preferring the binary predictor TSS proximity.

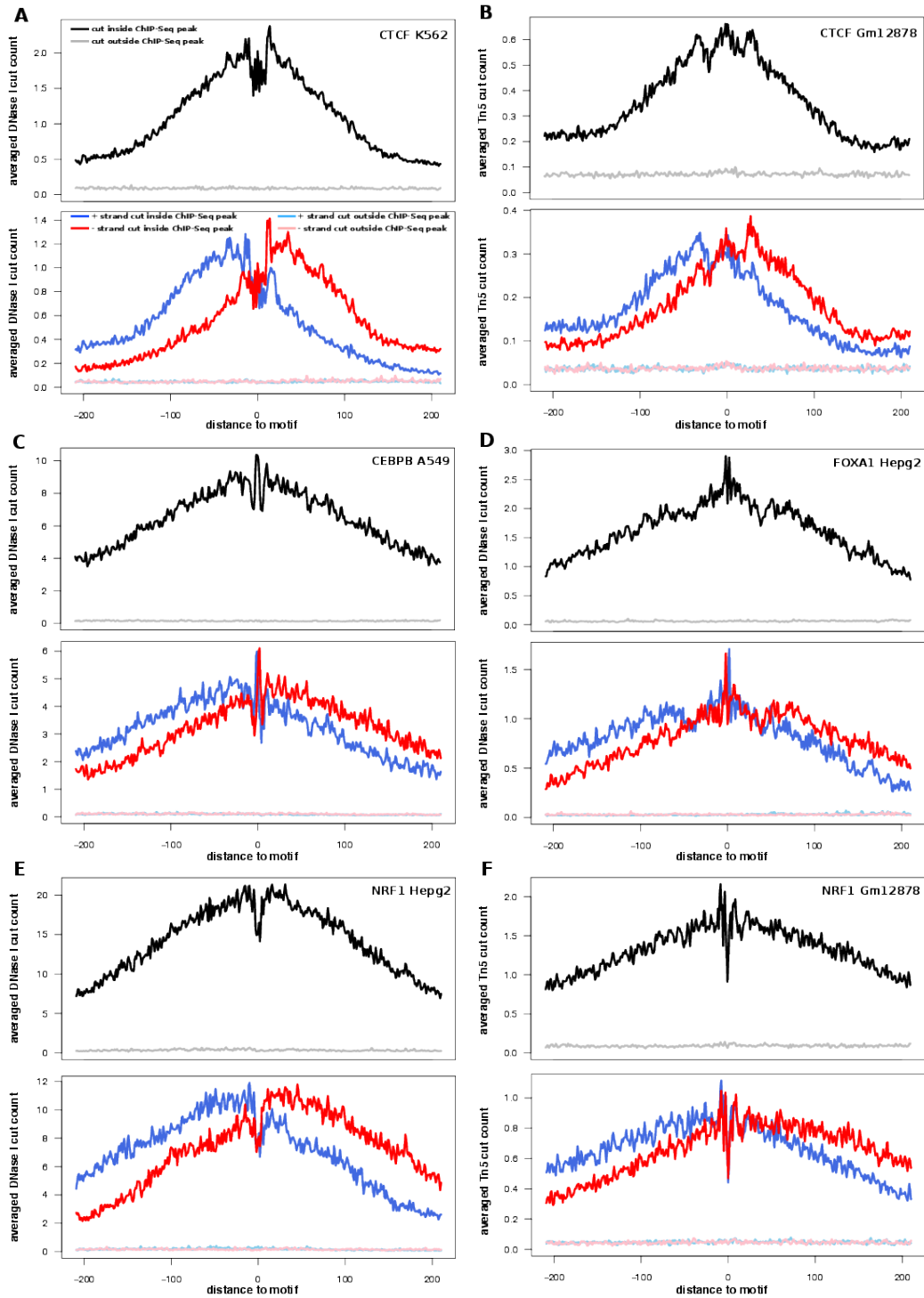


Figure S6: Footprint profiles vary across factors, cell types and experimental protocols. Black: cut count pattern around motifs overlapping ChIP-Seq peaks; grey: cut count pattern around motifs not overlapping ChIP-Seq peaks; blue/red: cut count pattern on the positive/negative strand around motifs overlapping ChIP-Seq peaks; light blue/pink: cut count pattern on the positive/negative strand around motifs not overlapping ChIP-Seq peaks.

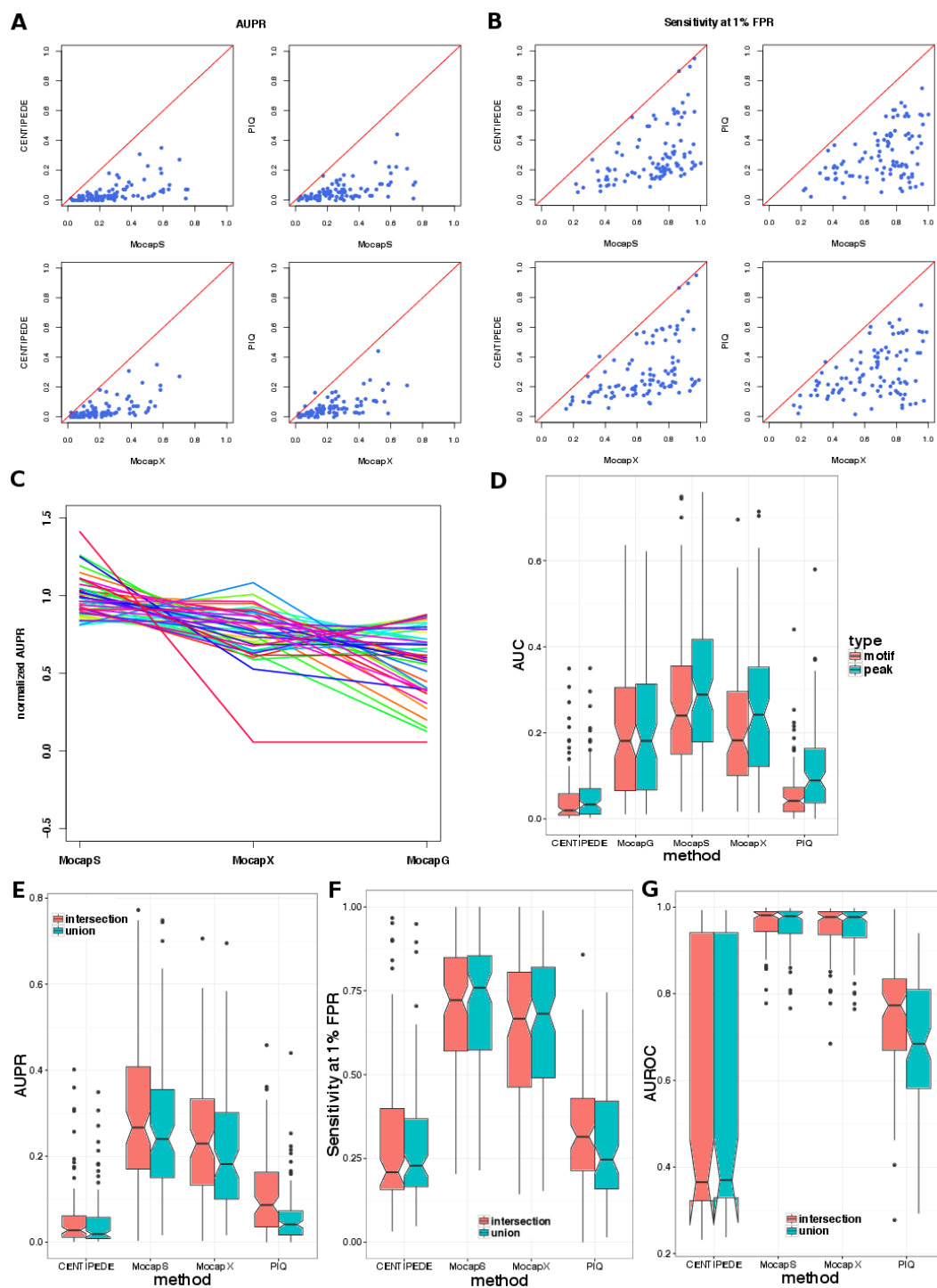


Figure S7: Performance comparison between methods. (A-B) Scatterplots showing pairwise comparisons between MocapS/MocapX and CENTIPEDE/PIQ in AUPRC and Sensitivity at 1% FPR. (C) Line plot contrasting differences in performance between MocapS, MocapX and MocapG for each of the 35 TF-cell type conditions with mapping to another TF-cell type trained sparse logistic regression model using MocapX. (D) Grouped boxplots showing similar results comparing the areas under precision over motif site recall curves (red) or precision over ChIP-Seq peak recall curves (green) across methods. (E-G) Grouped boxplots showing similar performance comparison results using the intersection (red) or union (green) of all method-predicted motif sets as the scope of comparison for AUPR, Sensitivity at 1% FPR and AUROC respectively.



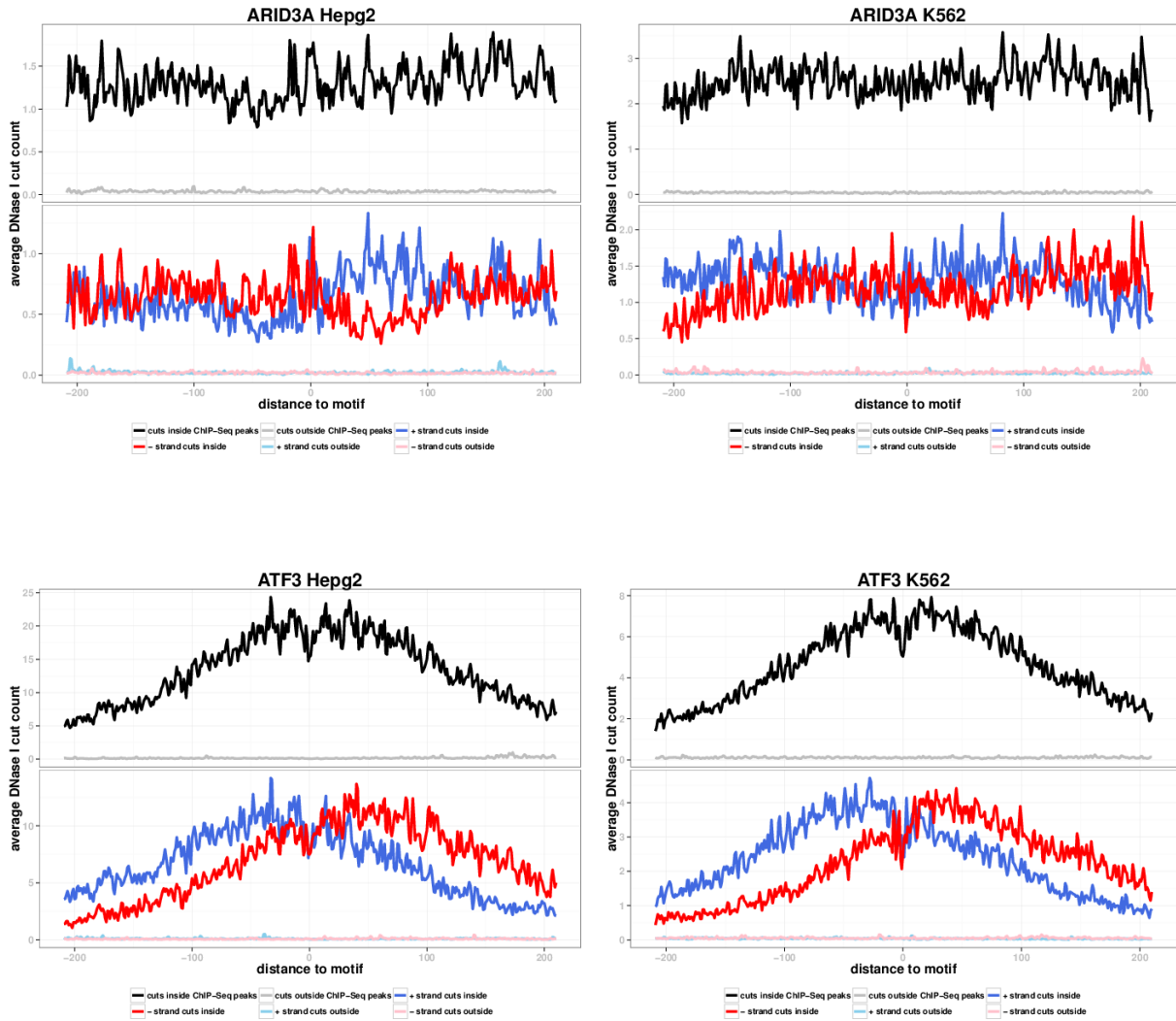


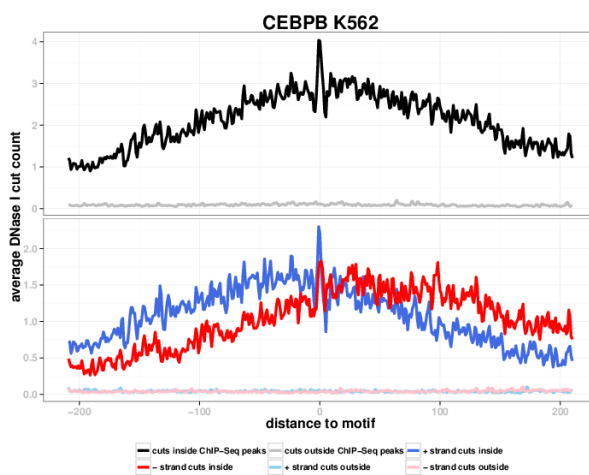
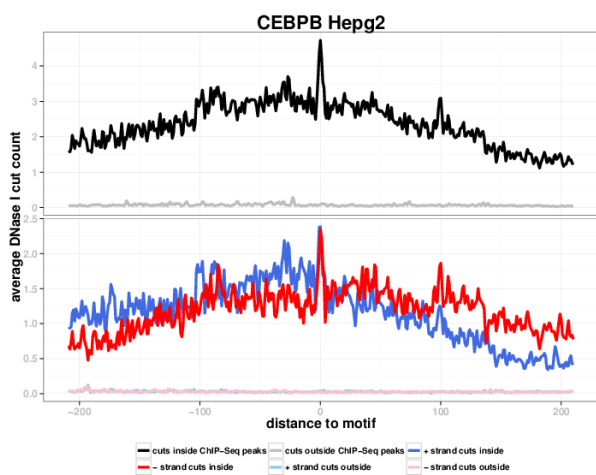
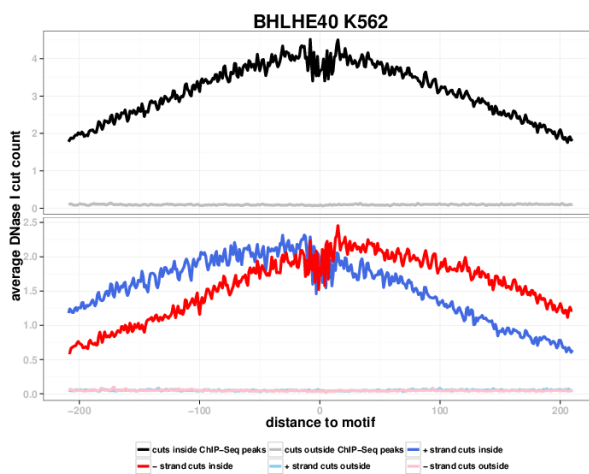
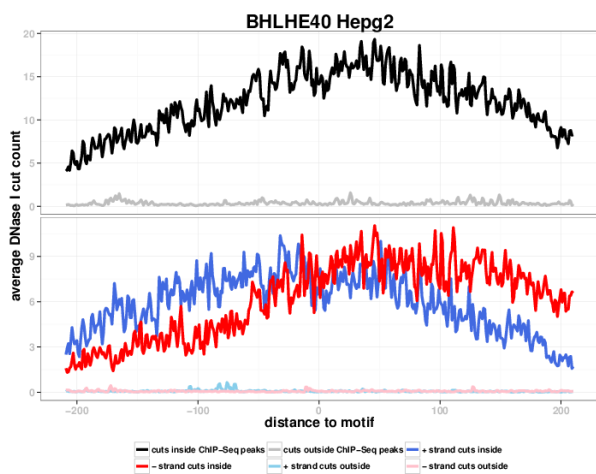
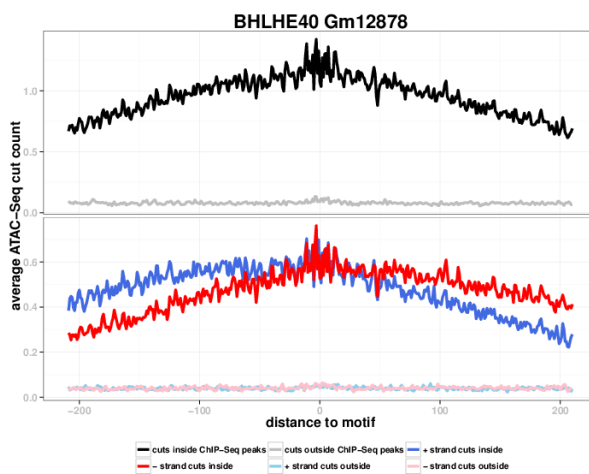
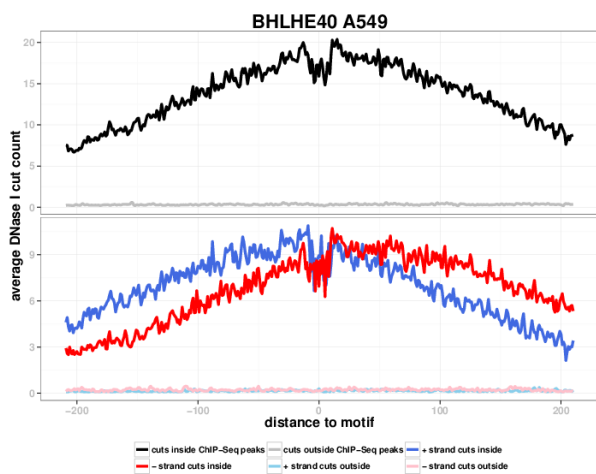
Table S2: Performance comparison for ATAC-Seq samples

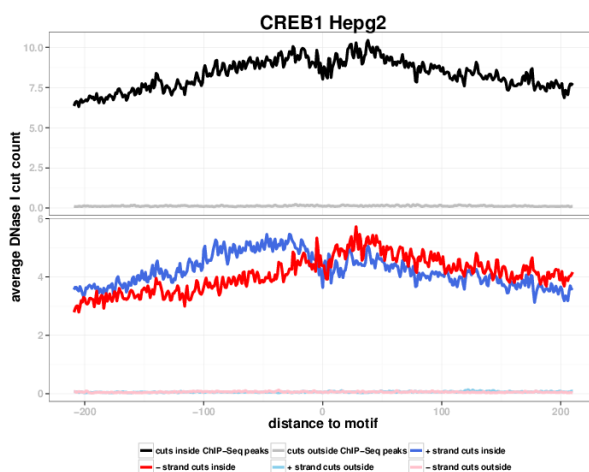
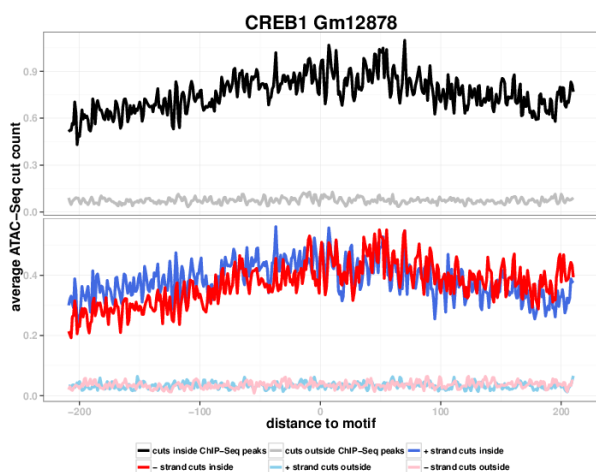
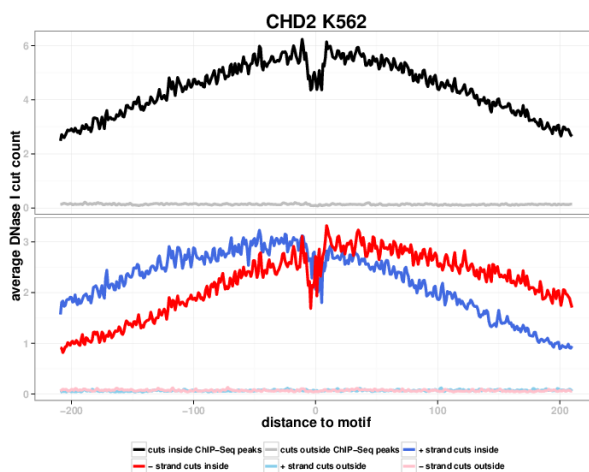
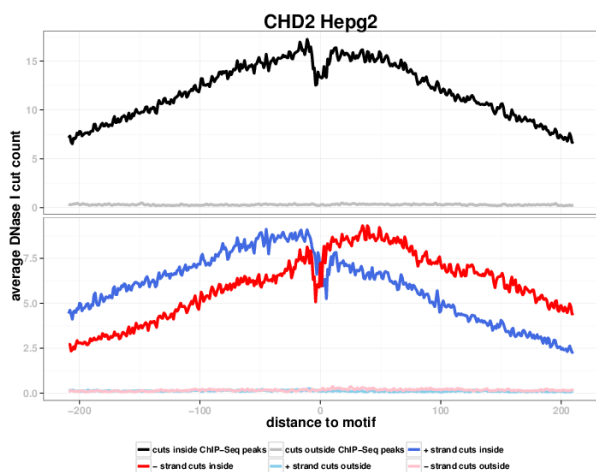
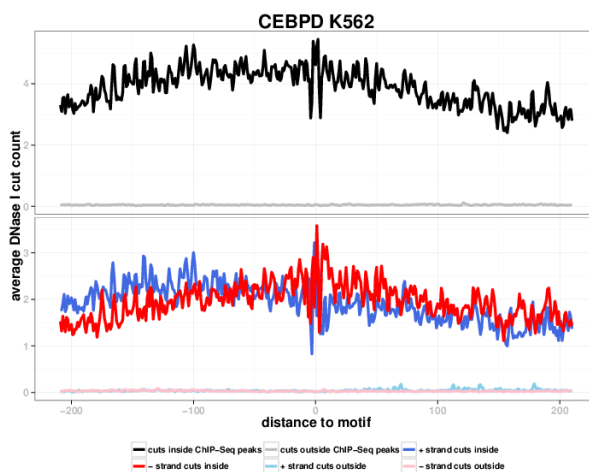
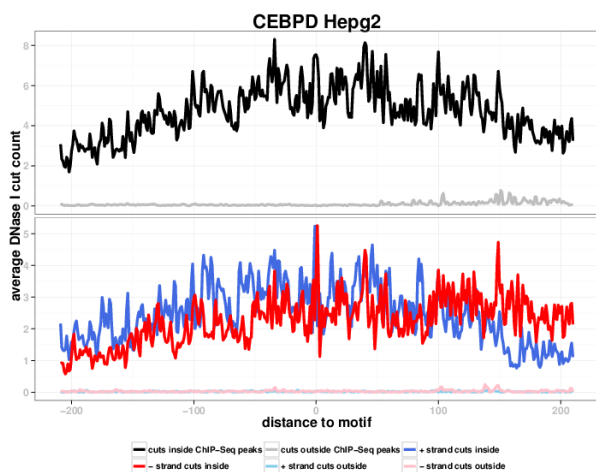
TF*	Sample mapping	MocapG			MocapX			CENTIPEDE		
		AUPR	AUROC	Sensitivity	AUPR	AUROC	Sensitivity	AUPR	AUROC	Sensitivity
BHLHE40	default	0.24	0.95	0.59	0.24	0.95	0.59	0.05	0.39	0.26
CREB1	default	0.08	0.91	0.54	0.08	0.91	0.54	0.03	0.4	0.26
CTCF	CTCFK562	0.05	0.81	0.31	0.12	0.80	0.32	0.01	0.44	0.12
E2F4	default	0.14	0.99	0.87	0.14	0.99	0.87	0.02	0.34	0.26
ELF1	NR2C2Hepg2	0.44	0.98	0.79	0.46	0.97	0.8	0.06	0.34	0.25
EP300	default	0.17	0.95	0.8	0.17	0.95	0.8	0.03	0.35	0.25
ETS1	default	0.15	0.98	0.82	0.15	0.98	0.82	0.02	0.37	0.29
GABPA	default	0.13	0.98	0.84	0.13	0.98	0.84	0.02	0.39	0.3
MAFK	MAFKK562	0.002	0.56	0.07	0.02	0.74	0.17	0.001	0.46	0.04
MAZ	MAZK562	0.52	0.97	0.76	0.51	0.97	0.77	0.1	0.37	0.26
MXI1	default	0.41	0.96	0.79	0.41	0.96	0.79	0.04	0.26	0.15
NFYA	default	0.32	0.98	0.9	0.32	0.98	0.9	0.02	0.39	0.3
NRF1	NRF1K562	0.25	0.99	0.8	0.31	0.99	0.84	0.03	0.35	0.25
RAD21	RAD21Hepg2	0.05	0.83	0.32	0.12	0.86	0.4	0.01	0.44	0.13
SIN3A	SIN3AK562	0.28	0.98	0.8	0.28	0.99	0.85	0.06	0.36	0.28
SIX5	default	0.06	0.98	0.93	0.06	0.98	0.93	0.004	0.34	0.25
SP1	default	0.36	0.96	0.81	0.36	0.96	0.81	0.04	0.37	0.28
SRF	SRFHepg2	0.19	0.95	0.76	0.17	0.93	0.72	0.03	0.4	0.3
TCF12	TCF12Hepg2	0.40	0.91	0.58	0.36	0.95	0.51	0.1	0.47	0.2
USF1	USF1Hepg2	0.07	0.88	0.5	0.18	0.92	0.59	0.01	0.37	0.16
YY1	YY1K562	0.3	0.94	0.73	0.28	0.95	0.71	0.05	0.38	0.25
ZBTB33	default	0.06	0.99	0.73	0.06	0.99	0.73	0.005	0.32	0.23
ZNF143	ZNF143K562	0.15	0.9	0.5	0.15	0.9	0.5	0.04	0.42	0.24
ZNF384	default	0.15	0.68	0.33	0.15	0.68	0.33	0.02	0.36	0.08

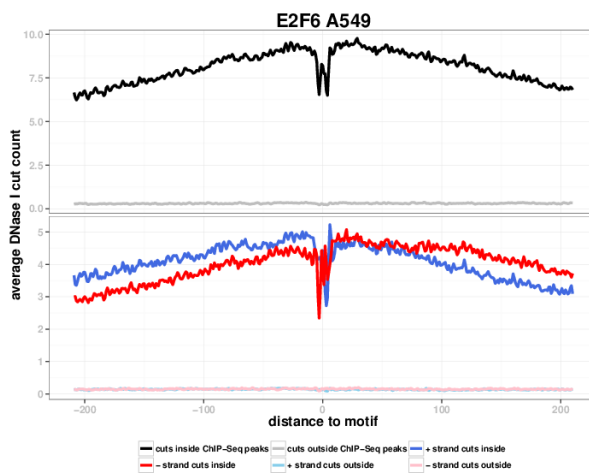
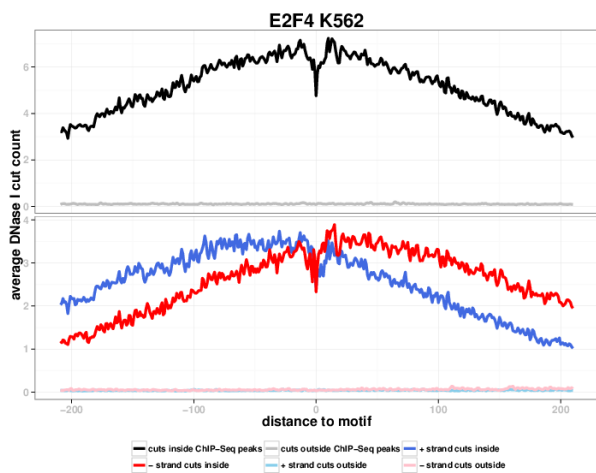
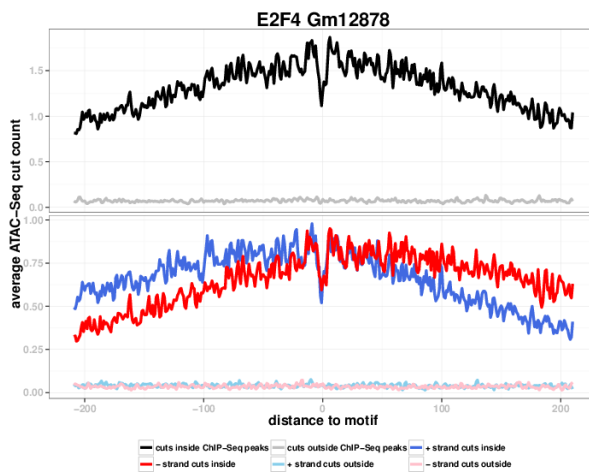
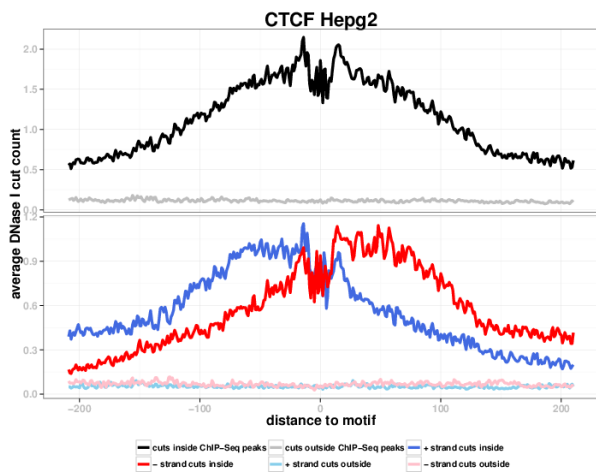
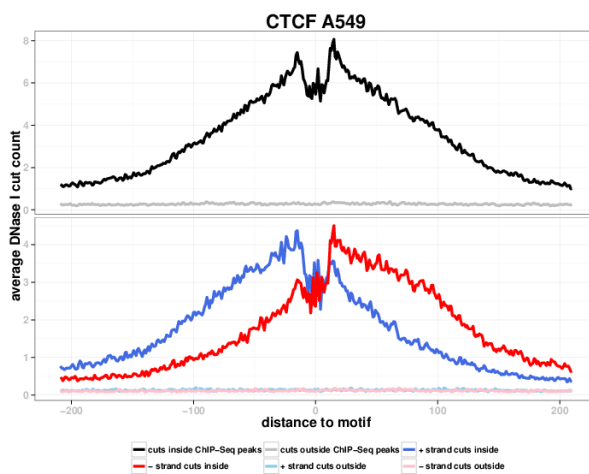
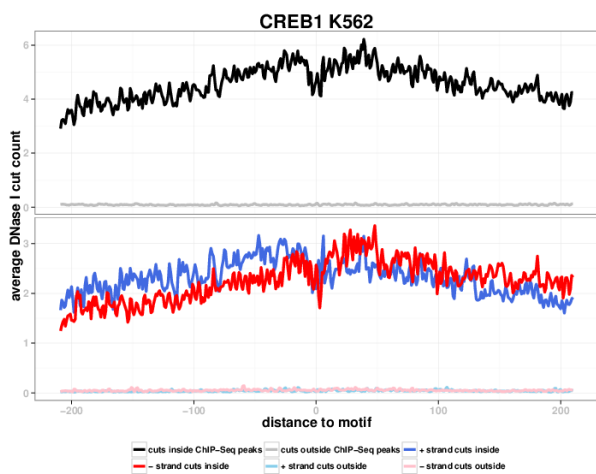
\*Cell type: Gm12878.

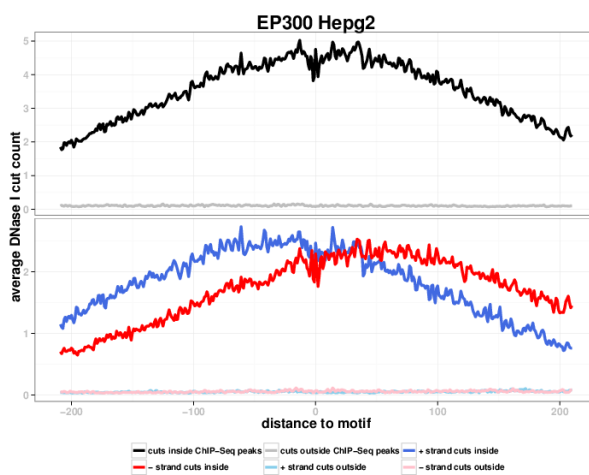
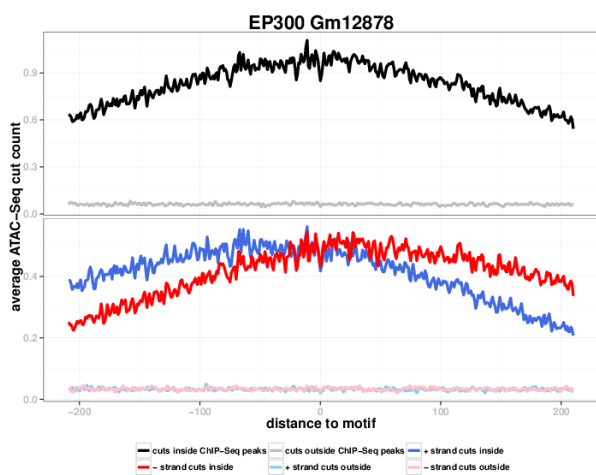
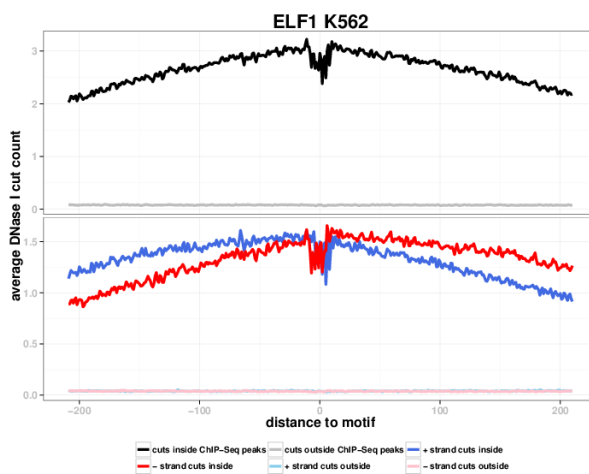
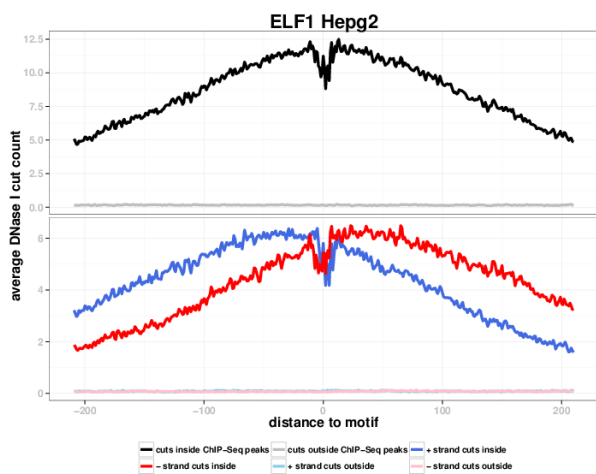
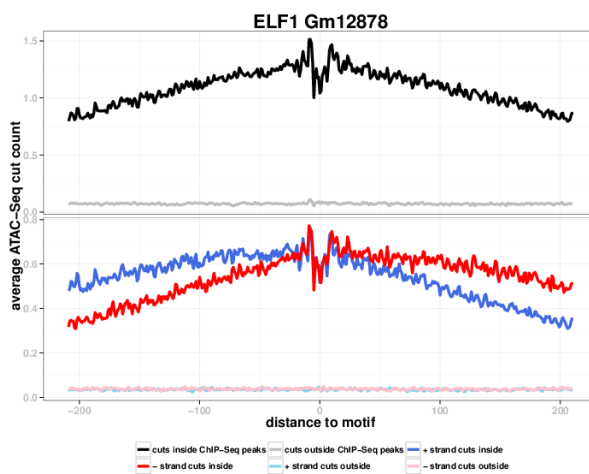
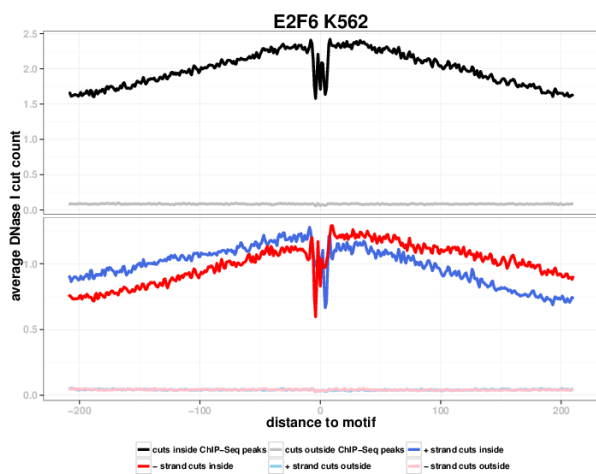
## Appendix A Footprint profile plots

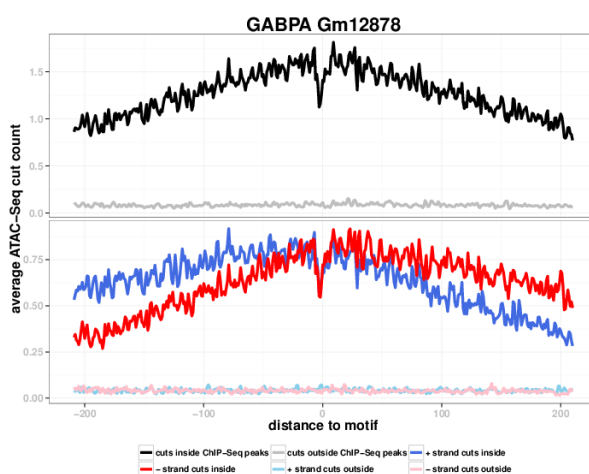
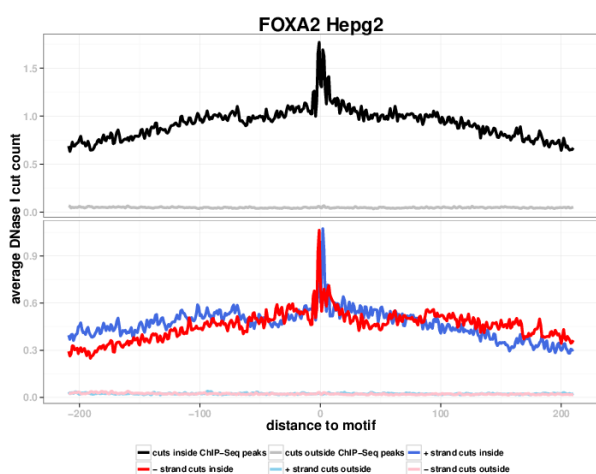
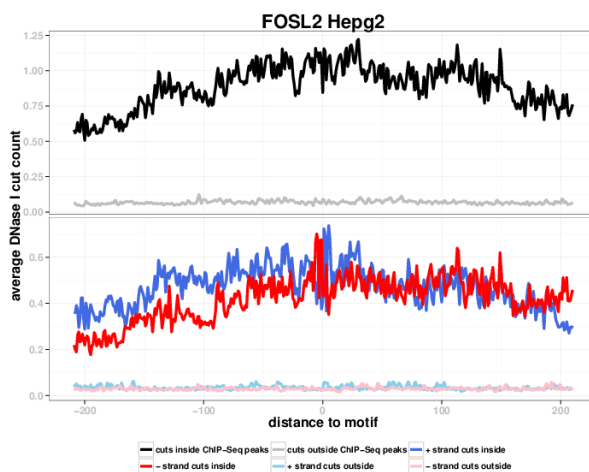
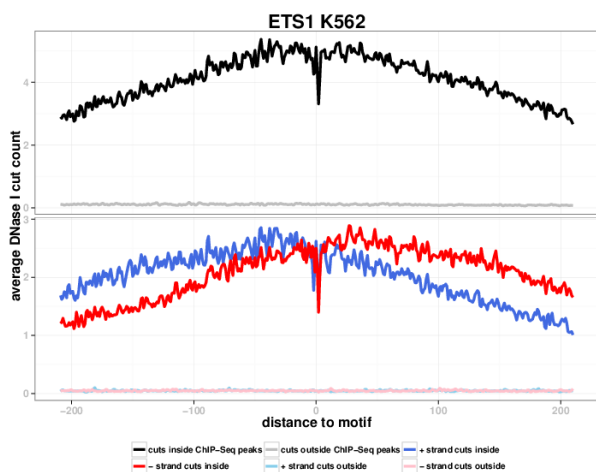
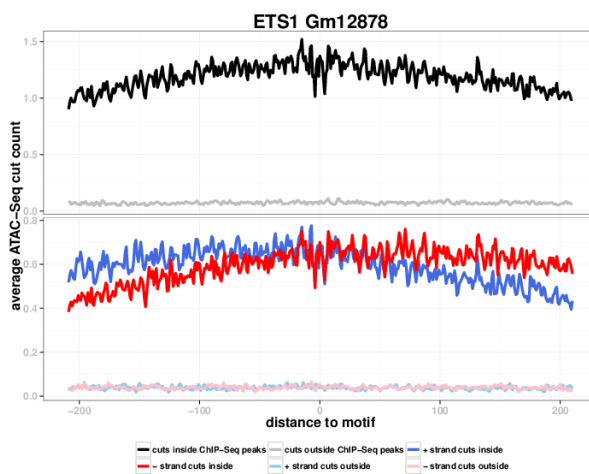
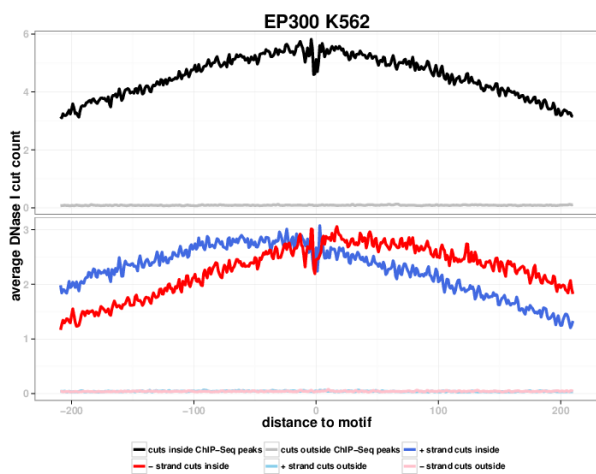


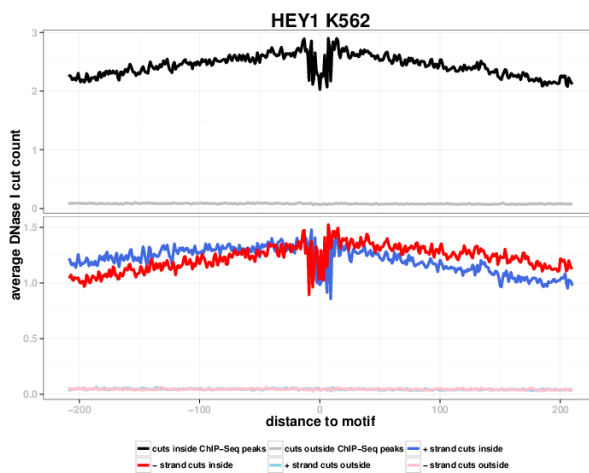
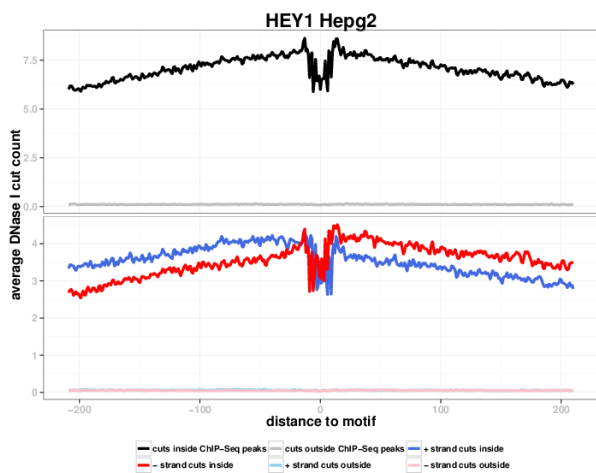
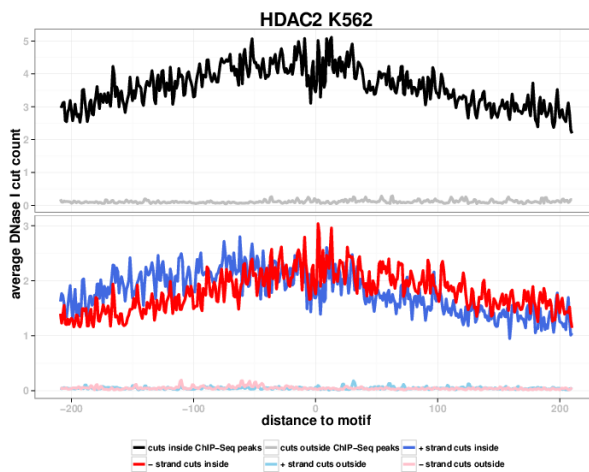
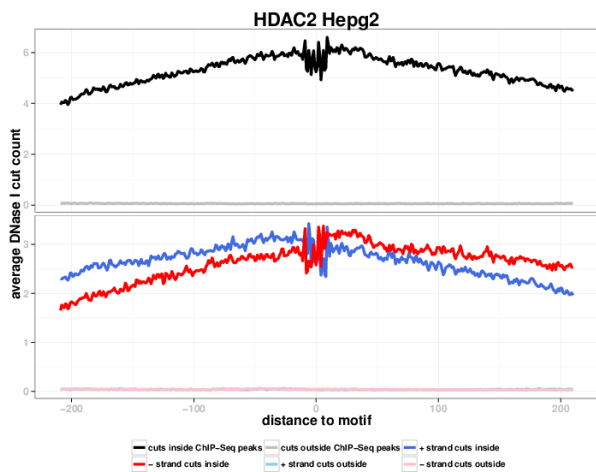
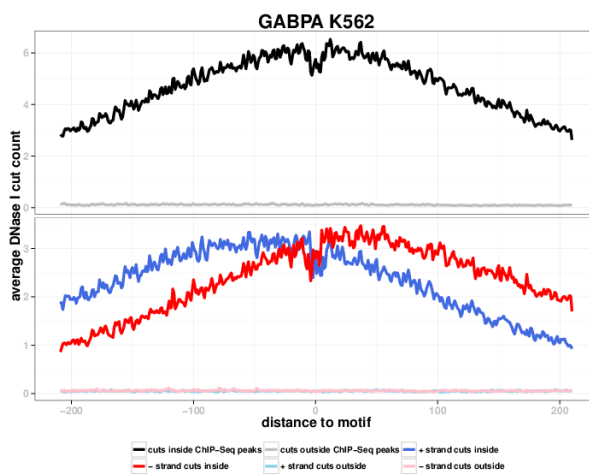
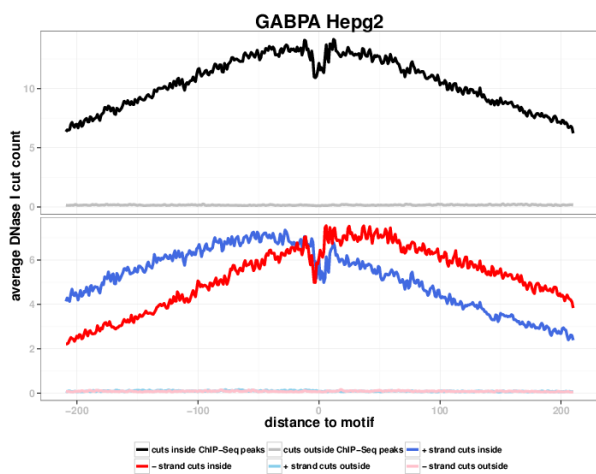




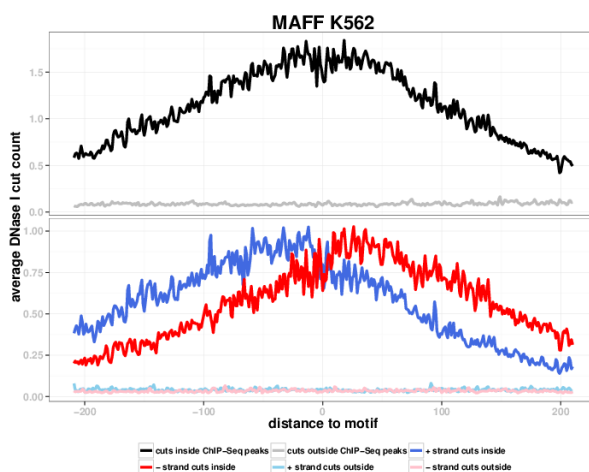
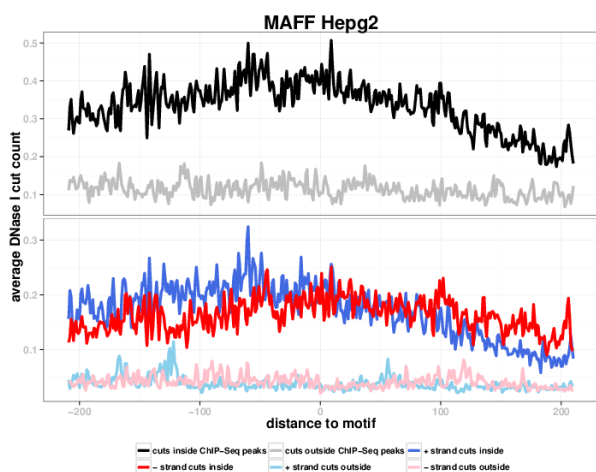
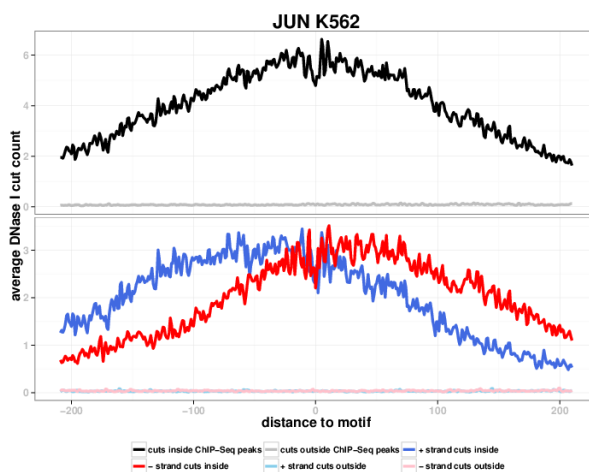
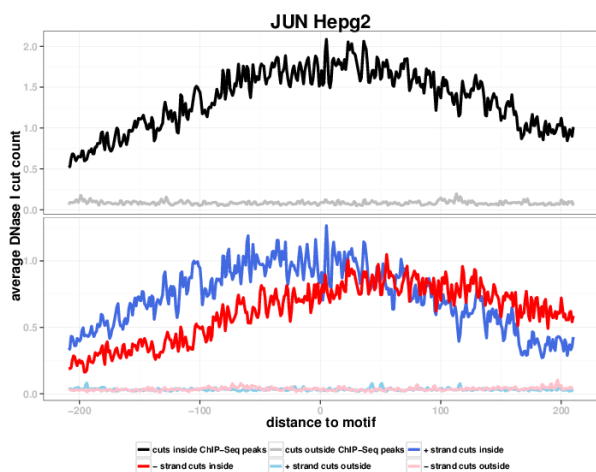
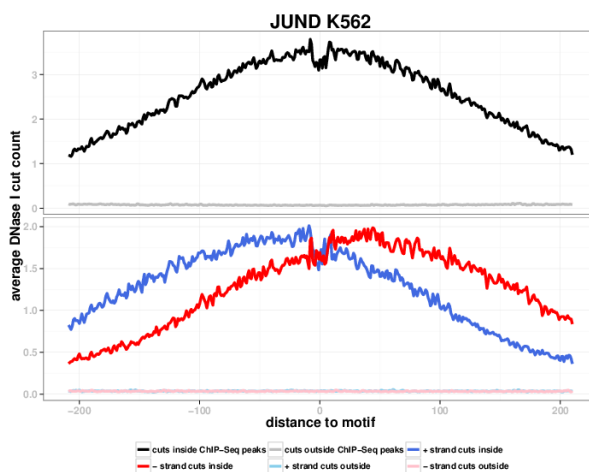
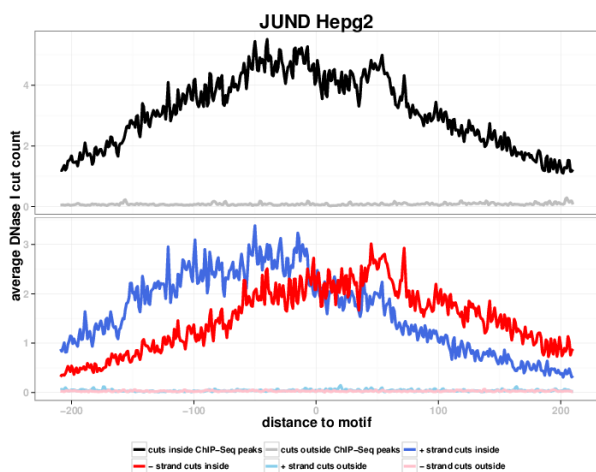


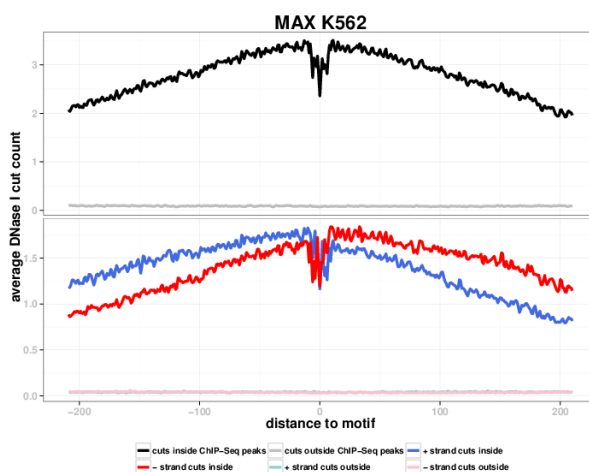
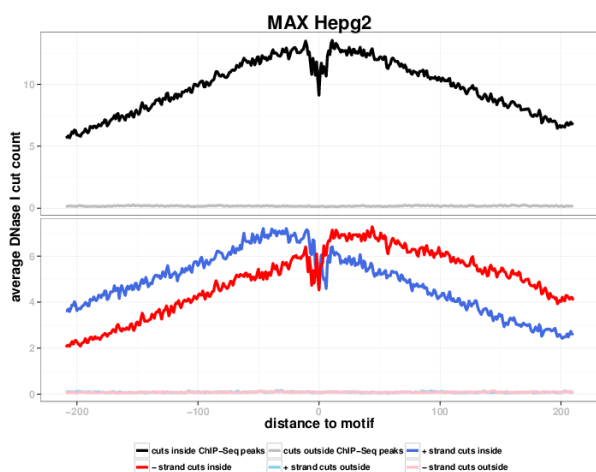
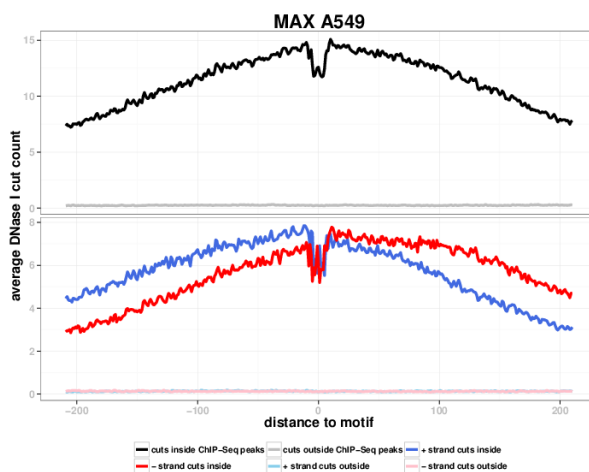
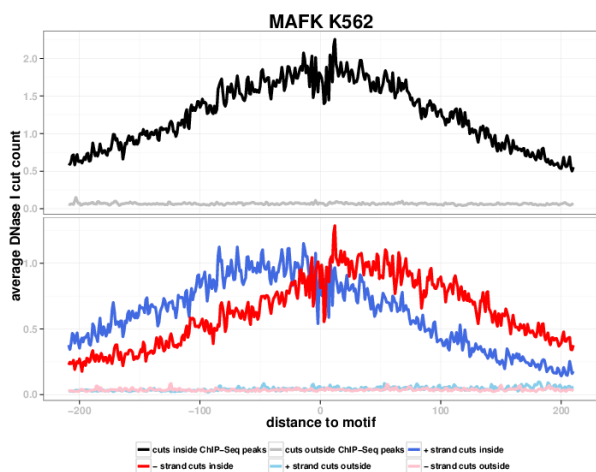
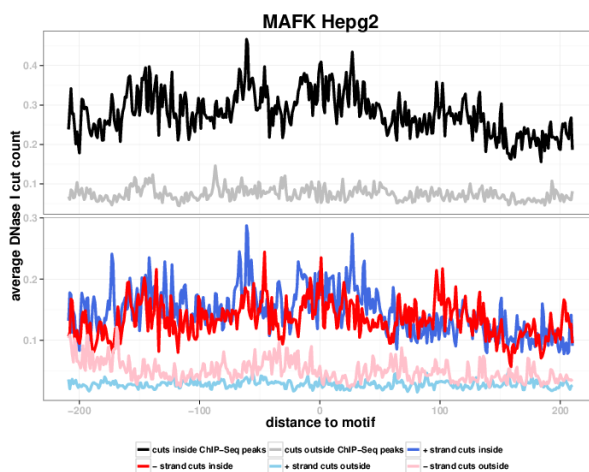
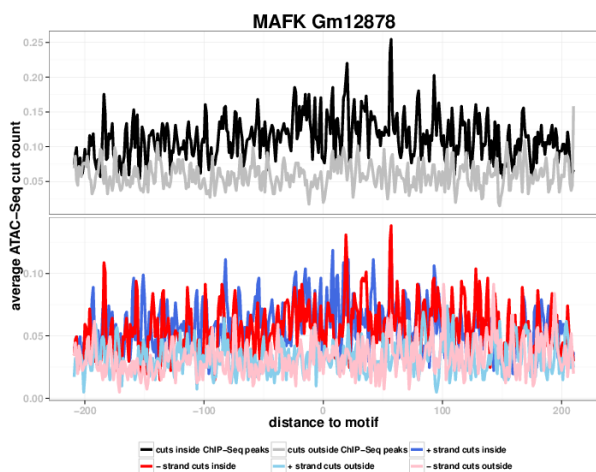


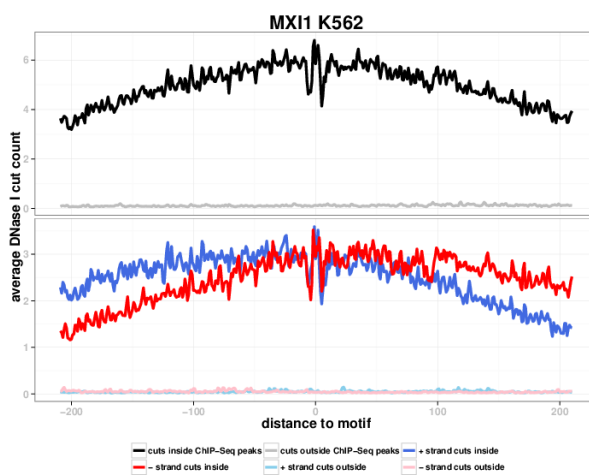
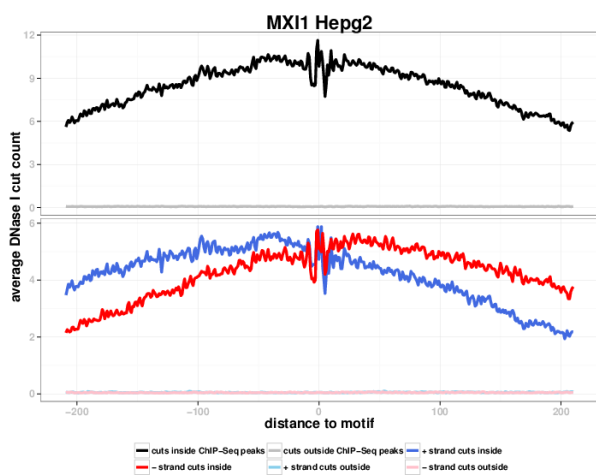
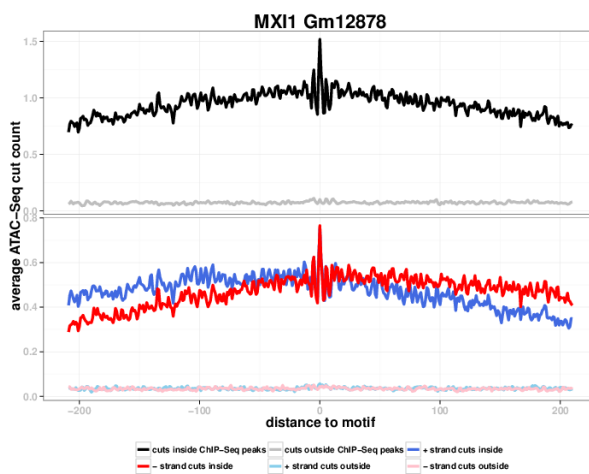
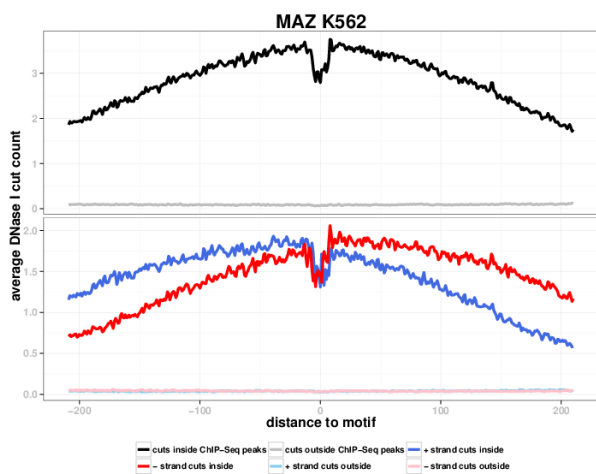
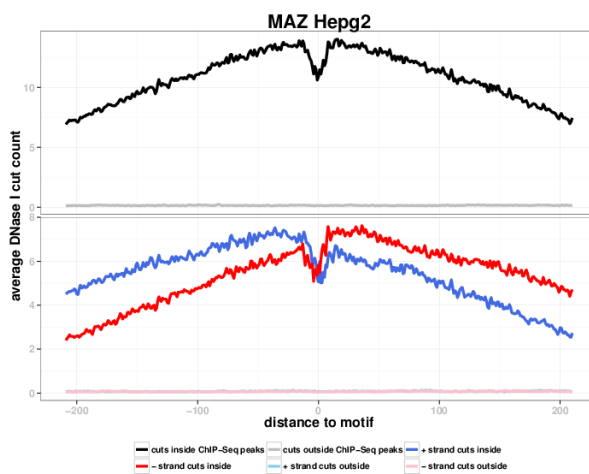
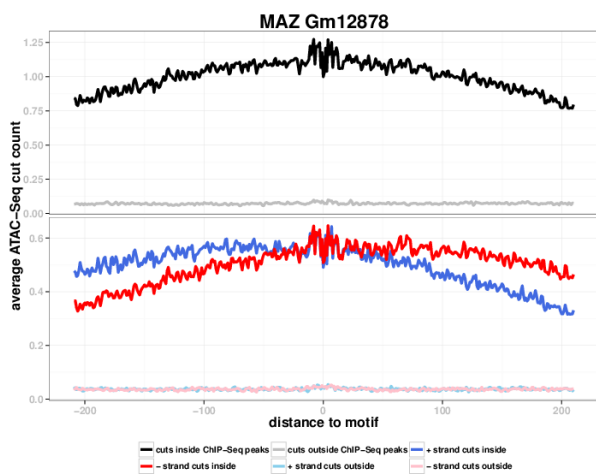


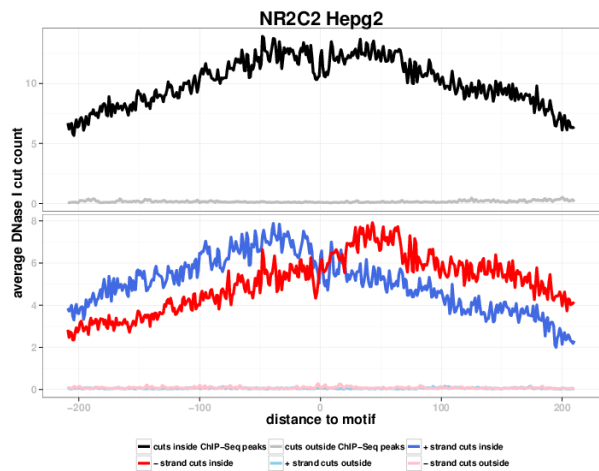
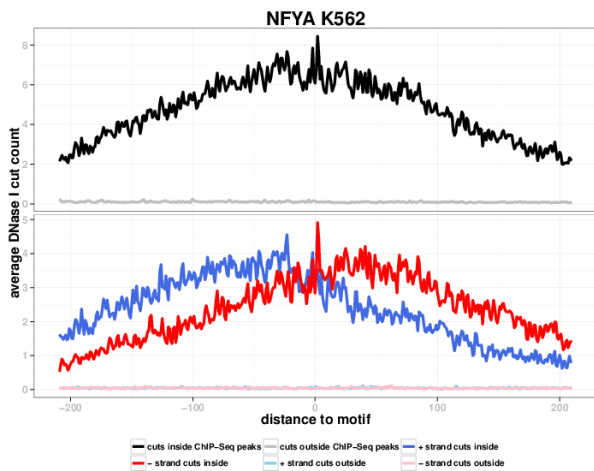
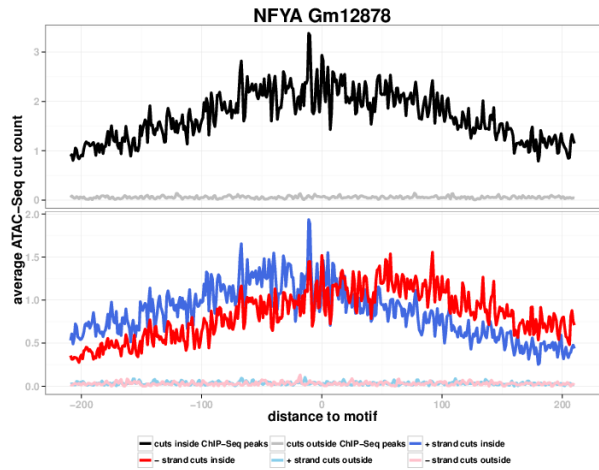
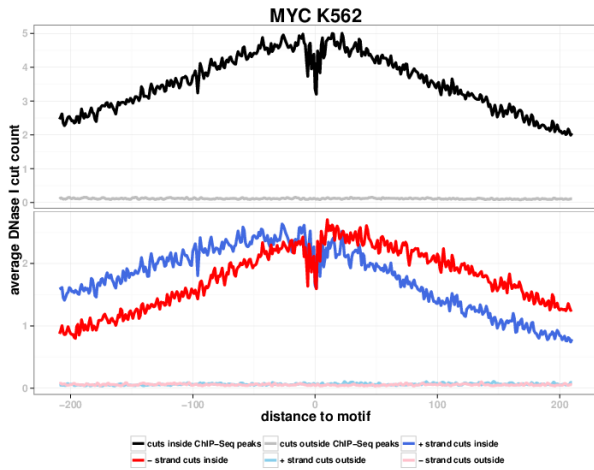
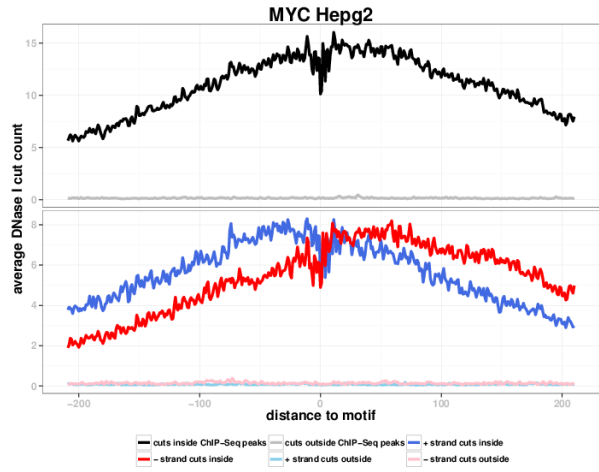
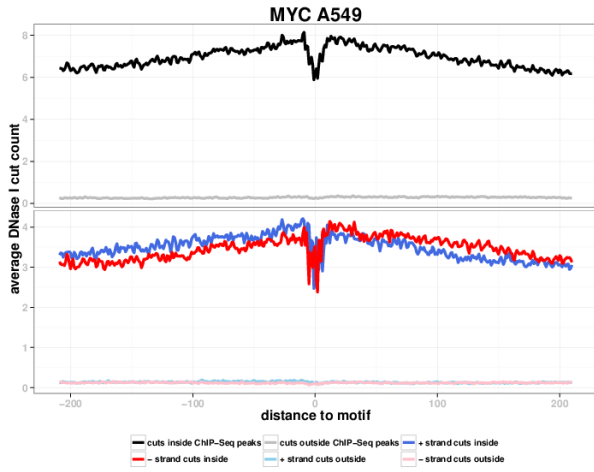


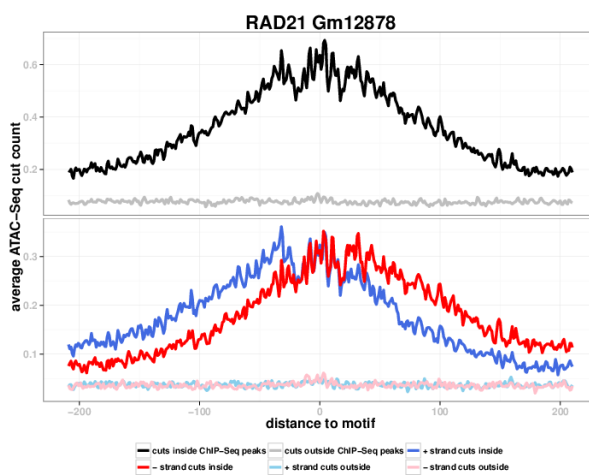
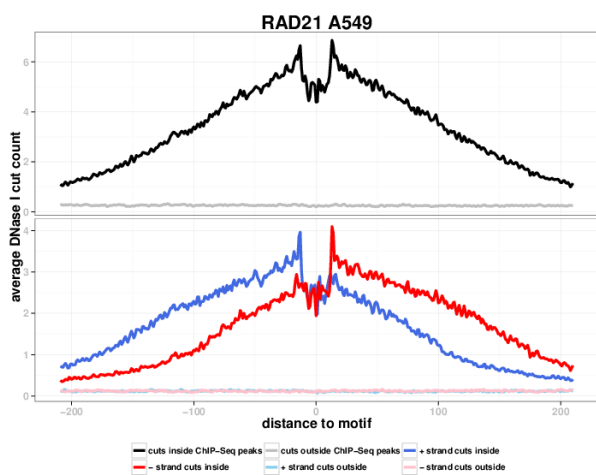
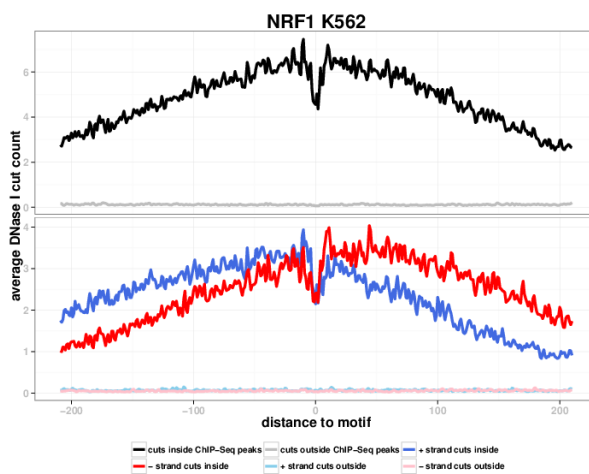
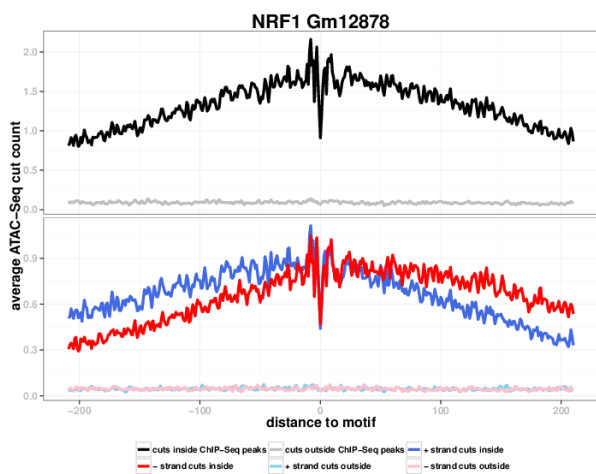
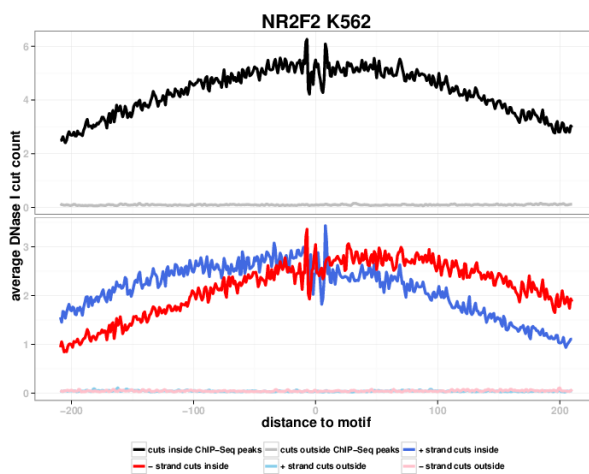
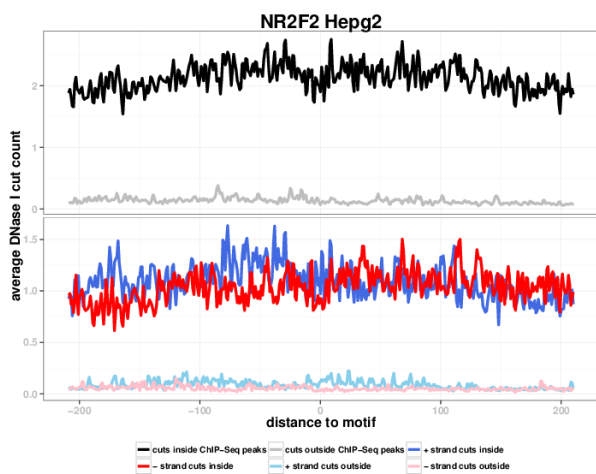


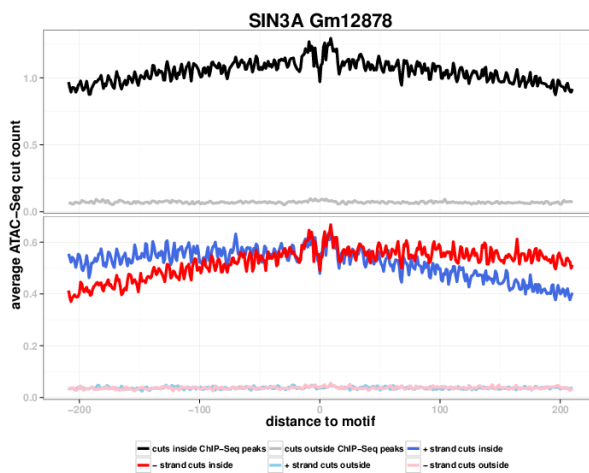
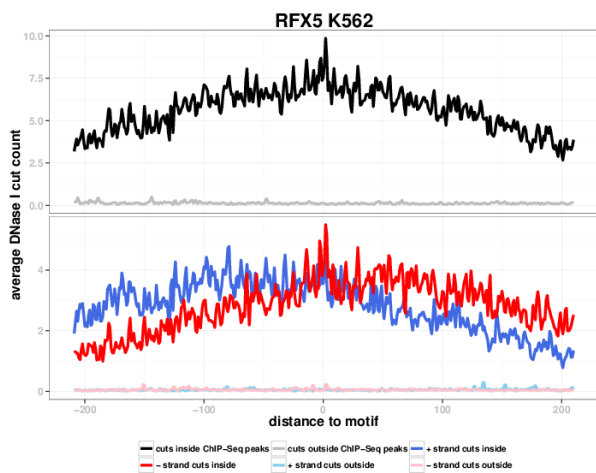
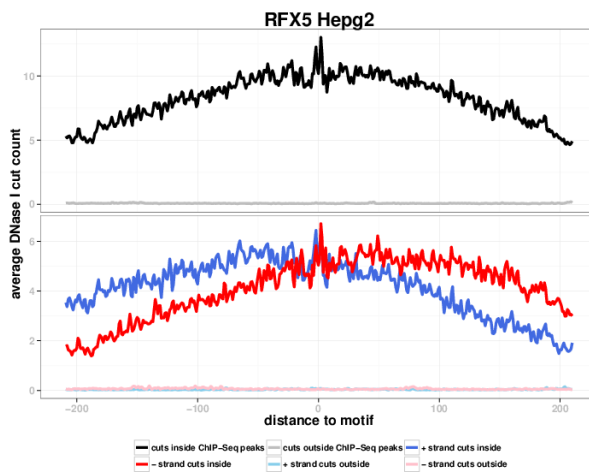
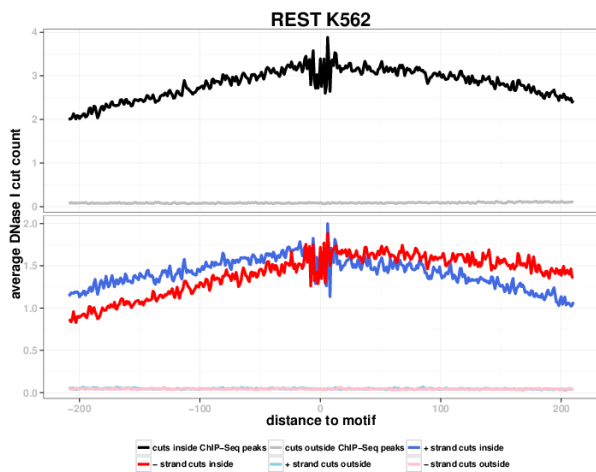
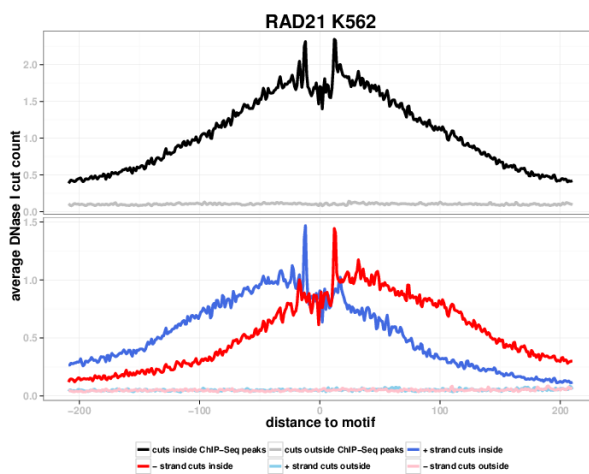
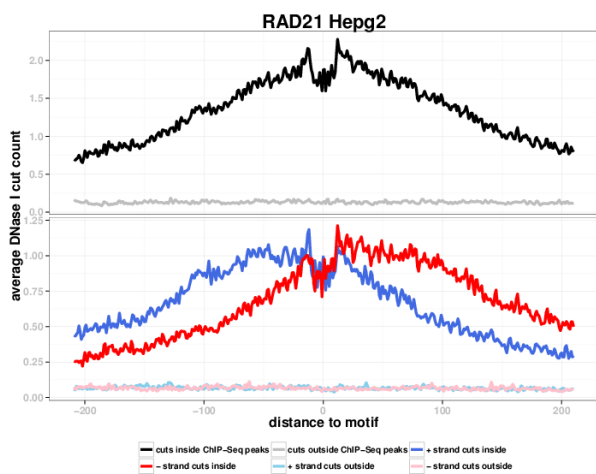


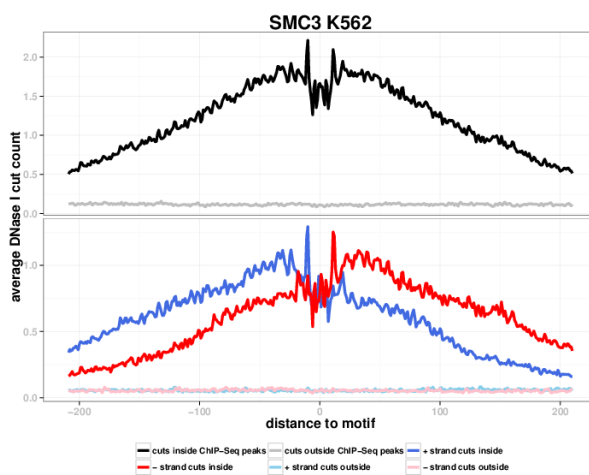
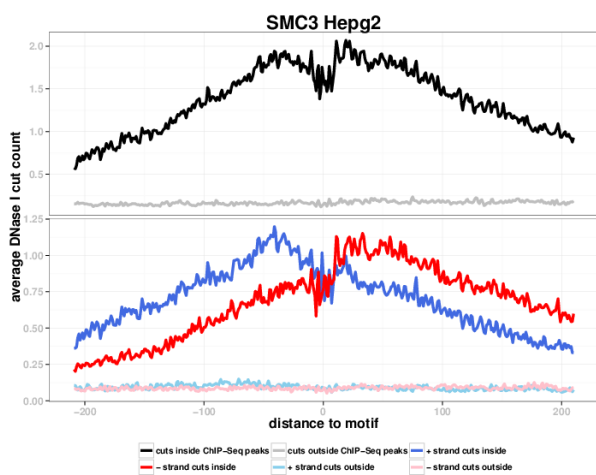
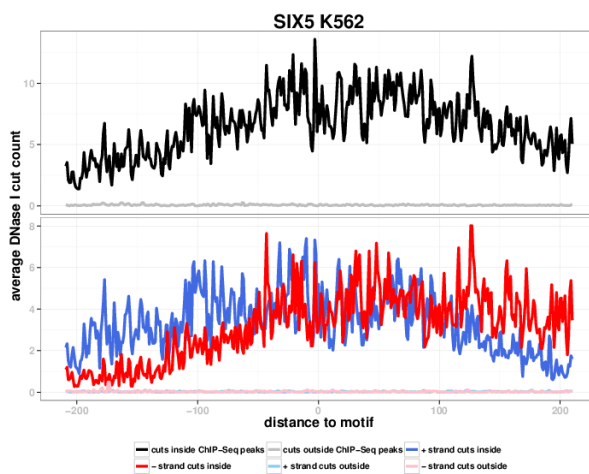
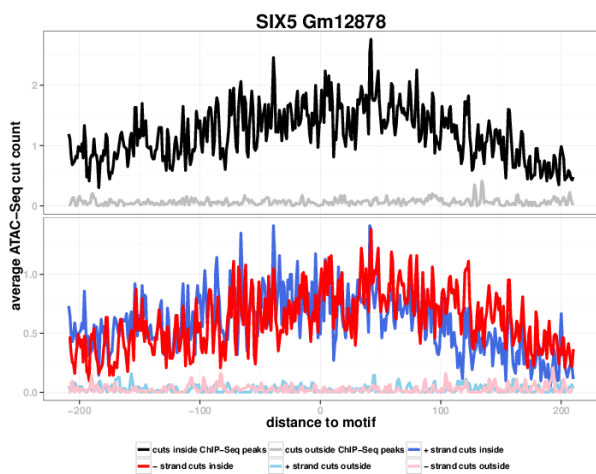
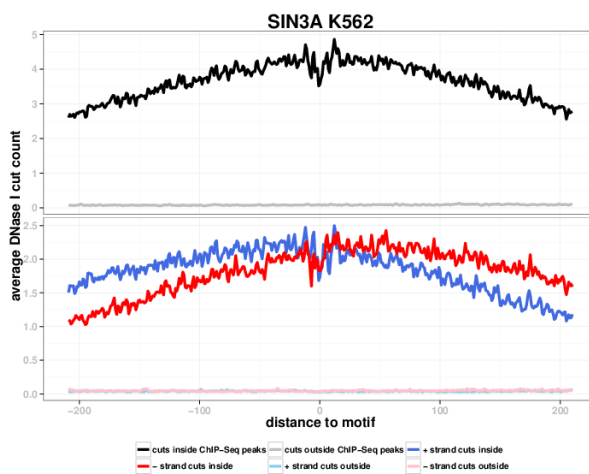
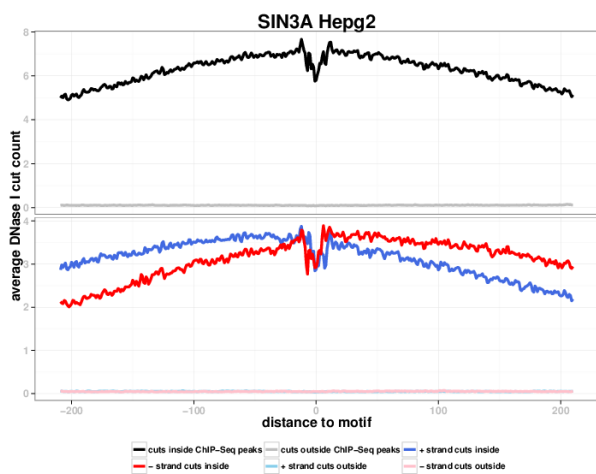


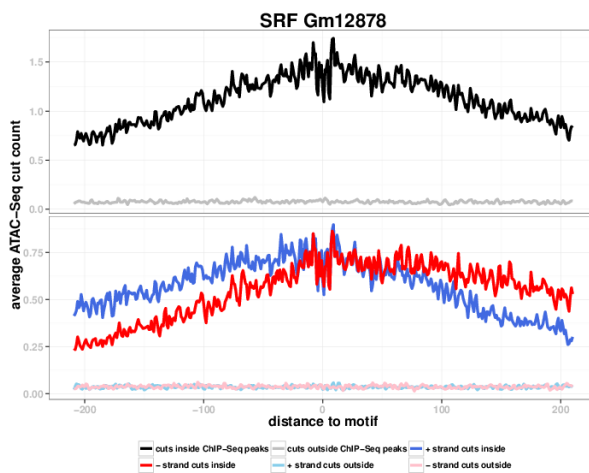
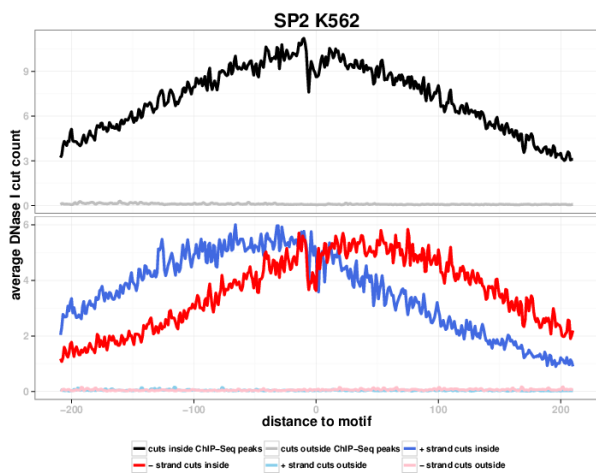
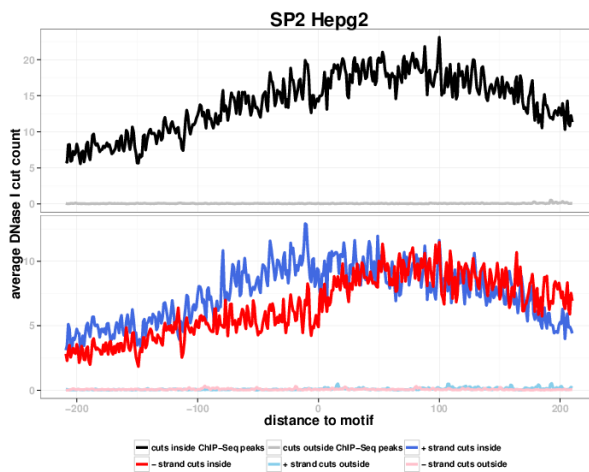
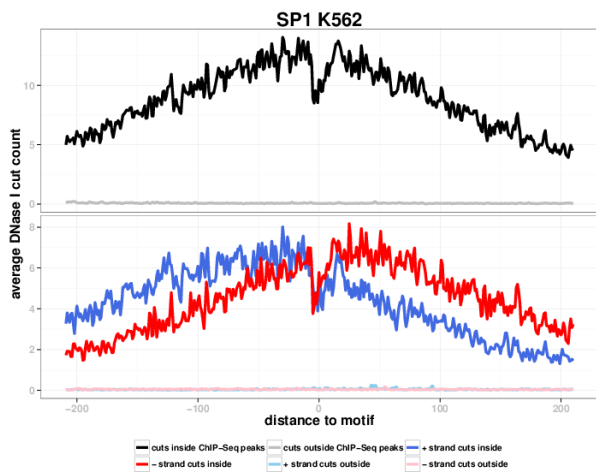
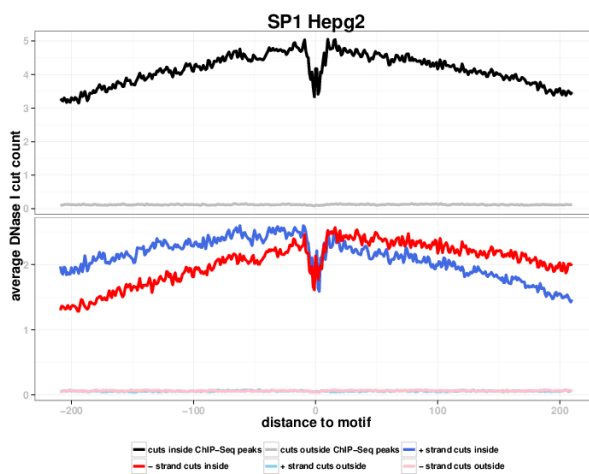
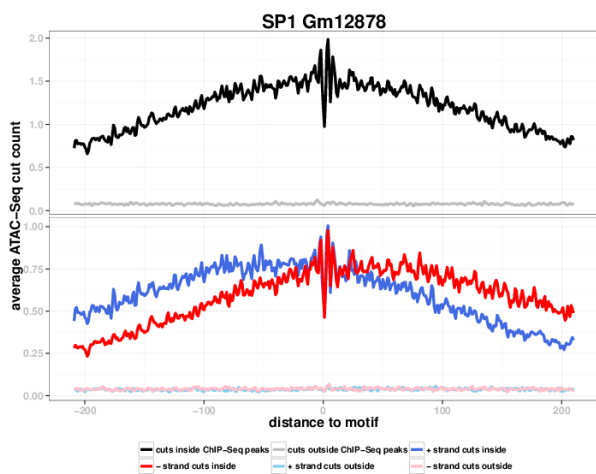




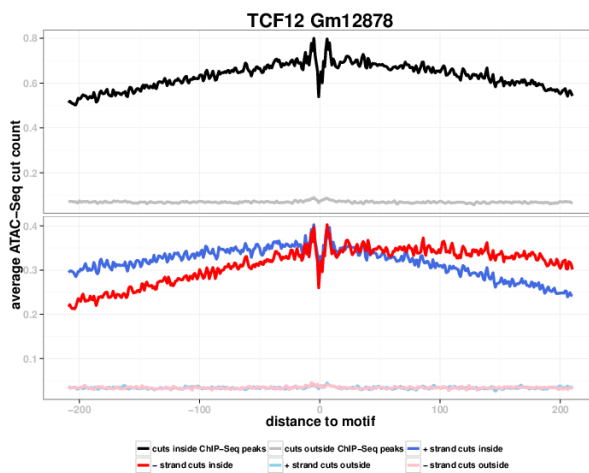
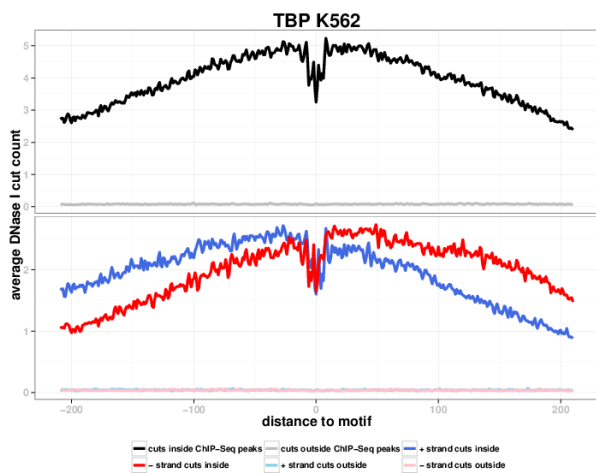
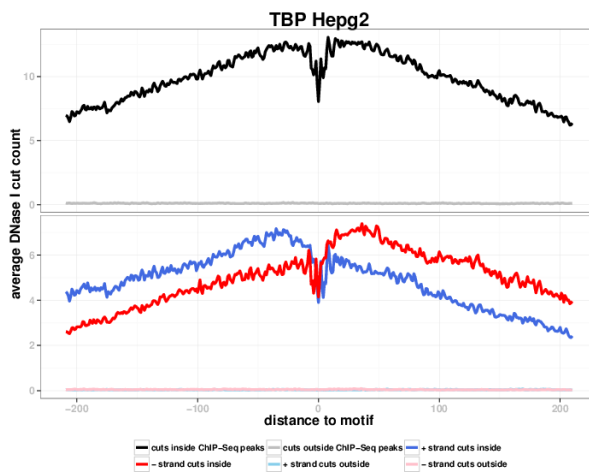
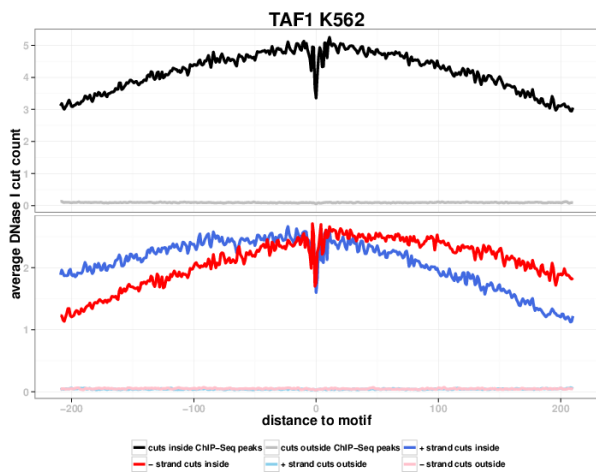
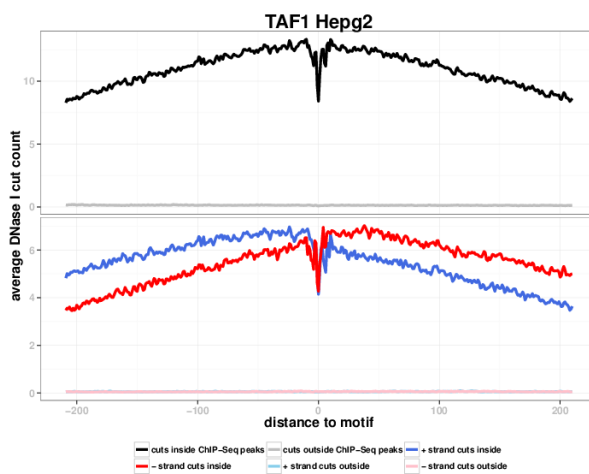
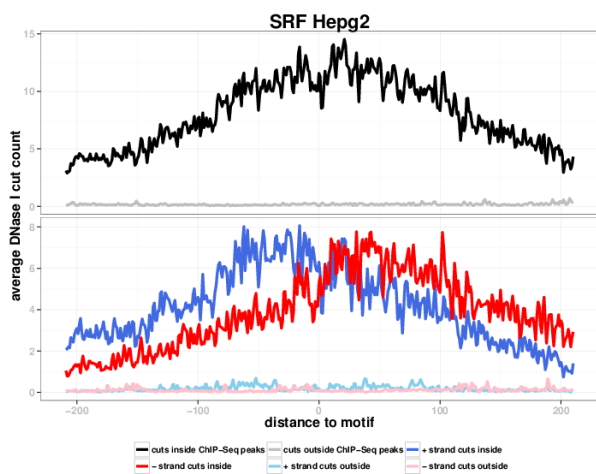


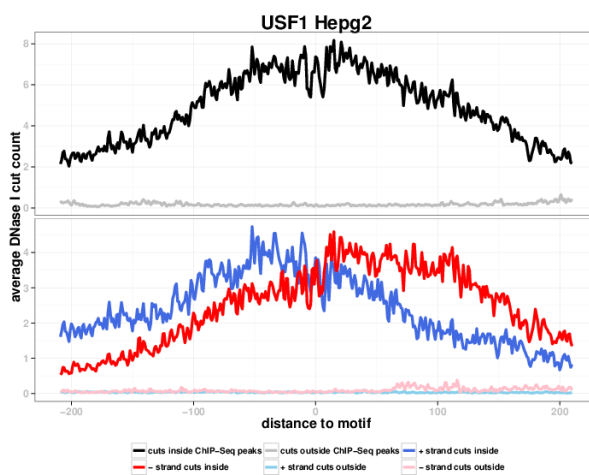
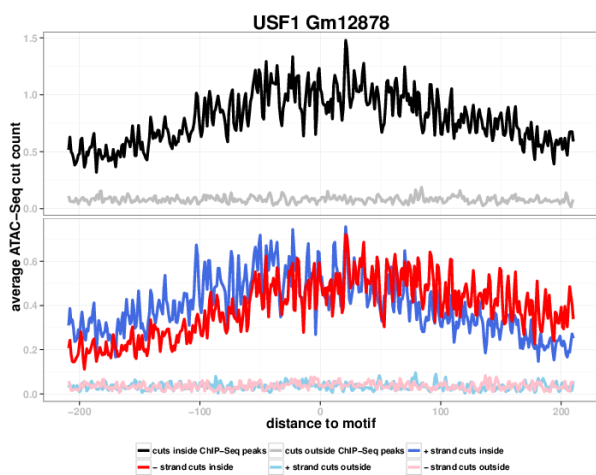
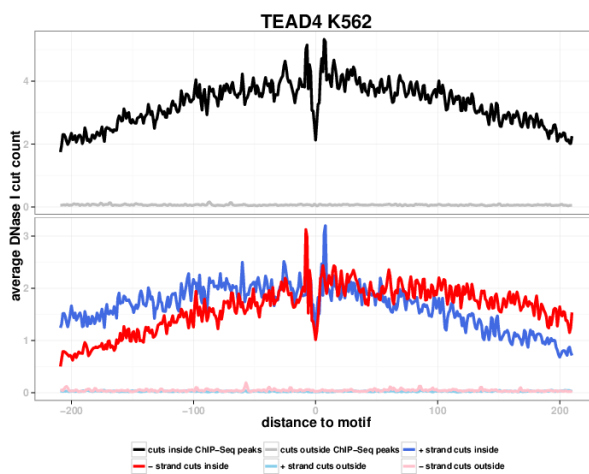
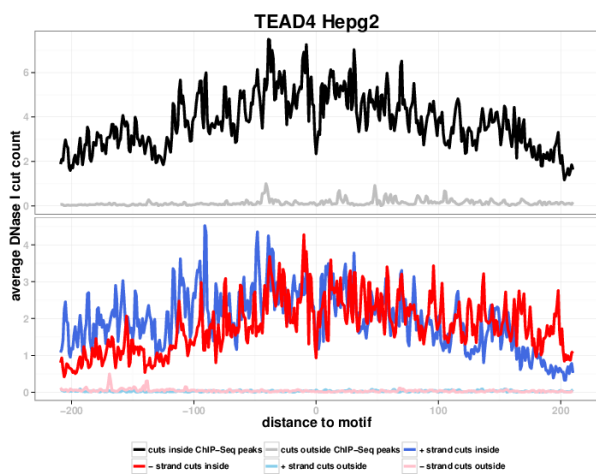
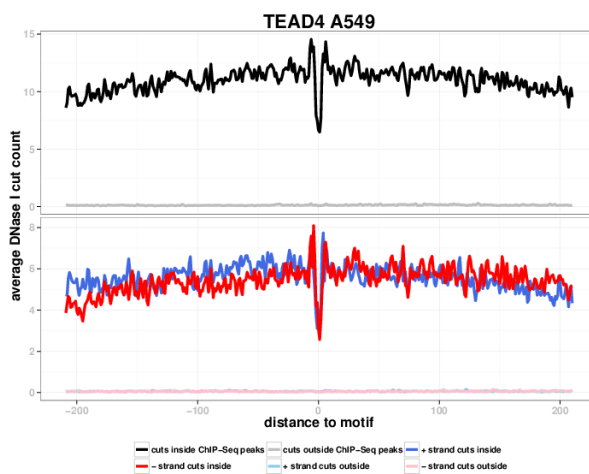
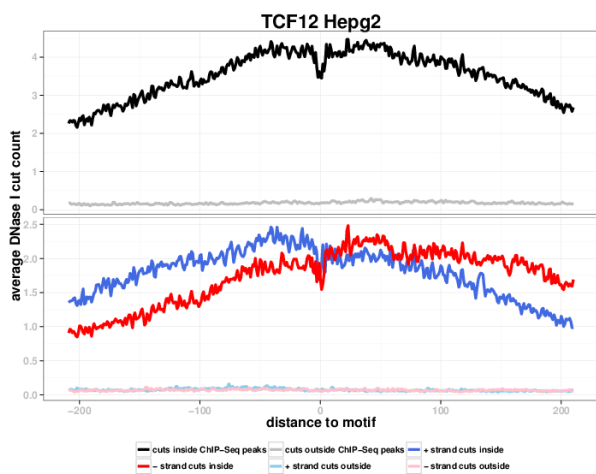


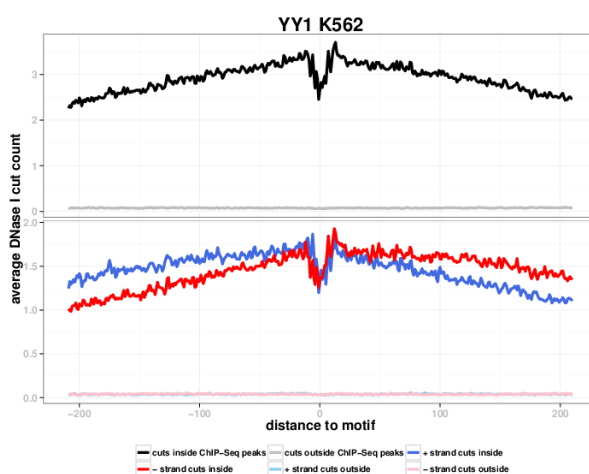
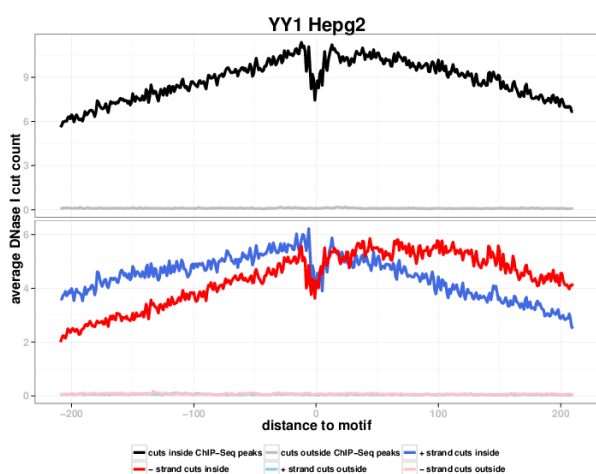
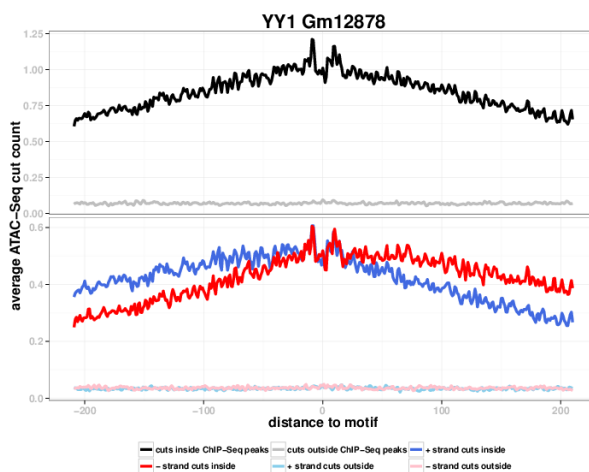
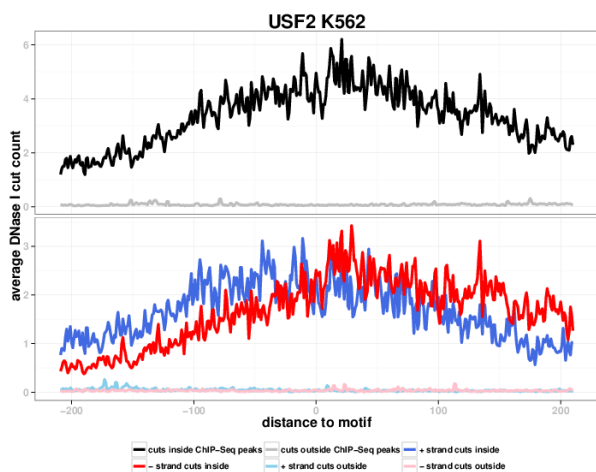
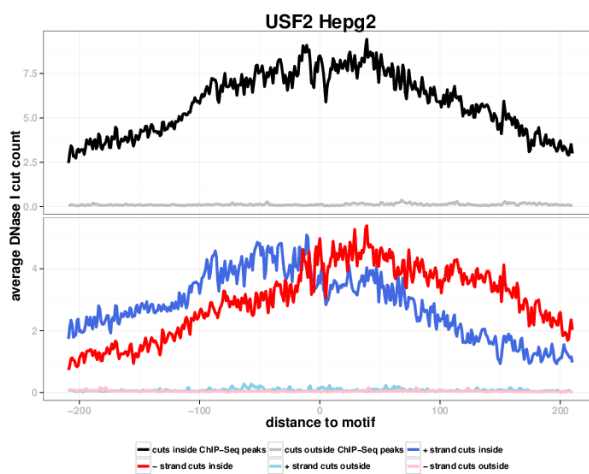
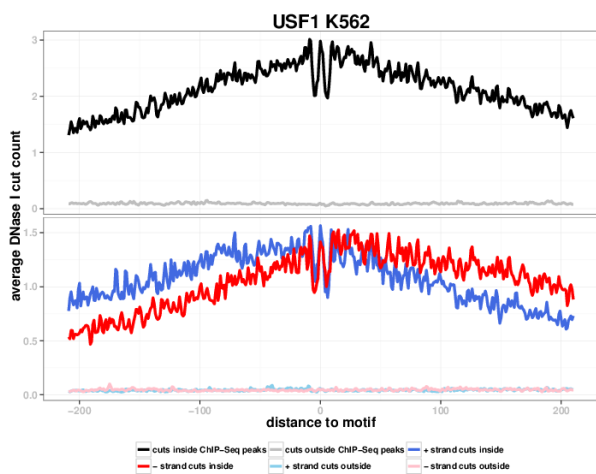


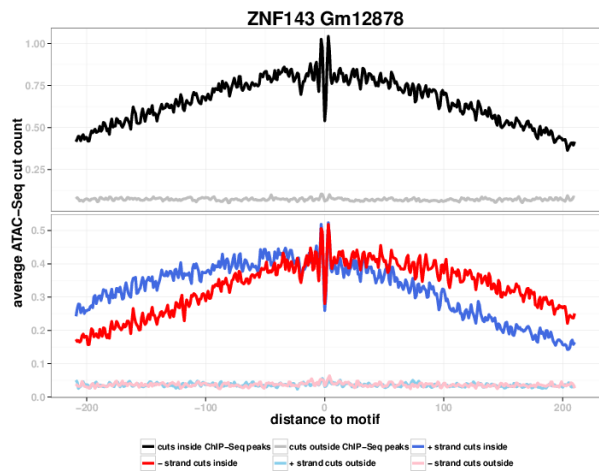
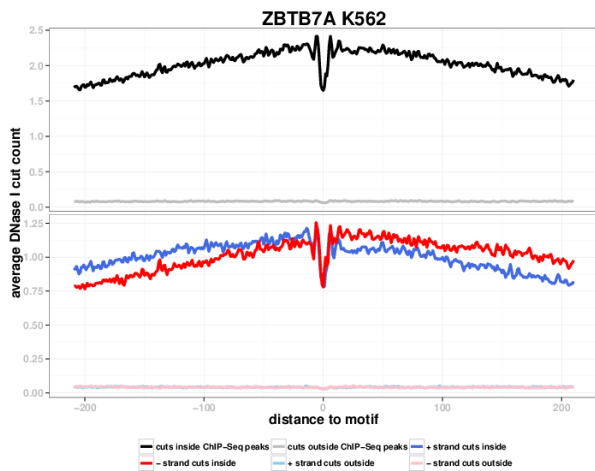
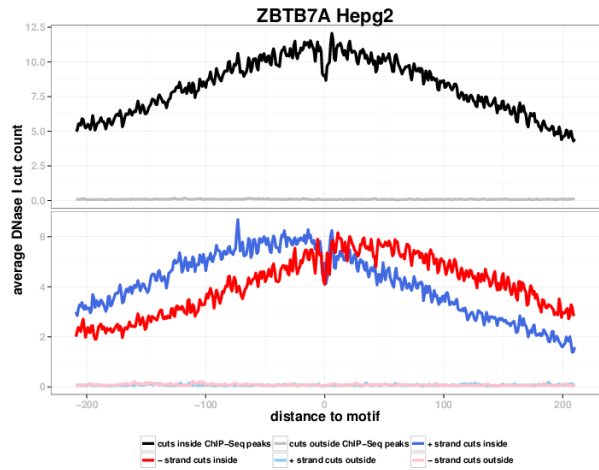
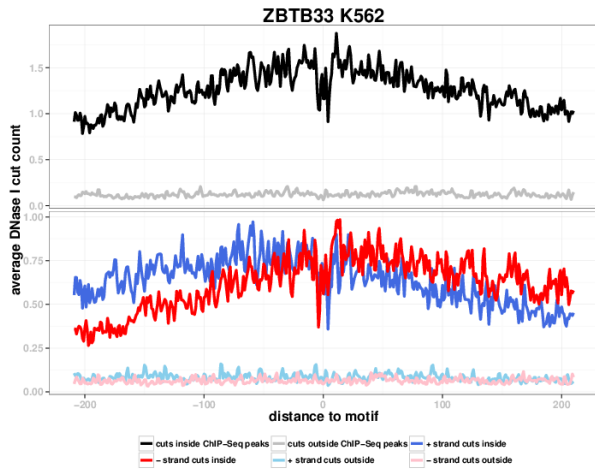
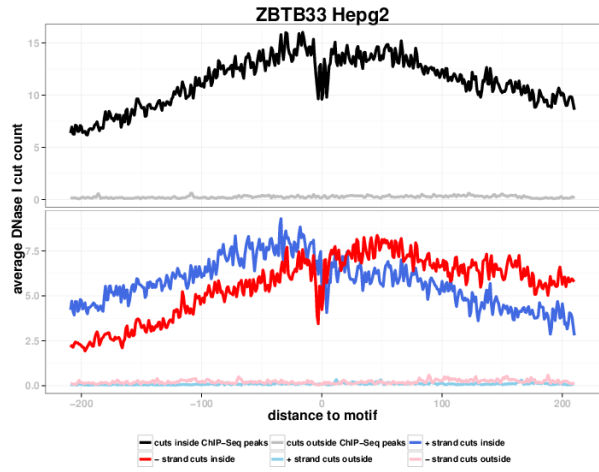
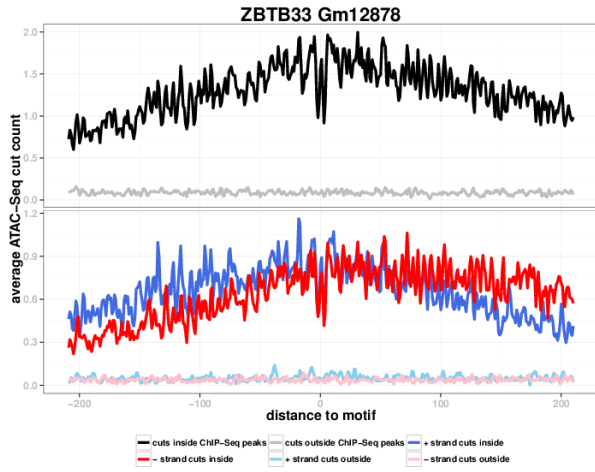


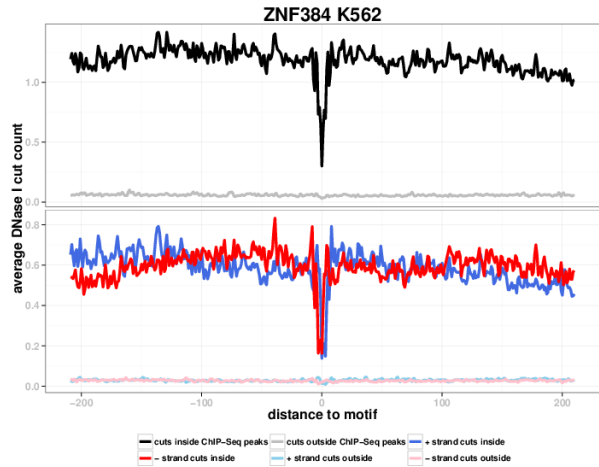
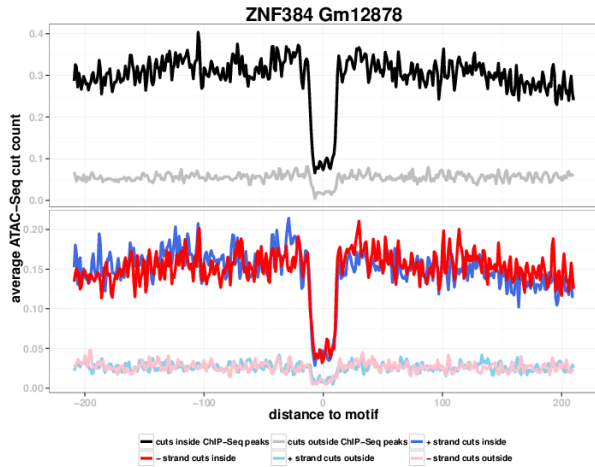
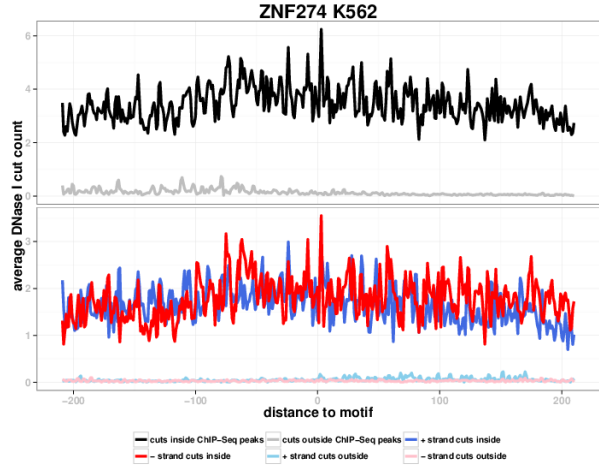
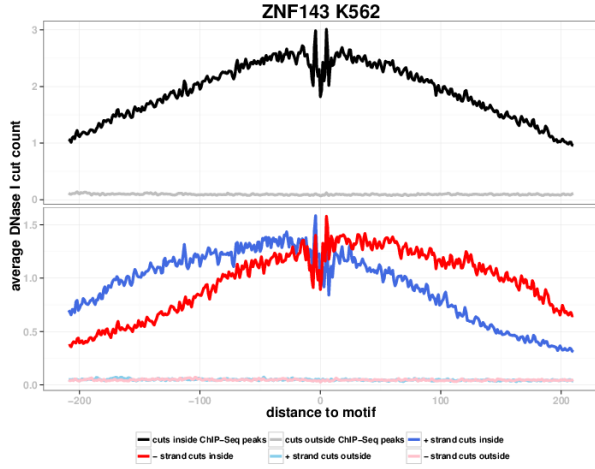












## Appendix B Sparse logistic regression training plots

