

Natural genetic variation affecting transcription factor spacing at regulatory regions is generally well tolerated

Zeyang Shen^{1,2}, Jenhan Tao¹, Gregory J. Fonseca¹ and Christopher K. Glass^{1,3,4}

¹Department of Cellular and Molecular Medicine, School of Medicine

²Department of Bioengineering, Jacobs School of Engineering

³Department of Medicine, School of Medicine, University of California, San Diego, La Jolla, CA, 92093, USA

⁴ckg@ucsd.edu

Abstract

Regulation of gene expression requires the combinatorial binding of sequence-specific transcription factors (TFs) at promoters and enhancers. Single nucleotide polymorphisms (SNPs) and short insertions and deletions (InDels) can influence gene expression by altering the sequences of TF binding sites. Prior studies also showed that alterations in the spacing between TF binding sites can influence promoter and enhancer activity. However, the relative importance of altered TF spacing has not been systematically analyzed in the context of natural genetic variation. Here, we exploit millions of InDels provided by five diverse strains of mice to globally investigate the effects of altered spacing on TF binding and local histone acetylation in macrophages. We find that spacing alterations resulting from InDels are generally well tolerated in comparison to genetic variants that directly alter TF binding sites. These findings have implications for interpretation of non-coding genetic variation and comparative analysis of regulatory elements across species.

Introduction

Genome-wide association studies (GWASs) have identified thousands of genetic variants associated with diseases and other traits (MacArthur *et al.*, 2017, Visscher *et al.*, 2017). Many of these variants fall into regulatory regions of the genome, implicating their effects on gene regulation (GTEx Consortium, 2015, Farh *et al.*, 2015). Gene expression is regulated in a cell-type-specific manner by transcription factors (TFs) that bind to short, degenerate sequences in promoters and enhancers referred to as TF binding motifs. Active promoters and enhancers are selected by combinations of sequence-specific TFs that bind in an inter-dependent manner to closely spaced motifs. Genetic variation that creates or disrupts TF binding motifs is a well-established mechanism for altering gene expression and biological function (Grossman *et al.*, 2017, Deplancke *et al.*, 2016, Heinz *et al.*, 2015). Collaborative binding of TFs required for enhancer or promoter selection can interact over a relatively broad range of spacing (e.g., 100-200 bp; Slattery *et al.*, 2014, Jiang and Singh, 2014, Heinz *et al.*, 2010). Consistent with this, flexibility in motif spacing relationships has been demonstrated using reporter assays in *Drosophila* (Menoret *et al.*, 2013) and HepG2 cells (Smith *et al.*, 2013).

On the contrary, substantial evidence also showed that specific spacing relationships between motifs can be important for TF binding and function (Boeva, 2016). A special category is provided by TFs that form ternary complexes recognizing composite binding sites, exemplified by CAP-SELEX studies of 9,400 TF pairs (Jolma *et al.*, 2015), MyoD and other muscle-specific factors in muscle cells (Nandi *et al.*, 2013), Sox2 and Oct4 in embryonic stem cells (Rodda *et al.*, 2005), Ets and E-box in haematopoietic cells (Ng *et al.*, 2014), etc. Similar constrained spacing between independent motifs are required for the optimal binding

and function of interacting TFs at the interferon- β enhanceosome (Panne, 2008). In addition, reporter assays examining synthetic alterations of motif spacing between collaborative factors revealed examples for high sensitivity of gene expression on spacing in Ciona (Farley *et al.*, 2015). However, these studies did not distinguish the impact of altered spacing on transcription factor binding or subsequent recruitment of co-activators required for gene activation. Moreover, it remains unknown the extent to which these findings are relevant to altered spacing resulting from natural genetic variation in human population or between animal species.

Bone marrow-derived macrophages (BMDMs) from genetically diverse strains of mice provide a powerful system for studying the genome-wide impact of natural genetic variation on gene regulation. Single nucleotide polymorphisms (SNPs) and short insertions and deletions (InDels) represent common forms of genetic variation in the genomes of different mouse strains (Keane *et al.*, 2011) and are associated with strain-specific variation in gene expression. SNPs and InDels could affect motif sequence and mutate a motif, while InDels could additionally change spacing between motifs. Initial studies in the BMDMs from two strains of mice used naturally occurring motif mutations to support a collaborative binding model between LDTFs (e.g., PU.1 and C/EBP β) and a hierarchical binding model between LDTFs and signal dependent transcription factors (SDTFs) (e.g., PU.1 and p65) (Heinz *et al.*, 2013). Subsequent studies leveraging more than 50 million SNPs and 5 million InDels from five mouse strains linked ~60% of strain-specific TF binding sites to mutated motifs (Link *et al.*, 2018a), suggesting a possibility for the remaining strain-specific sites to be impacted by InDels that alter motif spacing.

To investigate the effects of altered spacing on TF binding and function, we first characterized the genome-wide binding patterns of macrophage LDTFs and SDTFs based on their binding sites determined by chromatin immuno-precipitation sequencing (ChIP-seq). By leveraging the local genetic variation at the TF binding sites from the five strains of mice, we linked the alteration of motif spacing to the change of TF binding activity and local acetylation of histone H3 lysine 27 (H3K27ac), which is a histone modification that is highly correlated with enhancer and promoter function (Creyghton *et al.*, 2010). We find that InDels altering spacing between specific pairs of LDTFs and SDTFs can be associated with significant changes in their respective binding, but this relationship can largely be explained by effects of these InDels on the binding motifs of other collaborative factors, suggesting a general tolerance of spacing alterations resulting from natural genetic variation. These findings have implications for understanding mechanisms underlying enhancer selection, interpretation of non-coding variants associated with phenotypic variation, and comparisons of regulatory elements between species.

Results

Characterization of the spacing between macrophage LDTFs

As a starting point, we characterized the spacing relationship between the macrophage LDTFs, PU.1 and C/EBP β (**Figure 1A**), which have been found to bind in a collaborative manner at regulatory regions of macrophage-specific genes (Heinz *et al.*, 2010). We first determined reproducible PU.1 and C/EBP β binding sites from the replicate ChIP-seq data of C57BL/6J (C57) mice (Link *et al.*, 2018a) and then categorized them into three groups: co-bound by both factors, bound by PU.1 only, and bound by C/EBP β only (**Figure 1B**). For every binding site, we identified the DNA sequence best matching the motifs of PU.1 and C/EBP β as determined by position weight matrices (PWMs) (Stormo, 2000; **Materials and Methods**). We then computed the spacing (i.e., distance) between the centers of best-

matching sequences and plotted its distribution for sites within the same group. Co-binding sites showed a preference, but not strictness, for PU.1 and C/EBP β motifs to occur within ± 75 bp of each other (**Figure 1C**; **Figure 1—figure supplement 1**), in agreement with prior studies (Heinz *et al.*, 2010). Noticeably, a discontinuity occurred at where the two motifs overlap (spacing < 12 bp), potentially due to a steric inhibition of co-binding in these instances. For the sites bound by PU.1 or C/EBP β alone, the spacing relationship between PU.1 and C/EBP β motifs was statistically similar to the background distribution, consistent with few collaborative interactions between PU.1 and C/EBP β at these sites.

After observing the overall proximity between PU.1 and C/EBP β motifs at their co-binding sites, we investigated whether this spacing preference had an impact on TF binding. The binding activities of PU.1 and C/EBP β were quantified by ChIP-seq reads at the co-binding sites. We correlated the number of reads with either motif spacing or motif score, which represents the similarity of a sequence in comparison to PWMs. Both PU.1 and C/EBP β binding activities were positively correlated with the motif scores of their respective motifs but showed a much weaker correlation with spacing (**Figure 1D**). Interestingly, PU.1 binding activity is negatively correlated with C/EBP β motif score, implicating a synergistic binding model between these two TFs, which would allow the recognition of more degenerate motif sequence when they bind together to DNA.

Effect of altered spacing on transcription factor binding based on natural genetic variation across mouse strains

Table 1. P-values and effect sizes for the effects of different genetic variation between C57 and PWK on PU.1 binding, C/EBP β binding, and H3K27ac.

	Mutated PU.1 motif		Mutated C/EBP β motif		Altered spacing (unfiltered)		Altered spacing (filtered by collabor. factors)		Altered spacing (filtered by unrelated factors)	
	<i>p</i>	<i>d</i>	<i>p</i>	<i>d</i>	<i>p</i>	<i>d</i>	<i>p</i>	<i>d</i>	<i>p</i>	<i>d</i>
PU.1 binding	$<1e-4$	[1.07, 1.12]	$<1e-4$	[0.71, 0.78]	$<1e-4$	[0.38, 0.58]	0.002	[0.16, 0.52]	$<1e-4$	[0.37, 0.66]
C/EBP β binding	$<1e-4$	[0.76, 0.84]	$<1e-4$	[1.03, 1.10]	$<1e-4$	[0.27, 0.50]	0.16	[0.04, 0.37]	$6e-4$	[0.20, 0.55]
H3K27ac	$<1e-4$	[0.56, 0.65]	$<1e-4$	[0.55, 0.65]	$<1e-4$	[0.24, 0.45]	0.17	[0.04, 0.29]	$9e-4$	[0.16, 0.45]

P-values are calculated based on 10,000 iterations of permutation tests by comparing the log fold changes of ChIP-seq reads against the variant-free category. P-value below $1e-4$ is beyond the specified testing power. Effect sizes are represented by Cohen's *d*, displayed by its 90% confidence interval, which is based on 10,000 iterations of sampling the variant-free regions. For the comparison of different filters for InDels altering spacing, we computed the statistics for the complete set (i.e., unfiltered) and the filtered sets after excluding InDels simultaneously mutating motifs of collaborative factors or non-collaborative factors.

To investigate the requirement and tolerance of spacing proximity observed for collaborative factors, we leveraged the natural genetic variation across genetically diverse strains of mice as a mutagenesis screen. We selected five strains from which the ChIP-seq data of macrophage LDTFs and SDTFs were previously published (Link *et al.*, 2018a): C57BL/6J (C57), BALB/cJ (BALB), NOD/ShiLtJ (NOD), PWK/PhJ (PWK), and SPRET/EiJ (SPRET). Independent comparisons were conducted between C57 and one of the other four strains, which provide 4-40 million SNPs and 1-4 million InDels with respect to C57 (Keane *et al.*, 2011). We first identified the co-binding sites of PU.1 and C/EBP β for every strain based on ChIP-seq data. For each pairwise analysis, the co-binding sites from C57 and the compared strain were pooled and distributed into four categories based on the impacts of local variants: altered spacing, mutated PU.1 motif, mutated C/EBP β motif, and variant-free (**Figure 2A**; **Materials and Methods**). We quantified the effects of genetic variation on TF binding based on log2 fold changes of ChIP-seq reads between the compared strains. PU.1 binding is

significantly affected by mutated PU.1 motif and mutated C/EBP β motif, which have the largest and second largest effect sizes represented by Cohen's d (Table 1; Figure 2B; Figure 2—figure supplement 1). Similarly, C/EBP β binding is most significantly affected by mutated C/EBP β motif, followed by mutated PU.1 motif (Table 1; Figure 2C; Figure 2—figure supplement 1). Altered spacing resulting from InDels showed a much smaller, but still significant effect on both PU.1 and C/EBP β binding (Table 1 “unfiltered”). In many cases, an alteration of several nucleotides between PU.1 and C/EBP β motifs has no observable effect on TF binding, while one SNP that alters the core sequence of a motif can disrupt TF binding (Figure 2—figure supplement 2).

After observing a significant effect of altered spacing on the binding of PU.1 and C/EBP β , we investigated whether the effect size is influenced by the scale or direction of spacing alterations. By correlating the change of binding activity with the size of InDels (positive for insertions and negative for deletions), spacing alteration demonstrated an effect independent from the scale or direction of InDels (Figure 2D; Figure 2—figure supplement 3). On the contrary, changes of motif score are strongly correlated with changes of ChIP-seq reads, consistent with the important role of motif for TF binding. The invariable effects of InDels were unexpected because, based on the spacing relationship of PU.1 and C/EBP β , we expected a preference for closer spacing and a larger effect from longer InDels. However, InDels altering 1 or 2 bp between motifs can often have an effect as large as relatively long InDels, and such effect is not affected by the initial spacing between motifs (Figure 2—figure supplement 4), suggesting that the significant effect of InDels might not be directly resulted from the alteration of spacing but from other reasons.

InDels that alter spacing may simultaneously mutate motifs of other collaborative factors

One possibility for seeing a significant effect of InDels that reside between PU.1 and C/EBP β motifs could be alterations of motifs recognized by other collaborative factors. To test this hypothesis, we developed a computational framework to confidently identify collaborative factors. Considering that it would be a vast undertaking to perform ChIP-seq on all expressed TFs, our framework leverages TF binding sites identified from single ChIP-seq data and predicts collaborative factors based on high-score and closely spaced motifs (Figure 3A). This design is supported by our observations on PU.1 and C/EBP β binding sites where an increasing threshold on the motif score of collaborative factors recovered a larger proportion of co-binding sites (Figure 3—figure supplement 1) and also recovered the spacing relationships previously identified from ChIP-seq data (Figure 3—figure supplement 2). To compare the spacing distribution predicted by our framework to the distribution identified from co-binding sites, we tested on PU.1 binding sites measured by ChIP-seq and searched for high-score motif of cJun (i.e., FOS::JUN or TGAG/CTCA), which is a known LDTF of macrophages and a collaborative factor of PU.1. The predicted distribution of cJun around PU.1 aligned well with experimentally determined distribution based on cJun ChIP-seq (Link *et al.*, 2018a), showing the utility of identifying collaborative factors based on closely spaced high-score motifs (Figure 3B). Therefore, we applied this approach to uncover the collaborative factors of PU.1 and C/EBP β from over five hundred TFs whose motifs are available in the JASPAR database (Fornes *et al.*, 2020). To facilitate the comparison across motifs, we used the top 4,000 regions ranked by the motif score of every computed motif to obtain spacing distribution and compared each distribution against the background distribution using KS tests. P-values from KS tests were given signs to distinguish positive or negative associations with proximal spacing. Most TFs indicate no spacing relationship with either PU.1 or C/EBP β (Figure 3C; Figure 3—figure supplement 3). Motifs with proximal

spacing relationships tend to have relatively high expression based on RNA-seq data (Link *et al.*, 2018a).

Based on our computational framework, we selected twelve predicted collaborative factors, which are closely spaced with PU.1 or C/EBP β (KS test p-value < 1e-5) and also highly expressed in mouse macrophages (TPM value > 16). We refined the testing of co-binding sites by filtering out those with motif mutations of any collaborative factors on at least one core position (roughly equivalent to a change of motif score greater than 1). The remaining sites with InDels between PU.1 and C/EBP β motifs, which should represent a clean set of spacing alterations, showed a diminished effect on TF binding (**Table 1** “filtered by collabor. factors”; **Figure 3D**). When we filtered on unrelated factors identified as non-collaborative by our framework, the effect sizes were not affected (**Table 1** “filtered by unrelated factors”; **Figure 3—figure supplement 4**).

To investigate whether the effects of altered spacing on PU.1 and C/EBP β binding can be generalized to hierarchical interactions with signal-dependent transcription factors, we repeated our analyses on another pair of TFs, PU.1 and the NF κ B subunit p65. Upon macrophage activation with the TLR4-specific ligand Kdo2 lipid A (KLA), p65 enters the nucleus and primarily binds to poised enhancer elements that are selected by pioneering factors including PU.1. By leveraging the ChIP-seq data of the two TFs from C57 and PWK macrophages treated for one hour with KLA (Link *et al.*, 2018a), we observed a preference for proximal spacing between PU.1 and p65 motifs at their co-binding sites (**Figure 3—figure supplement 5**) and a diminished effect of altered spacing after excluding InDels that affect motifs of the predicted collaborative factors of PU.1 (**Figure 3E**; **Figure 3—figure supplement 6**), consistent with our finding from PU.1 and C/EBP β that spacing alterations are well tolerated by TF binding.

Effect of altered spacing on promoter and enhancer function

Although alterations in motif spacing were generally well tolerated at the level of DNA binding, it remained possible that changes in motif spacing could influence subsequent steps in enhancer and/or promoter activation. To examine this, we extended our analysis to local histone acetylation as a surrogate of promoter and enhancer function. We leveraged the H3K27ac ChIP-seq data for the five strains of mice (Link *et al.*, 2018a) and calculated the log fold changes of H3K27ac level within the extended 1000-bp regions of the PU.1 and C/EBP β co-binding sites. Testing all the co-binding sites demonstrated significant effects of both spacing alteration and motif mutation (**Table 1**; **Figure 4A**; **Figure 4—figure supplement 1**). However, the significance for altered spacing disappeared after filtering out sites potentially having motif mutations for the previous twelve collaborative factors (**Table 1**; **Figure 4B**; **Figure 4—figure supplement 2**). Again, filtering for unrelated factors did not influence the effect size (**Table 1**; **Figure 4—figure supplement 3**). The tolerance of spacing alteration was further reflected by a weak correlation between the change of acetylation level and the size of InDels, in comparison to a much stronger correlation with motif scores of both PU.1 and C/EBP β (**Figure 4C**). Similar to what was observed for TF binding, altered spacing demonstrated trivial effects on histone acetylation, which is supported by the high consistency between change of TF binding and change of acetylation (**Figure 4—figure supplement 4**). Noticeably, the acetylation level has an overall smaller scale of change compared to TF binding activity, reflecting its more complex dependency on TF binding (Reiter *et al.*, 2017).

Consideration of gap penalties in cross species sequence alignments

The insignificant effect of InDels between TF binding sites on TF binding and local histone acetylation at a genome-wide scale suggested that evolutionary pressure on enhancer selection and function would be relatively tolerant of these forms of genetic variation in comparison to InDels that directly affect the sequences of TF binding motifs. This is in contrast to effects of InDels in protein coding regions of the genome, in which insertions or deletions of bases other than multiples of three would result in frame-shift mutations. To explore this possibility, we performed sequence alignments using BLAST (Boratyn *et al.*, 2013) for well-established regulatory elements of macrophage-specific genes in the mouse with the human genome using i) standard parameters that impose significant penalties for gaps or ii) lenient parameters in which gap penalties were diminished. These comparisons frequently resulted in relatively short sequence alignments when standard gap penalties were applied but much more extended alignments that contained multiple relevant TF binding motifs using lenient gap penalties (**Figure 5A; Figure 5—figure supplements 1 and 2**). Examples are provided for putative regulatory elements of genes with known functions in macrophages (**Figure 5B; Figure 5—figure supplement 3**), including *Anxa7* (Li *et al.*, 2013), *Fos* (Hop *et al.*, 2018), *Vmp1* (Dziuba *et al.*, 2012), *Max* (Ayer *et al.*, 1993), and *Sema4d* (Li *et al.*, 2018). These regions are bound by PU.1 and C/EBP β in mouse BMDMs and are acetylated at H3K27 in both mouse BMDMs and human monocytes. A standard alignment of the human genome using 300-bp sequences from the mouse genome resulted in homologies that contain neither PU.1 nor C/EBP β motif (**Figure 5C; Figure 5—figure supplement 4**). In contrast, lenient gap penalties captured much more extended regions, containing high-score motifs of both PU.1 and C/EBP β . The motif sequences of PU.1 and C/EBP β are well preserved between human and mouse, but the motif spacing is altered by 1-6 bp, further supporting the general tolerance of spacing alterations.

Discussion

We investigated the global dependencies of collaborative TFs on spacing, using LDTFs and SDTFs of macrophages as the study model. PU.1 and C/EBP β demonstrated a preference for proximal motif spacing at their co-binding sites, but this preference for proximal spacing is not a strong modifier of TF binding in comparison to the high correlation between motif scores and TF binding activities. By leveraging natural genetic variation across genetically diverse strains of mice, we revealed the effects of spacing alterations and motif mutations on TF binding and function. InDels that alter spacing between PU.1 and C/EBP β motifs were associated with a smaller, but significant, change of TF binding and histone acetylation compared to motif mutations. However, by excluding InDels that potentially affect motifs of other collaborative factors identified by our newly developed framework, we observed an insignificant effect of the remaining sites. This finding suggests that the significant effects observed for InDels at some sites are very likely due to the motif mutations of other collaborative factors instead of spacing alterations between PU.1 and C/EBP β motifs. This result is consistent with the slope seen in the spacing distributions of PU.1 and C/EBP β at their co-binding sites (**Figure 1C**). For example, an InDel resulting in a change in spacing from 20 bp to 30 bp would still place the motifs well within the range of collaborative interactions. Similar relationships were observed for PU.1 and cJun, and for PU.1 and p65. Although these relationships are likely to be general, studies of additional LDTFs and SDTFs in other cell types will be required to establish this point.

These findings provide evidence that a subset of transcriptional regulatory elements does not require strict spacing relationships between transcription factors, in contrast to the examples provided by functional and structural studies of the interferon- β enhanceosome (Panne, 2008) and demonstrated in vivo in the case of synthetically modified enhancer elements in Ciona

(Farley *et al.*, 2015). However, these two examples represent regulatory elements in which key TF motifs are tightly spaced in their native contexts (i.e., 6-13 bp between motif centers). Direct protein-protein interactions are observed between bound TFs at the interferon- β enhanceosome, analogous to interactions defined for cooperative TFs that form ternary complexes (Morgunova and Taipale, 2017, Reményi *et al.*, 2003). Insertions or deletions between these tightly spaced motifs may result in sequence alterations as well as the potential for steric inhibition of DNA binding. Consistent with this, spacing distributions for most collaborative TFs exhibit a discontinuity at spacings of less than 12 bp between motif centers due to overlap of their sequences (**Figure 1C**). The present studies were thus not able to distinguish effects of spacing from effects of motif mutations below this motif distance threshold.

Another question raised by the discrepancy between the spacing dependencies discovered by previous studies and the spacing tolerance concluded by the present studies is the relative proportion of regulatory elements overall in which strict spacing relationships have functional importance. The current studies are limited by the ~5 million InDels provided by five strains of mice. Of the approximately 14,000 genomic locations co-bound by PU.1 and C/EBP β and associated with local histone acetylation, informative InDels to test for impact of spacing (i.e., between PU.1 and C/EBP β motifs, not affecting other collaborative TF motifs, and not complicated by other variants) were present at ~300 sites, representing ~2% of these regions. While this set of genomic locations enabled clear conclusions based on comparisons to ~4000 variant free sites, the extent to which this set of binding sites is representative of all regulatory elements is unclear. In particular, the interferon- β enhancer is among many regulatory elements that have no InDels across the five mouse strains examined. It thus remains possible that a subset of enhancers is dependent on strict spacing relationships.

Regardless of the extent of potential spacing-dependent regulatory elements, the present studies provide strong evidence that naturally occurring alterations in spacing between TF binding sites within putative regulatory elements are generally well tolerated. The conclusions are likely transferrable to explain the effects of InDels observed in human genomes, considering the similar number and size of InDels observed in human population (Mills *et al.*, 2011). To leverage an additional source of genetic variation, we compared the regulatory elements of mouse macrophages lacking InDels to human genomic sequences. Standard gap penalties generally resulted in short sequence fragments, whereas more lenient penalties recovered extended regions of homology containing corresponding LDTF motifs. These findings support that InDels are tolerated by a large fraction of regulatory elements and provide a basis for decreasing gap penalties for sequence comparisons of putative regulatory elements across species. Nevertheless, these studies rely on natural genetic variation, which is subject to natural selection. It will therefore be of interest to systematically introduce variable sizes of InDels between LDTFs in representative variant free enhancers to obtain an unbiased answer to the generality of the tolerance of spacing alterations.

Materials and Methods

Sequencing data processing

The mouse sequencing data used in this study were downloaded from the GEO database with accession number GSE109965 (Link *et al.*, 2018a). We mapped the ChIP-seq reads using Bowtie2 v2.3.5.1 (Langmead and Salzberg, 2012) and mapped the RNA-seq reads using STAR v2.5.3a (Dobin *et al.*, 2013) all with default parameters. Data from C57BL/6J mice were mapped to mm10 genome. Reads from BALB/cJ, NOD/ShiLtJ, PWK/PhJ, and SPRET/EiJ were mapped to their respective genomes built by MMARGE v1.0 with default

variant filters and were then shifted to mm10 genome using MMARGE v1.0 “shift” function (Link *et al.*, 2018b) to facilitate comparison at homologous regions. The reproducible TF binding sites were identified from mapped ChIP-seq data by first using HOMER v4.9.1 (Heinz *et al.*, 2010) to call unfiltered 300-bp peaks (command “findPeaks -style factor -L 0 -C 0 -fdr 0.9 -size 200”) and then running IDR v2.0.3 (Li *et al.*, 2011) on replicates with default parameters. Gene expression was quantified by TPM to represent normalized RNA-seq reads mapped to exons using HOMER v4.9.1 (command “analyzeRepeats.pl rna mm10 -count exons -condenseGenes -tpm”). Activity of TF binding was quantified by the number of TF ChIP-seq reads within 300-bp TF binding sites normalized by library size using HOMER v4.9.1 (command “annotatePeaks.pl mm10 -norm 1e7”). Activity of promoter and enhancer was quantified by normalized H3K27ac ChIP-seq reads within extended 1000-bp regions around TF binding sites.

Motif score and motif spacing calculation

We extracted the DNA sequences of TF binding sites from the genomes of different mouse strains using the MMARGE v1.0 “extract_sequences” function (Link *et al.*, 2018b). Based on DNA sequences, we computed motif scores and identified TF binding motifs as previously described (Shen *et al.*, 2020). Generally, we first calculated dot products between position weight matrices (PWMs) and sequence vectors using Biopython package (Cock *et al.*, 2009). PWMs for PU.1, C/EBP β , and over 500 other TFs were obtained from the JASPAR vertebrate core database (Fornes *et al.*, 2020). Then the highest score for each PWM and its position across 300 bp were recorded to represent the entire sequence. Changes of motif scores were computed between the highest motif scores in two compared strains at the same regions. To obtain the confident binding positions of the measured TFs, we excluded TF binding sites whose corresponding motifs are larger than 40 bp away from the peak centers or have a score lower than zero (i.e., less likely to occur than random chance). Approximately 70% of total peaks passed these criteria for both PU.1 and C/EBP β . Motif spacing was calculated from center of one motif to another, but only for sites whose highest motif scores are greater than zero.

Background sequence generation

We generated background sequences by shuffling the sequences of TF binding sites in a unit of dimers, which can well preserve the GC content. We then manually replaced the central part of each background sequence with a TF binding motif by sampling nucleotides based on the probabilities in its PWM. Motif score and motif spacing were calculated in the same way for these shuffled sequences as for the TF binding sites.

Categorization of regions based on genetic variation

To investigate the effects of genetic variation, we separated the PU.1 and C/EBP β co-binding sites into four categories. “Mutated PU.1” and “Mutated CEBPB” include sites with variants that change the motif scores of PU.1 and C/EBP β motifs, respectively. “Altered spacing” category includes sites where InDels exist between PU.1 and CEBPB motifs, which are not altered by any other variant. Co-binding sites classified into these three categories all experience a single impact from their local variants (either altered spacing or mutated motif, not both) so that the effect size can be clearly traced. “Variant free” is the control category, which contains sites with no genetic variation. The information about genetic variation across mouse strains were extracted using MMARGE v1.0 “mutation_info” function (Link *et al.*, 2018b).

Statistical testing of effect size

Effect size of genetic variation was computed by the ratio of ChIP-seq read counts between two compared strains followed by log2 transformation. We conducted permutation tests with 10,000 iterations to compare the absolute log ratios between “Variant free” and other categories. During every iteration, we randomly selected a comparable size of regions from the “Variant free” category and computed the mean of the selected set. Based on 10,000 mean values, we generated the null distribution and computed the percentile of the mean from the testing category on the null distribution as p-value. We also obtained the Cohen’s *d* between the sampled variant-free set and the testing category as the effect size (Sullivan and Feinn, 2012) and summarized the 90% confidence interval from 10,000 *d* values.

Identification of collaborative factors based on motif score and spacing

The TF binding sites identified from ChIP-seq data were first centered around the corresponding motif based on the highest motif score within 300-bp regions. Again, we filtered out those with motif score below zero or motif located more than 40 bp away from peak center. Next, we searched for the motifs of other TFs within ± 150 bp. If the motif has a score greater than zero and does not overlap with the motif of the bound TF, we compute the distance from motif center to region center (i.e., center of the bound TF motif) and obtain a predicted spacing distribution by aggregating all the distances for each motif. The predicted distribution is further smoothed by a sliding average window of 8 bp for visualization. Each spacing distribution is compared to the distribution obtained from background sequences with the Kolmogorov-Smirnov test (KS test) using Scipy package (Virtanen *et al.*, 2019). We conducted KS test for both halves of the distribution, upstream and downstream, generating two p-values for each motif. The mean p-values are used to represent the significance of dissimilarity from background distribution. Additionally, we gave signs to the p-values depending on whether more distances occur within or beyond 75 bp. Positive sign shows a preference for close spacing while negative sign represents inhibition of close spacing. Collaborative factors are predicted to have preference for close spacing. During the analyses of PU.1 and C/EBP β binding, the twelve predicted collaborative TFs used to filter for InDels include IRF3, E2F6, SP1, ATF4, USF family (USF1, USF2), ETS family (ELF4, ETV6, ELK4), and AP-1 family (FOS::JUN, FOSL2::JUN, JDP2), while the unrelated factors used as controls include EGR1, OLIG1, NEUROD2, STAT1, KLF13, CTCF, and BARHL2. During the analyses of PU.1 and p65 binding, ten out of the twelve predicted collaborative TFs were used after excluding USF family, which was only predicted to be collaborators of C/EBP β .

Sequence alignment between mouse and human

Among 3,917 variant-free PU.1 and C/EBP β co-binding sites merged from C57 and PWK, we quantified the H3K27ac level within the extended 1,000-bp regions and set a cutoff of H3K27ac ChIP-seq reads at 16 to obtain active regulatory elements. We extracted 300-bp sequences of these co-binding sites from the mm10 genome and aligned them to the hg38 genome using BLASTn algorithm (Boratyn *et al.*, 2013). Except for the different gap penalties (“Gap Costs” on the BLAST web interface) tested in our studies, the other parameters were used as default settings.

Acknowledgements

We would like to thank Leslie Van Ael for assistance with manuscript preparation.

439

440 **Competing interests**

441 None declared

442

Figure 1. Spacing relationship of PU.1 and C/EBP β . (A) Schematic of the collaborative binding model between PU.1 and C/EBP β , which recognize their own motifs spaced in macrophage-specific enhancers. (B) Numbers of singly binding and co-binding sites of PU.1 and C/EBP β identified from ChIP-seq data. (C) Distributions of C/EBP β motif around PU.1 binding sites. The distributions for non-overlapping sites (spacing > 12 bp) of each category were compared against the background distribution generated from shuffled sequences using Kolmogorov–Smirnov test (KS test). P-values from KS test are displayed in brackets. The spacing distributions were smoothed by an 8-bp sliding window for visualization purpose. (D) Hexbin plots showing the correlation between TF binding activity and motif spacing or motif score for the 9849 co-binding sites. Log₂ ChIP-seq reads were calculated within 300 bp to quantify the binding activity of PU.1 and C/EBP β . The color gradients represent the density of sites. Spearman correlation coefficients together with p-values are displayed to show the level of correlation.

The following figure supplements are available for figure 1.

Figure supplement 1. Spacing relationship of PU.1 and C/EBP β . (A) Spacing distributions of PU.1 motif around C/EBP β motif at co-binding sites and C/EBP β -singly-binding sites. P-values display the comparison against the background distribution using KS tests. (B) Spacing distributions regarding different orientation of the motifs. Co-binding sites and PU.1-singly-binding sites were divided into two subgroups representing same or opposite orientation of the PU.1 and C/EBP β motifs. The overall distributions are very similar for both subgroups.

Figure 2. Effects of spacing alterations resulting from natural genetic variation across mouse strains. (A) Schematic showing impacts of genetic variation on motif sequence or motif spacing. PU.1 and C/EBP β co-binding sites can be classified into four categories based on the impacts of local variants: “altered spacing”, “mutated PU.1”, “mutated C/EBP β ”, and “variant free”. (B, C) Absolute log₂ fold changes of ChIP-seq reads between C57 and PWK for (B) PU.1 binding and (C) C/EBP β binding. Boxplot shows the median and quartiles of every distribution with its sample size displayed on top. (*) indicates a significant effect size with $p < 0.001$ from permutation tests compared against the “variant free” category (Materials and Methods). (D) Correlations between change of C/EBP β binding and change of motif spacing or motif score. The co-binding sites used for change of spacing, PU.1 motif score, and C/EBP β motif score are from the previously defined categories “altered spacing”, “mutated PU.1”, “mutated C/EBP β ”, respectively.

The following figure supplements are available for figure 2.

Figure supplement 1. Change of PU.1 and C/EBP β binding affected by genetic variation for the other three pairwise comparisons. (*) indicates significance value $p < 0.001$ based on permutation test of every category against “variant free” category. The results from C57 vs. BALB and C57 vs. NOD are similar to what we saw for C57 vs. PWK. C57-SPRET comparison did not show significant results for “altered spacing” category, likely due to much more genetic variants between these two strains than other pairs, which introduced stronger trans effects to the “variant-free” category making the baseline effects high and potentially complicating the effects from InDels altering motif spacing.

Figure supplement 2. Example sites of motif mutation and spacing alteration. (A) a 5-bp increase in spacing has little effect on TF binding. (B) an A-to-G mutation on PU.1 motif yields a loss of both PU.1 and C/EBP β binding.

Figure supplement 3. Correlations between change of PU.1 binding and change in spacing or motif score. The co-binding sites here are from the C57-PWK comparison and are the same as those in Figure 2. Motif scores showed high correlation, while scale of spacing alteration is not associated with change of PU.1 binding.

Figure supplement 4. Effect size of genetic variation in relation with the initial spacing between PU.1 and C/EBP β motif. Co-binding sites from C57-PWK comparison are binned based on the initial motif spacing and then used to calculate the absolute log2 fold change between the two strains, which were aggregated to compute mean values for each bin. The effect size of InDels altering spacing is overall not affected by the initial motif spacing.

Figure 3. Refining InDels to exclude those potentially mutating motifs of other collaborative factors. (A) Overview of our newly developed framework for identifying collaborative factors from single ChIP-seq data. Given the binding sites for TF of interest, our method searches for other motifs and uses regions with high-score motifs to compute the spacing distribution, which is further compared against the background distribution using KS test. Those with significant proximal distribution are predicted collaborative factors. (B) Comparison between the actual spacing relationship obtained from co-binding sites and the predicted spacing distribution of cJun and PU.1. P-values from KS test by comparing to the background distribution are shown in brackets. (C) Signed p-values of over five hundred motifs for PU.1 and C/EBP β binding sites. Color gradients indicate the level of gene expression measured by RNA-seq and quantified by TPM. The complete list of p-values is available in Figure 3—source data 1. (D) Effect size of a refined set of PU.1 and C/EBP β co-binding sites for C57-PWK comparison. About half of the original “altered spacing” sites were excluded due to their impacts on at least one of the twelve predicted collaborative factor motifs. (*) indicates $p < 0.001$ based on permutation tests against the “variant free” category. (E) Effect size of refined PU.1 and p65 co-binding sites for C57-PWK comparison. “Altered spacing” category has excluded InDels that impacts motifs of the collaborative factors of PU.1 identified from our framework.

The following figure supplements are available for figure 3.

Figure supplement 1. Fractions of recovered co-binding sites by filtering with different motif score thresholds. (A) PU.1 binding sites identified from PU.1 ChIP-seq data were filtered with different thresholds on C/EBP β motif. (B) C/EBP β binding sites identified from C/EBP β ChIP-seq data were filtered with different thresholds on PU.1 motif. Both demonstrated an increase in fraction of co-binding sites by a larger threshold.

Figure supplement 2. Predicted spacing distributions of PU.1 and C/EBP β . Recovered from (A) PU.1 binding sites with top C/EBP β motif (CEBPB), and (B) C/EBP β binding sites with top PU.1 motif (SPI1). Both predicted distributions are similar to the spacing distribution obtained from the actual co-binding sites identified from PU.1 and C/EBP β ChIP-seq data.

Figure supplement 3. Examples of predicted spacing distributions. (A) PU.1 and GFI1 as an example of no spacing relationship, and (B) PU.1 and ZEB1 as an example of distant spacing relationship. P-values shown in brackets are obtained from KS tests by comparing to the background distribution (shuffled sequences) without assigning signs to distinguish proximal and distant spacing relationship.

Figure supplement 4. Fold changes of TF binding after filtering out mutations on non-collaborative factors. The remaining sites in “Altered spacing” category still have a significant effect on TF binding based on permutation tests ($p < 0.001$).

Figure supplement 5. Spacing relationship of PU.1 and p65 in mouse macrophages at pro-inflammatory state induced by KLA treatment for 1 hour. Co-binding sites show clear preference for PU.1 and p65 motifs to be proximal, while p65-singly-binding sites do not have the same preference. The distributions exclude sites where PU.1 and p65 motifs overlap with a shift of 3 or 4 bp (overlapping “GGAA”/“TTCC”).

Figure supplement 6. Fold changes of TF binding for four categories of PU.1 and p65 co-binding sites. “Altered spacing” includes all co-binding sites where InDels occur between PU.1 and p65 motifs and alter the motif spacing without considering any impact on motifs of

other collaborative factors. (*) indicates significance value $p < 0.001$ based on permutation test of every category against "variant free" category.

The following source data are available for figure 3.

Source data 1. Complete list of signed p-values indicating predicted spacing relationships.

Figure 4. Effects of spacing alteration on promoter and enhancer activity measured by local histone acetylation. (A, B) Absolute log2 fold changes of H3K27ac ChIP-seq reads between C57 and PWK for (A) unfiltered co-binding sites and (B) refined co-binding sites after excluding InDels that mutate motifs of potential collaborative factors. (*) indicates $p < 0.001$ based on permutation tests against the "variant free" category. (C) Correlations between change of H3K27ac level and change of motif score or motif spacing. The co-binding sites used here are unfiltered.

The following figure supplements are available for figure 4.

Figure supplement 1. Results from the other three pairwise comparisons on the change of H3K27ac level for four categories of PU.1 and C/EBP β co-binding sites. (*) indicates significance value $p < 0.001$ based on permutation test of every category against "variant free" category. Again, C57-SPRET comparison did not show significant results for "altered spacing" category, likely due to the much larger genetic diversity between these two strains, which complicates the effects from InDels altering motif spacing with trans effects from the variants nearby.

Figure supplement 2. Results from the other three pairwise comparisons on the change of H3K27ac level after filtering out InDels in "Altered spacing" category that impact motifs of predicted collaborative factors. "Altered spacing" category no longer shows a significant effect on the acetylation level.

Figure supplement 3. Change of H3K27ac level affected by genetic variation after filtering out InDels that mutate motifs of non-collaborative factors. The remaining sites in "Altered spacing" category still have a significant effect on local acetylation of H3K27 based on permutation tests ($p < 0.001$).

Figure supplement 4. Correlation between change of TF binding and change of H3K27ac level at all PU.1 and C/EBP β co-binding sites. Fold changes were calculated by dividing C57 by PWK. The larger fold change between PU.1 and C/EBP β binding was used for plotting. Overall, a strong correlation exists between H3K27ac level and TF binding, represented by a Pearson correlation of 0.7.

Figure 5. Implications of reducing gap penalties in cross species sequence alignments. (A) Lengths of aligned sequences for the co-binding sites of PU.1 and C/EBP β in mouse BMDMs that are enriched with local H3K27ac and have no genetic variation between C57 and PWK. 300-bp sequences were aligned to the human genome using BLAST with either the standard or the lenient gap penalties. The complete list of regions together with alignment results are available in Figure 5—source data 1. (B) Example showing a co-binding site within intron of *Anxa7* that was successfully aligned to a homology region in human. The aligned regions are enriched with H3K27ac in both mouse BMDMs and human monocytes. (C) The alignment results of the example co-binding site using standard or lenient gap penalties. Lenient gap penalties resulted in the recovery of a PU.1 motif and a C/EBP β motif. The spacing between PU.1 and C/EBP β motif centers is 39 bp in mouse and 42 bp in human, which did not impact the binding of PU.1 or the activity of this region.

The following figure supplements are available for figure 5.

Figure supplement 1. Alignment results of all PU.1 and C/EBP β co-binding sites in mouse macrophages compared with human genome. All the co-binding sites were identified as 300-bp regions in mouse genome and compared to the human genome using BLAST. Lenient gap

penalties resulted in longer alignments than standard parameters. Despite that, the majority of these co-binding sites have less than a third of the complete sequences aligned. By comparing to Figure 5A, our results suggest that acetylated co-binding sites are much more conserved between mouse and human than the rest with less acetylation and potentially less functional importance.

Figure supplement 2. Alignment results of PU.1 and C/EBP β co-binding sites that are enriched with local H3K27ac and have no genetic variation between C57 and PWK using other possible gap penalties. All the co-binding sites were identified as 300-bp regions in mouse genome and compared to the human genome using BLAST. The lenient gap penalty ({2, 2}) still produced much longer aligned sequences than other penalty options.

Figure supplement 3. Examples of PU.1 and C/EBP β co-binding sites in BMDMs that are aligned to homology regions in human using lenient gap penalties. (A) Enhancer closest to *Fos*, which encodes AP-1 family transcription factor and is known to be important for macrophage function. (B) Enhancer proximal to *Vmpl*, which has been found to be associated with inflammatory response of macrophages. (C) Intron of *Max*, which encodes a basic-helix-loop-helix-zipper protein and is found to accumulate during macrophage differentiation. (D) Intron of *Sema4d*, which is found to be regulated by macrophages in tumor.

Figure supplement 4. The alignment results of the example sites shown in Figure supplement 3. (A) *Fos* enhancer, (B) *Vmpl* enhancer, (C) *Max* intron, and (D) *Sema4d* intron. The following source data are available for figure 5.

Source data 1. List of aligned regions.

References

- Ayer, D. E., and Eisenman, R. N. (1993). A switch from Myc: Max to Mad: Max heterocomplexes accompanies monocyte/macrophage differentiation. *Genes & development*, 7(11), 2110-2119.
- Boeva, V. (2016). Analysis of genomic sequence motifs for deciphering transcription factor binding and transcriptional regulation in Eukaryotic cells. *Frontiers in Genetics*, 7, 24.
- Boratyn, G. M., Camacho, C., Cooper, P. S., Coulouris, G., Fong, A., Ma, N., Madden, T. L., Matten, W. T., McGinnis, S. D., Merezuk, Y., et al. (2013). Blast: a more efficient report with usability improvements. *Nucleic acids research*, 41(W1):W29–W33.
- Cock, P. J., Antao, T., Chang, J. T., Chapman, B. A., Cox, C. J., Dalke, A., Friedberg, I., Hamelryck, T., Kauff, F., Wilczynski, B., and De Hoon, M. J. (2009). Biopython: Freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics*, 25(11):1422–1423.
- Creyghton, M. P., Cheng, A. W., Welstead, G. G., Kooistra, T., Carey, B. W., Steine, E. J., Hanna, J., Lodato, M. A., Frampton, G. M., Sharp, P. A., Boyer, L. A., Young, R. A., and Jaenisch, R. (2010). Histone H3K27ac separates active from poised enhancers and predicts developmental state. *Proceedings of the National Academy of Sciences of the United States of America*, 107(50):21931–21936.
- Deplancke, B., Alpern, D., and Gardeux, V. (2016). The genetics of transcription factor DNA binding variation. *Cell*, 166(3):538–554.
- Dobin, A., Davis, C. A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M., and Gingeras, T. R. (2013). Star: ultrafast universal rna-seq aligner. *Bioinformatics*, 29(1):15–21.
- Dziuba, N., Ferguson, M. R., O'Brien, W. A., Sanchez, A., Prussia, A. J., McDonald, N. J., Friedrich, B. M., Li, G., Shaw, M. W., Sheng, J., Hodge, T. W., Rubin, D. H., Murray, J. L. (2012). Identification of cellular proteins required for replication of human immunodeficiency virus type 1. *AIDS research and human retroviruses*, 28(10), 1329-1339.
- Farh, K. K.-H., Marson, A., Zhu, J., Kleinewietfeld, M., Housley, W. J., Beik, S., Shores, N., Whitton, H., Ryan, R. J., Shishkin, A. A., et al. (2015). Genetic and epigenetic fine mapping of causal autoimmune disease variants. *Nature*, 518(7539):337–343.
- Farley, E. K., Olson, K. M., Zhang, W., Brandt, A. J., Rokhsar, D. S., and Levine, M. S. (2015). Suboptimization of developmental enhancers. *Science*, 350(6258):325–328.
- Fornes, O., Castro-Mondragon, J. A., Khan, A., van der Lee, R., Zhang, X., Richmond, P. A., Modi, B. P., Correard, S., Gheorghe, M., Baranašić, D., et al. (2020). Jasp2020: update of the open-access database of transcription factor binding profiles. *Nucleic acids research*, 48(D1):D87–D92.
- Grossman, S. R., Zhang, X., Wang, L., Engreitz, J., Melnikov, A., Rogov, P., Tewhey, R., Isakova, A., Deplancke, B., Bernstein, B. E., Mikkelsen, T. S., and Lander, E. S. (2017). Systematic dissection of genomic features determining transcription factor binding and enhancer function. *Proceedings of the National Academy of Sciences*, 114(7):E1291–E1300.
- GTEX Consortium (2015). The genotype-tissue expression (gtex) pilot analysis: multitissue gene regulation in humans. *Science*, 348(6235):648–660.
- Heinz, S., Benner, C., Spann, N., Bertolino, E., Lin, Y. C., Laslo, P., Cheng, J. X., Murre, C., Singh, H., and Glass, C. K. (2010). Simple Combinations of Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for Macrophage and B Cell Identities. *Molecular Cell*, 38(4):576–589.
- Heinz, S., Romanoski, C. E., Benner, C., Allison, K. A., Kaikkonen, M. U., Orozco, L. D., and Glass, C. K. (2013). Effect of natural genetic variation on enhancer selection and function. *Nature*, 503:487.9
- Heinz, S., Romanoski, C. E., Benner, C., and Glass, C. K. (2015). The selection and function of cell type-specific enhancers. *Nature Reviews Molecular Cell Biology*, 16(3):144–154.
- Hop, H.T., Arayan, L.T., Huy, T.X., Reyes, A.W., Vu, S.H., Min, W., Lee, H.J., Rhee, M.H., Chang, H.H. and Kim, S., 2018. The key role of c-Fos for immune regulation and bacterial dissemination in brucella infected macrophage. *Frontiers in Cellular and Infection Microbiology*, 8, 287.
- Jiang, P. and Singh, M. (2014). Ccat: combinatorial code analysis tool for transcriptional regulation. *Nucleic acids research*, 42(5):2833–2847.
- Jolma, A., Yin, Y., Nitta, K. R., Dave, K., Popov, A., Taipale, M., Enge, M., Kivioja, T., Morgunova, E., and Taipale, J. (2015). DNA-dependent formation of transcription factor pairs alters their binding specificity. *Nature*, 527(7578):384–388.
- Keane, T. M., Goodstadt, L., Danecek, P., White, M. A., Wong, K., Yalcin, B., Heger, A., Agam, A., Slater, G., Goodson, M., Furlotte, N. A., Eskin, E., Nellåker, C., Whitley, H., Cleak, J., Janowitz, D., Hernandez-Pliego, P., Edwards, A., Belgard, T. G., Oliver, P. L., McIntyre, R. E., Bhomra, A., Nicod, J., Gan, X., Yuan, W., Van Der

- Weyden, L., Steward, C. A., Bala, S., Stalker, J., Mott, R., Durbin, R., Jackson, I. J., Czechanski, A., Guerra-Assunção, J., Donahue, L. R., Reinholdt, L. G., Payseur, B. A., Ponting, C. P., Birney, E., Flint, J., and Adams, D. J. (2011). Mouse genomic variation and its effect on phenotypes and gene regulation. *Nature*, 477(7364):289–294.
- Langmead, B. and Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nature Methods*, 9:357.
- Li, H., Huang, S., Wang, S., Zhao, J., Su, L., Zhao, B., Zhang, Y., Zhang, S., and Miao, J. (2013). Targeting annexin A7 by a small molecule suppressed the activity of phosphatidylcholine-specific phospholipase C in vascular endothelial cells and inhibited atherosclerosis in apolipoprotein E^{-/-} mice. *Cell death & disease*, 4(9), e806-e806.
- Li, H., Wang, J. S., Mu, L. J., Shan, K. S., Li, L. P., and Zhou, Y. B. (2018). Promotion of Sema4D expression by tumor-associated macrophages: Significance in gastric carcinoma. *World journal of gastroenterology*, 24(5), 593.
- Li, Q., Brown, J. B., Huang, H., and Bickel, P. J. (2011). Measuring reproducibility of high-throughput experiments. *Annals of Applied Statistics*, 5(3):1752–1779.
- Link, V. M., Duttke, S. H., Chun, H. B., Holtman, I. R., Westin, E., Hoeksema, M. A., Abe, Y., Skola, D., Romanoski, C. E., Tao, J., Fonseca, G. J., Troutman, T. D., Spann, N. J., Strid, T., Sakai, M., Yu, M., Hu, R., Fang, R., Metzler, D., Ren, B., and Glass, C. K. (2018a). Analysis of Genetically Diverse Macrophages Reveals Local and Domain-wide Mechanisms that Control Transcription Factor Binding and Function. *Cell*, 173(7):1796–1809.e17.
- Link, V. M., Romanoski, C. E., Metzler, D., and Glass, C. K. (2018b). MMARGE: Motif mutation analysis for regulatory genomic elements. *Nucleic Acids Research*, 46(14):7006–7021.
- MacArthur, J., Bowler, E., Cerezo, M., Gil, L., Hall, P., Hastings, E., Junkins, H., McMahon, A., Milano, A., Morales, J., et al. (2017). The new nhgri-ebi catalog of published genome-wide association studies (gwas catalog). *Nucleic acids research*, 45(D1):D896–D901.
- Menoret, D., Santolini, M., Fernandes, I., Spokony, R., Zanet, J., Gonzalez, I., Latapie, Y., Ferrer, P., Rouault, H., White, K. P., et al. (2013). Genome-wide analyses of shaven baby target genes reveals distinct features of enhancer organization. *Genome biology*, 14(8):R86.
- Mills, R. E., Pittard, W. S., Mullaney, J. M., Farooq, U., Creasy, T. H., Mahurkar, A. A., Kemeza, D. M., Strassler, D. S., Ponting, C. P., Webber, C., and Devine, S. E. (2011). Natural genetic variation caused by small insertions and deletions in the human genome. *Genome research*, 21(6), 830–839.
- Morgunova, E. and Taipale, J. (2017). Structural perspective of cooperative transcription factor binding. *Current Opinion in Structural Biology*, 47:1–8.
- Nandi, S., Blais, A., and Ioshikhes, I. (2013). Identification of cis-regulatory modules in promoters of human genes exploiting mutual positioning of transcription factors. *Nucleic Acids Research*, 41(19):8822–8841.
- Ng, F. S., Schütte, J., Ruau, D., Diamanti, E., Hannah, R., Kinston, S. J., and Göttgens, B. (2014). Constrained transcription factor spacing is prevalent and important for transcriptional control of mouse blood cells. *Nucleic Acids Research*, 42(22):13513–13524.10
- Panne, D. (2008). The enhanceosome. *Current Opinion in Structural Biology*, 18(2):236–242.
- Reiter, F., Wienerroither, S., and Stark, A. (2017). Combinatorial function of transcription factors and cofactors. *Current Opinion in Genetics & Development*, 43:73–81.
- Reményi, A., Lins, K., Nissen, L. J., Reinbold, R., Schöler, H. R., and Wilmanns, M. (2003). Crystal structure of a POU/HMG/DNA ternary complex suggests differential assembly of Oct4 and Sox2 on two enhancers. *Genes and Development*, 17(16):2048–2059.
- Rodda, D. J., Chew, J. L., Lim, L. H., Loh, Y. H., Wang, B., Ng, H. H., and Robson, P. (2005). Transcriptional regulation of Nanog by OCT4 and SOX2. *Journal of Biological Chemistry*, 280(26):24731–24737.
- Shen, Z., Hoeksema, M., Ouyang, Z., Benner, C., and Glass, C. (2020). MAGGIE: leveraging genetic variation to identify DNA sequence motifs mediating transcription factor binding and function. *bioRxiv*.
- Slattery, M., Zhou, T., Yang, L., Machado, A. C. D., Gordán, R., and Rohs, R. (2014). Absence of a simple code: how transcription factors read the genome. *Trends in biochemical sciences*, 39(9):381–399.
- Smith, R. P., Taher, L., Patwardhan, R. P., Kim, M. J., Inoue, F., Shendure, J., Ovcharenko, I., and Ahituv, N. (2013). Massively parallel decoding of mammalian regulatory sequences supports a flexible organizational model. *Nature genetics*, 45(9):1021.
- Stormo, G. D. (2000). DNA binding sites: representation and discovery. *Bioinformatics*, 16(1):16–23.
- Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S. J., Brett, M., Wilson, J., Millman, K. J., Mayorov, N., Nelson, A.

726 R. J., Jones, E., Kern, R., Larson, E., Carey, C., Polat, , Feng, Y., Moore, E. W., VanderPlas, J., Laxalde, D.,
727 Perktold, J., Cimrman, R., Henriksen, I., Quintero, E. A., Harris, C. R., Archibald, A. M., Ribeiro, A. H.,
728 Pedregosa, F., van Mulbregt, P., and Contributors, S. . . (2019). SciPy 1.0—Fundamental Algorithms for
729 Scientific Computing in Python. pages 1–22.

730 Sullivan, G. M., and Feinn, R. (2012). Using effect size—or why the P value is not enough. *Journal of graduate*
731 *medical education*, 4(3), 279–282.

732 Visscher, P. M., Wray, N. R., Zhang, Q., Sklar, P., McCarthy, M. I., Brown, M. A., and Yang, J. (2017). 10
733 years of gwas discovery: biology, function, and translation. *The American Journal of Human Genetics*,
734 101(1):5–22.

Figure 1

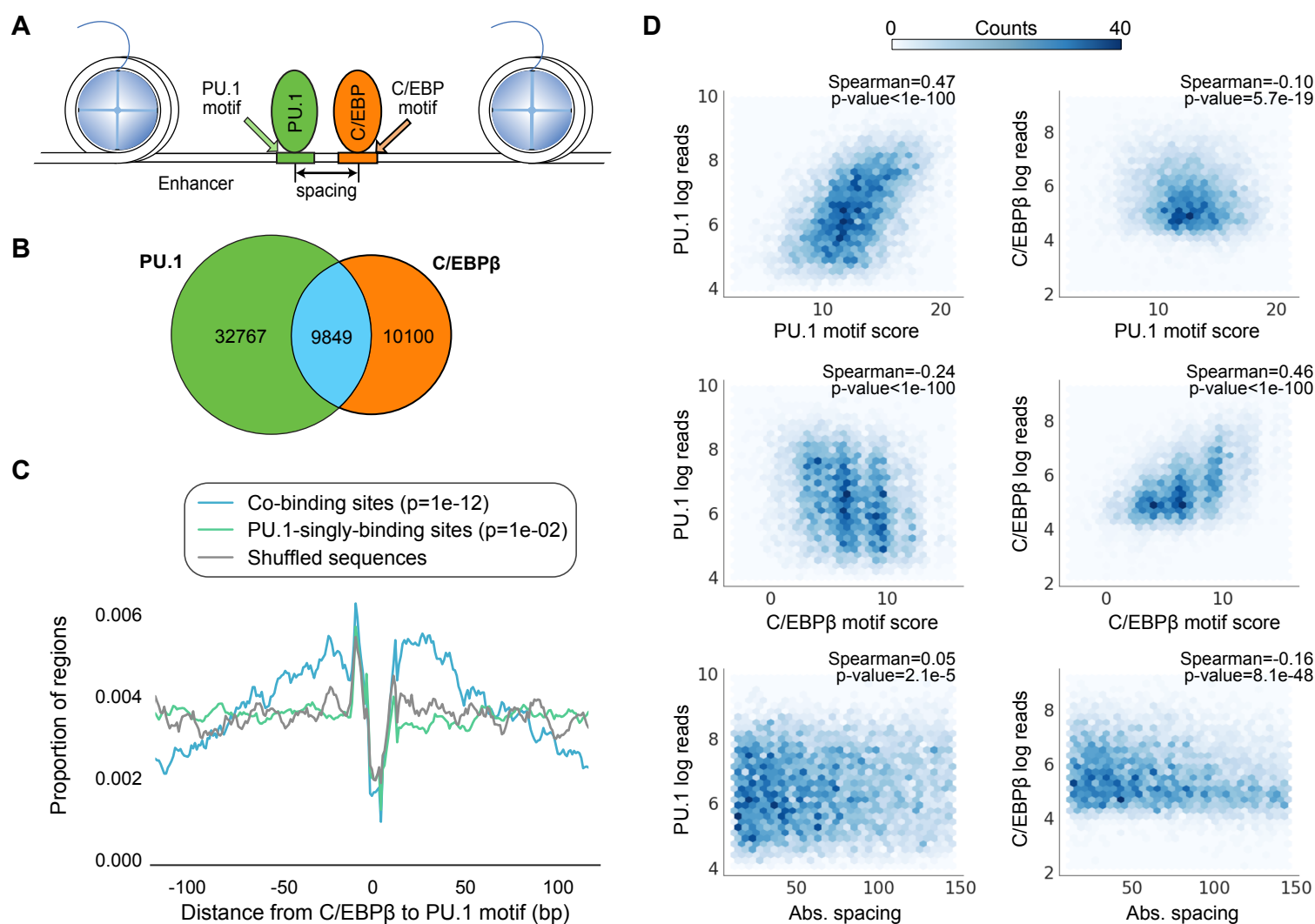


Figure 2

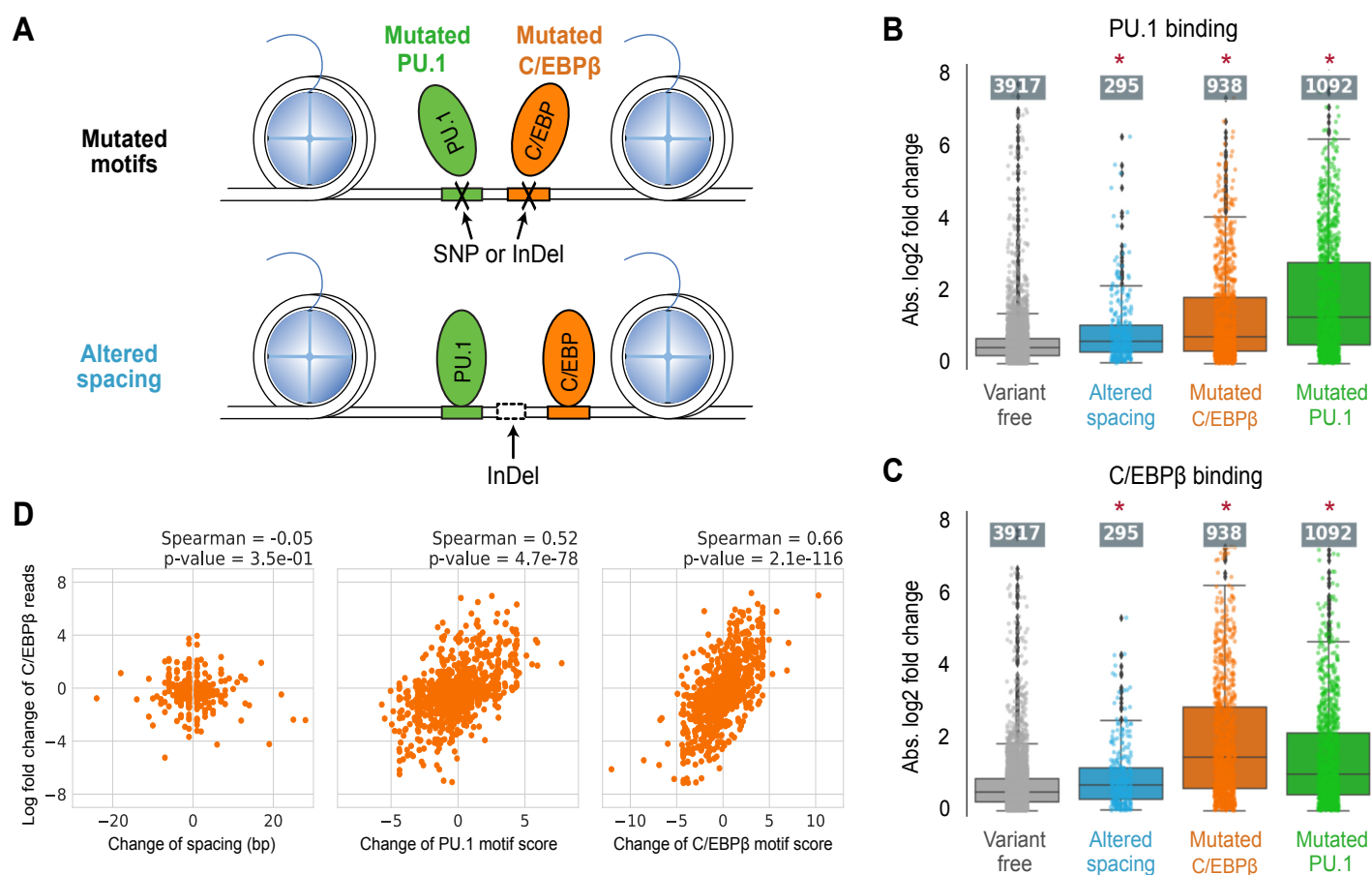


Figure 3

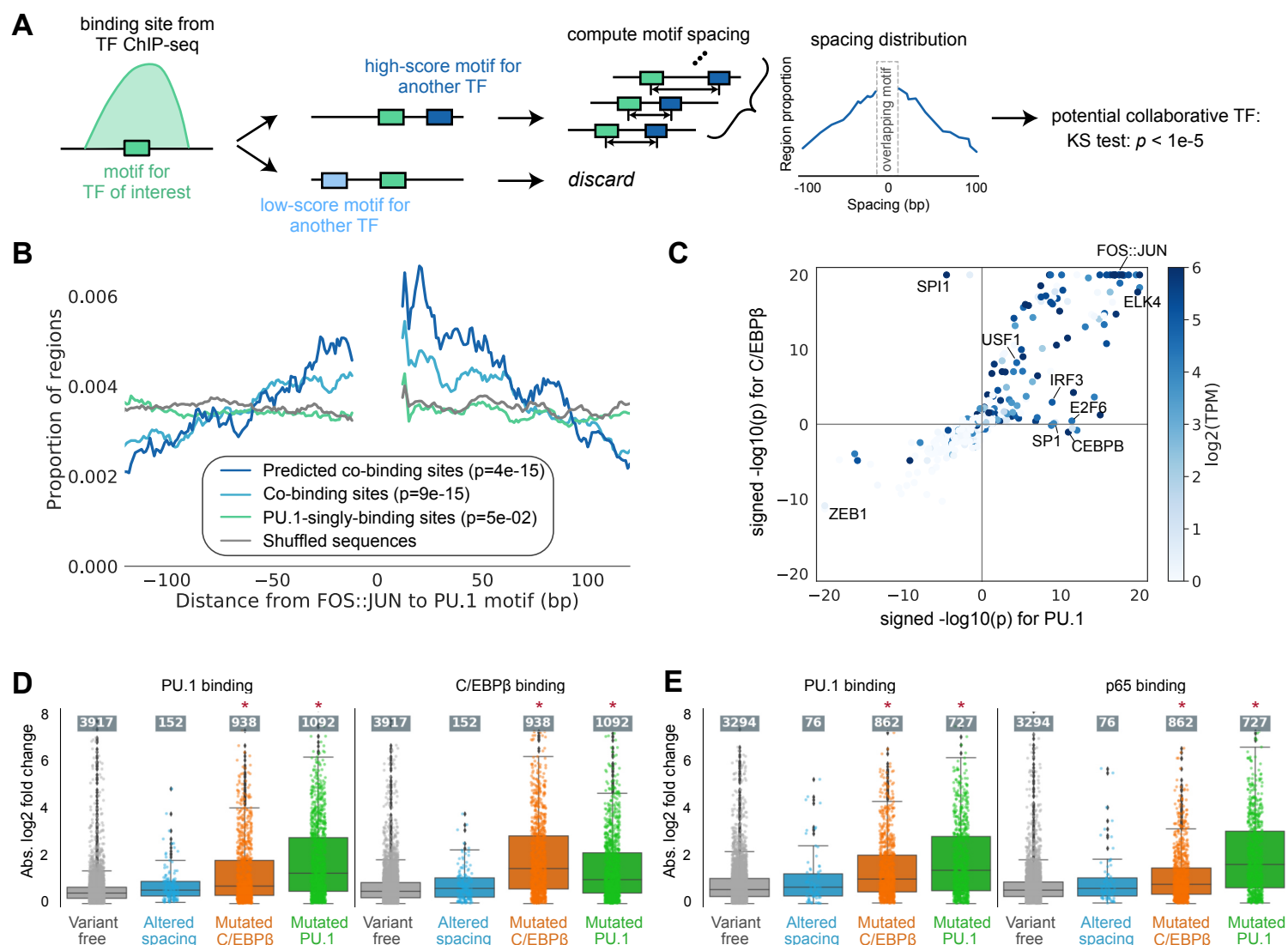


Figure 4

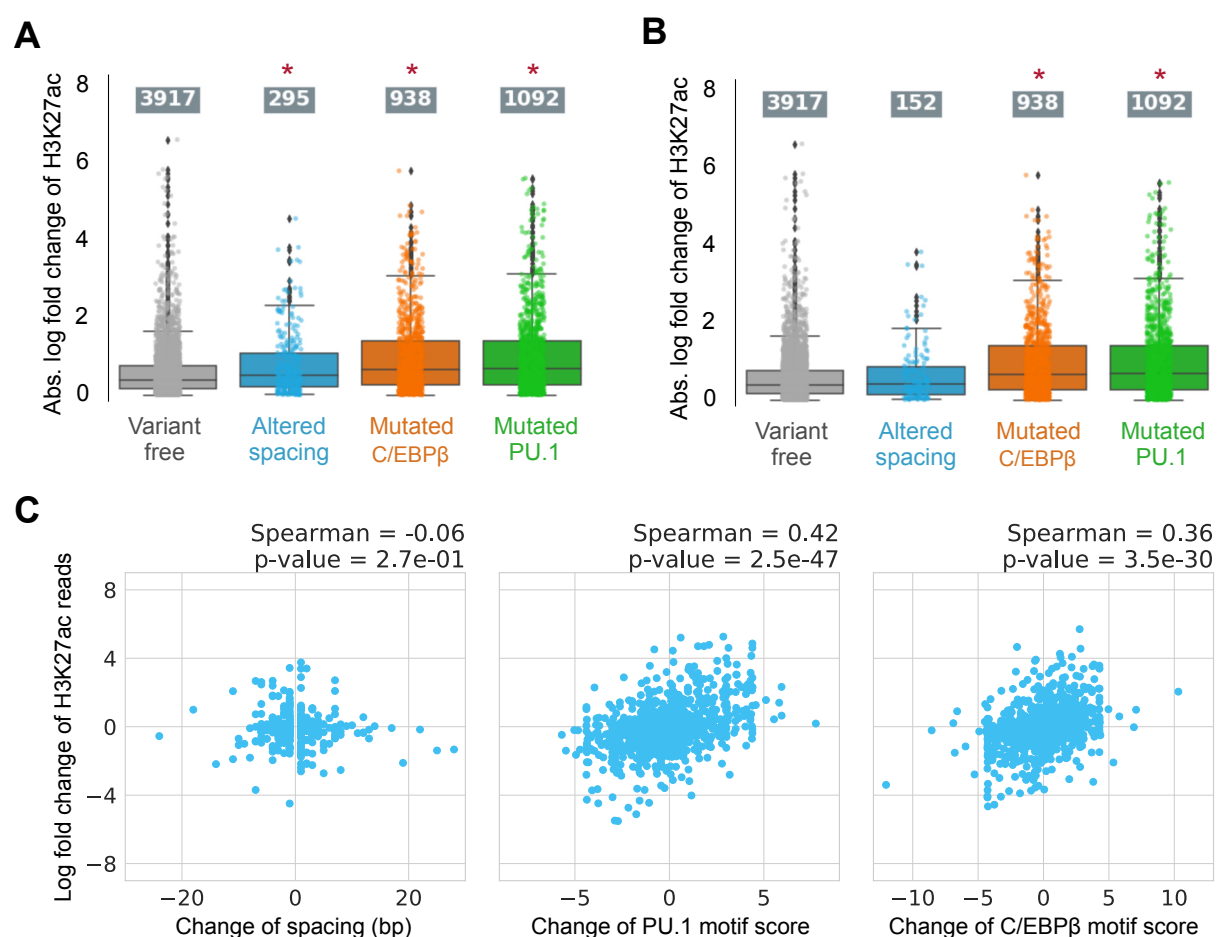


Figure 5

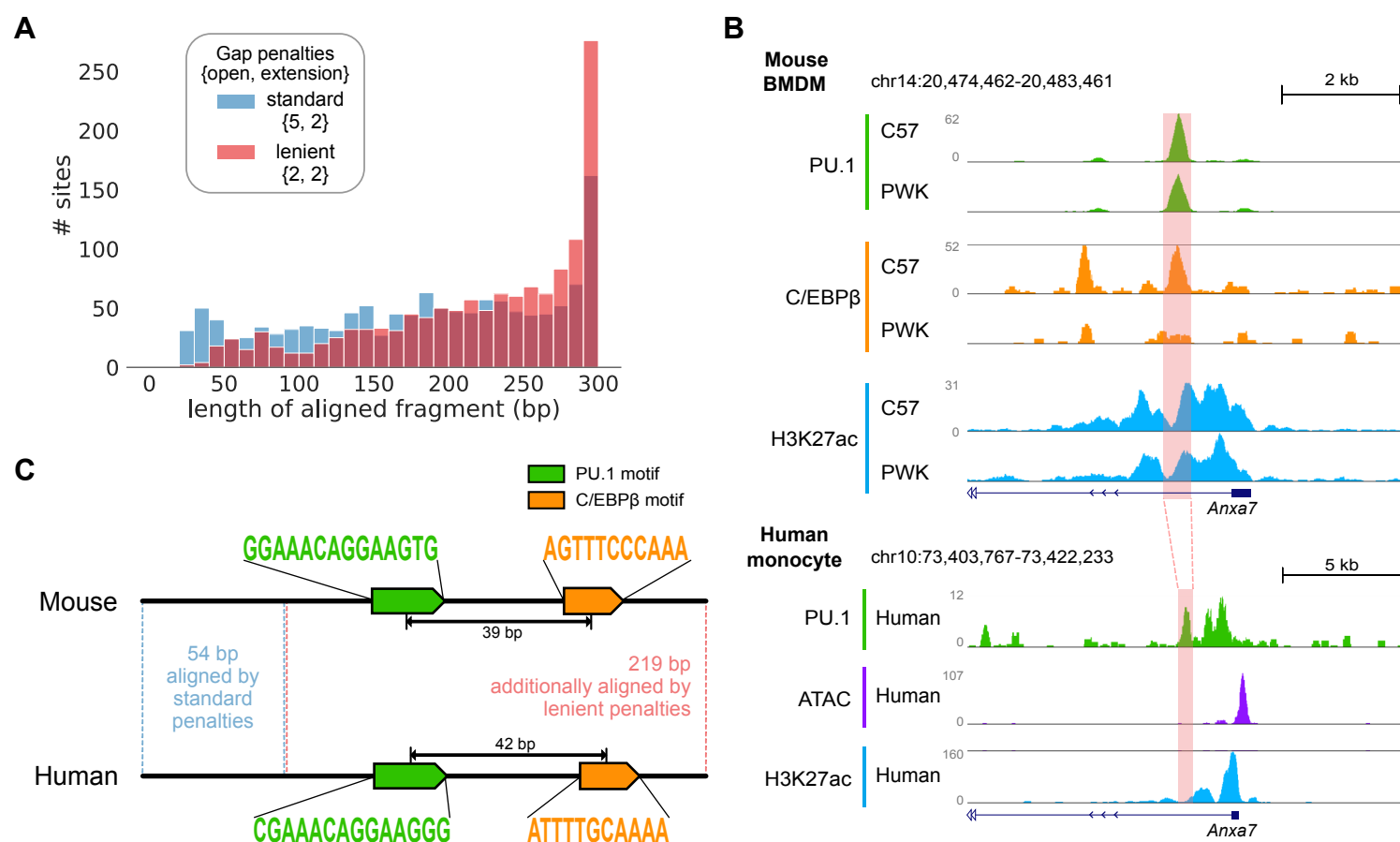


Figure 1—figure supplement 1

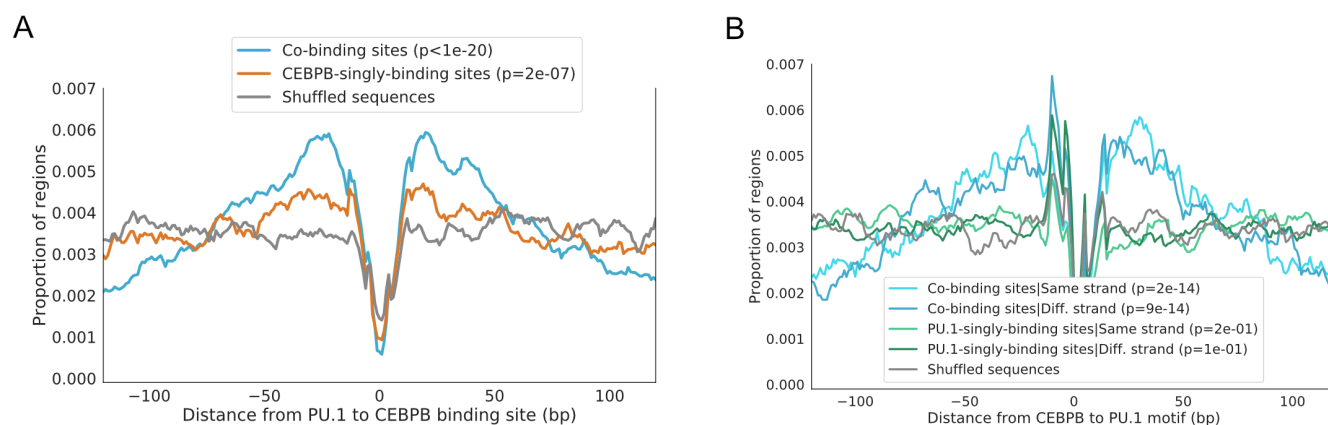


Figure 2—figure supplement 1

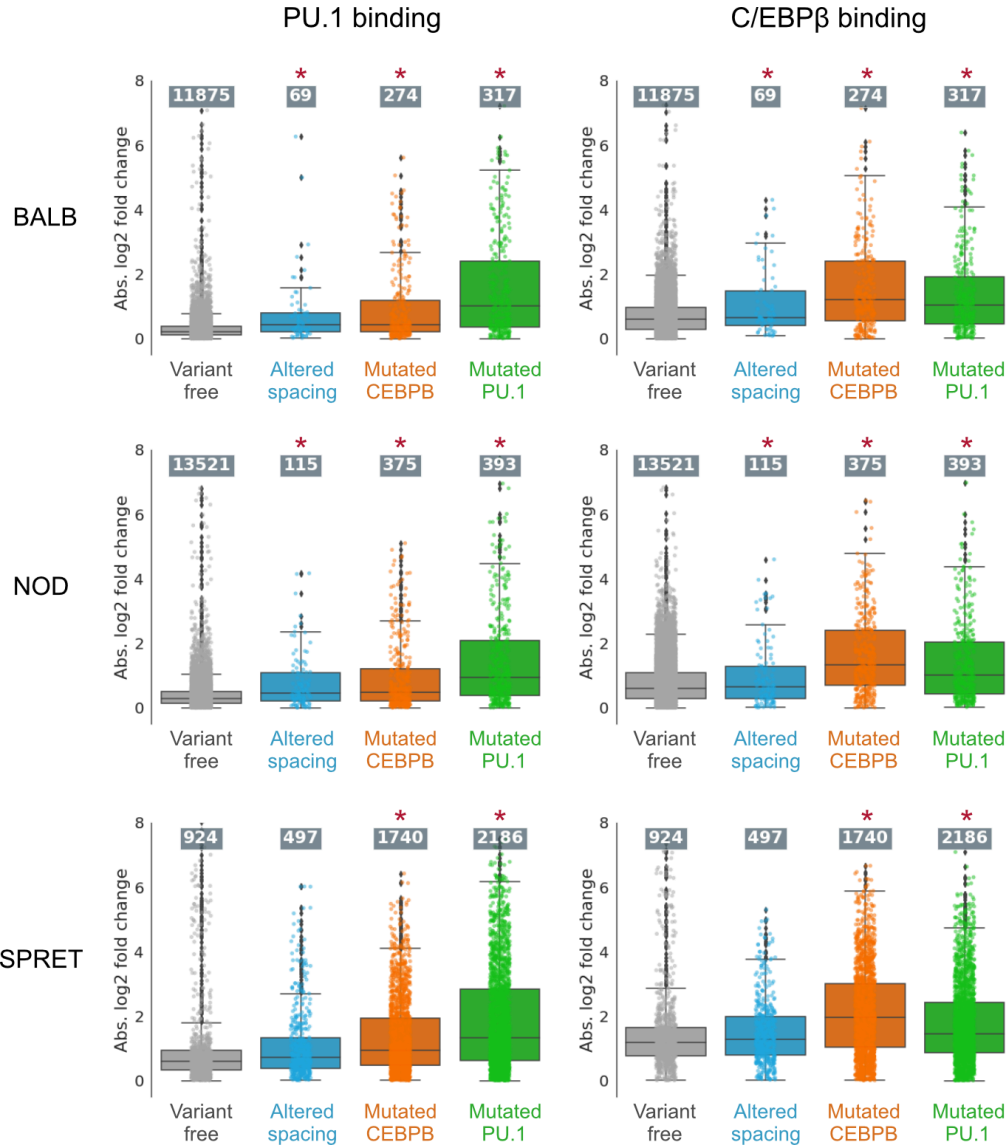


Figure 2—figure supplement 2

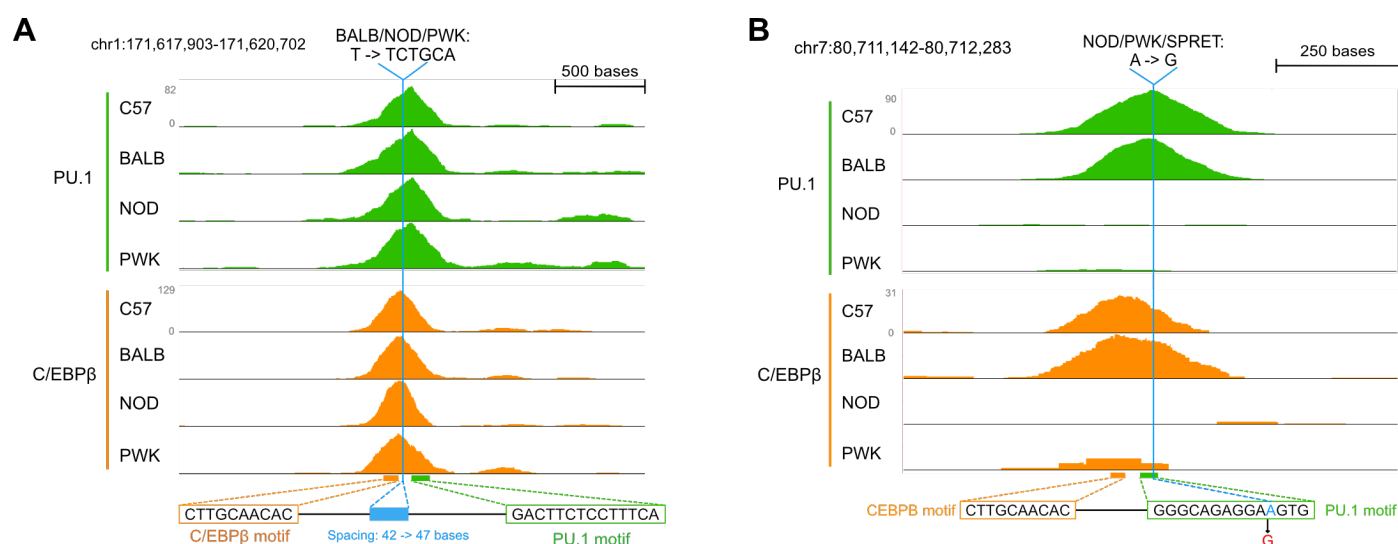


Figure 2—figure supplement 3

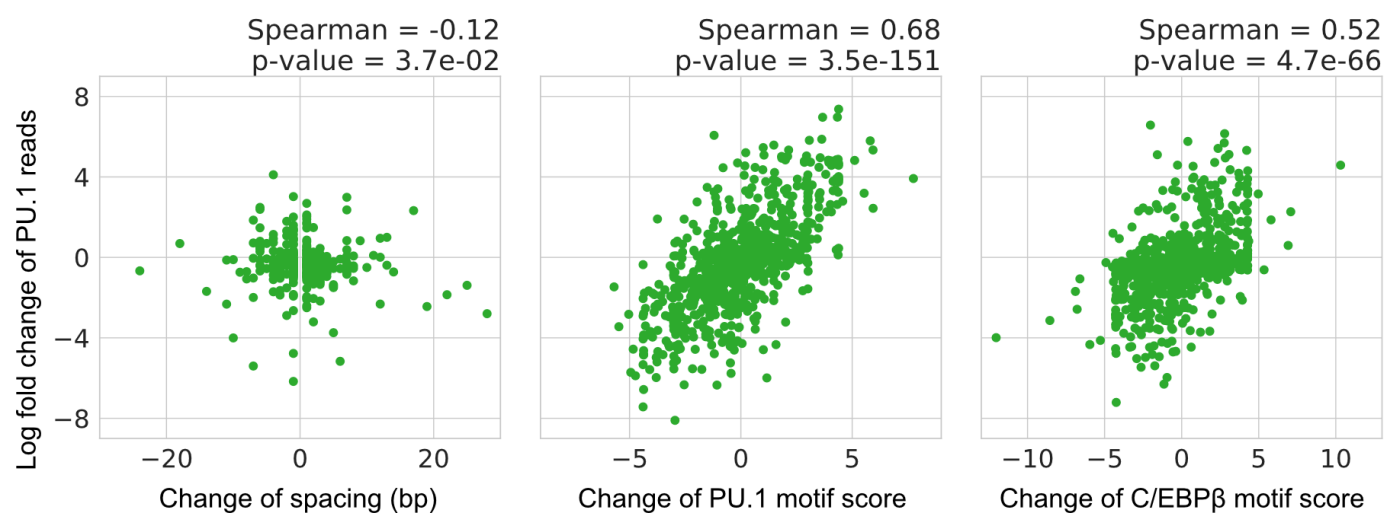


Figure 2—figure supplement 4

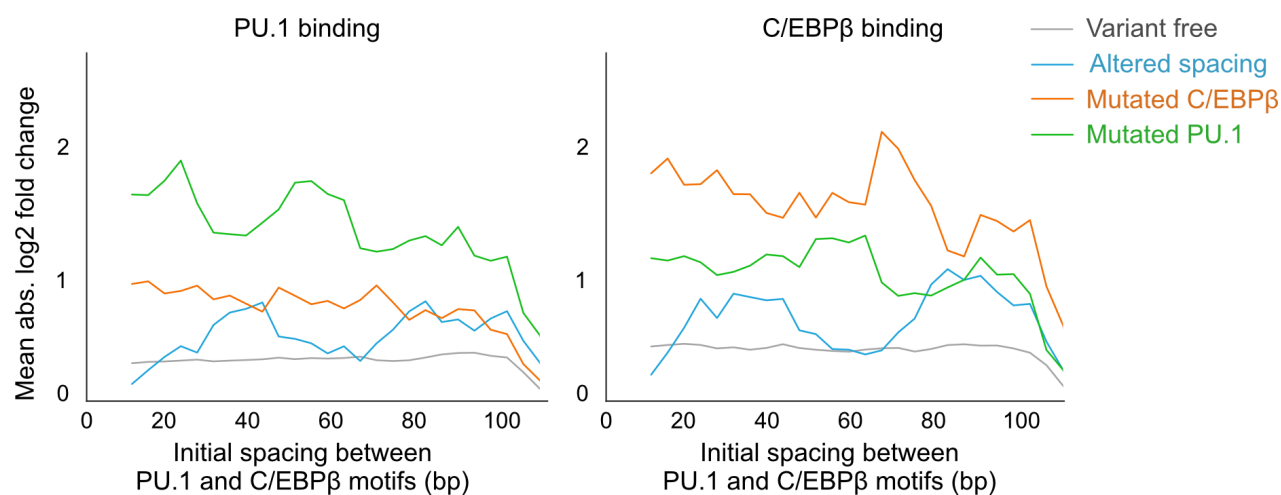


Figure 3—figure supplement 1

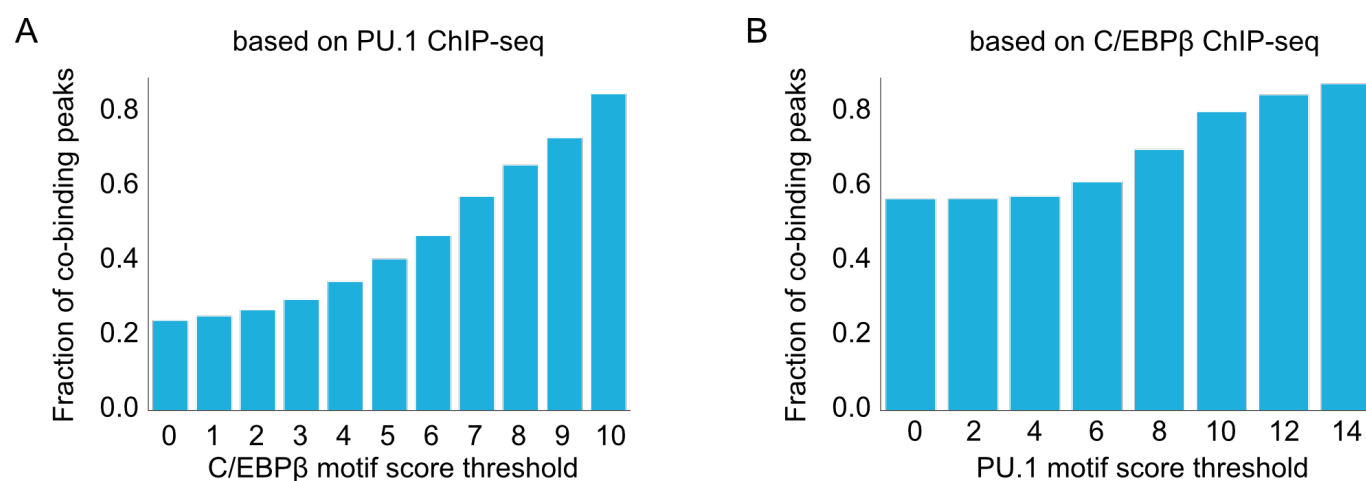


Figure 3—figure supplement 2

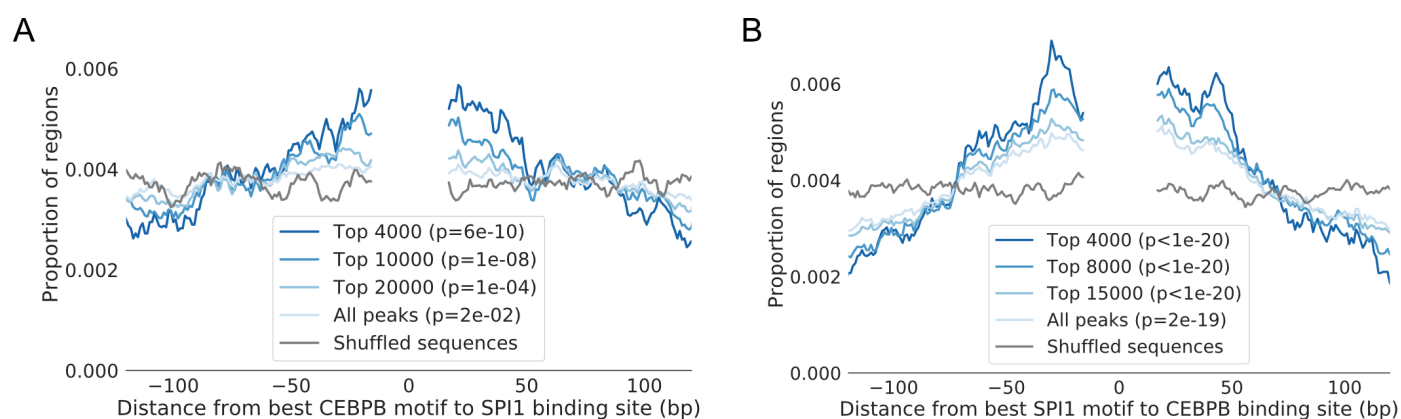


Figure 3—figure supplement 3

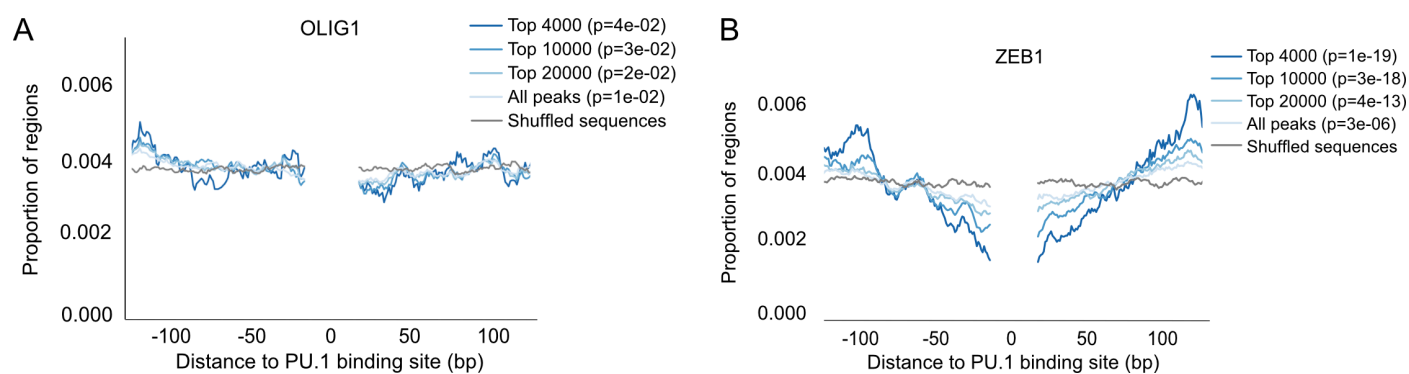


Figure 3—figure supplement 4

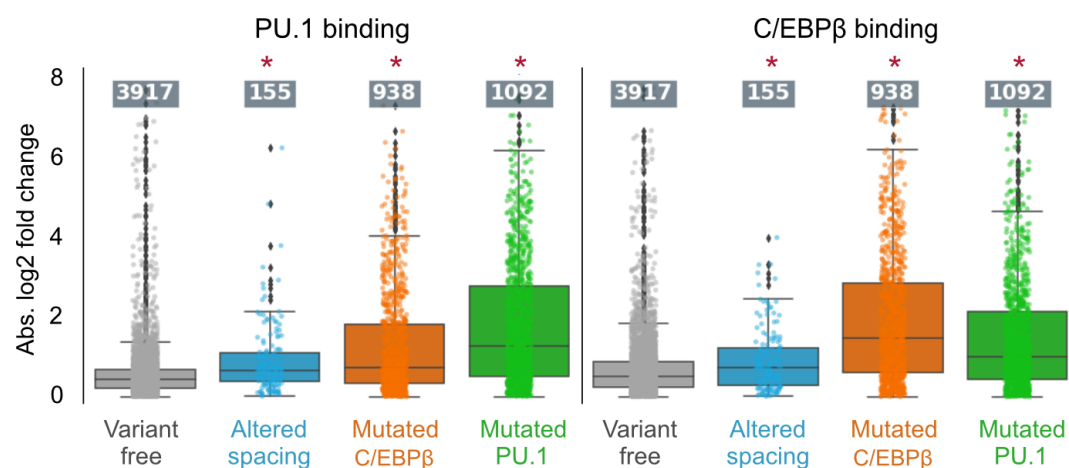


Figure 3—figure supplement 5

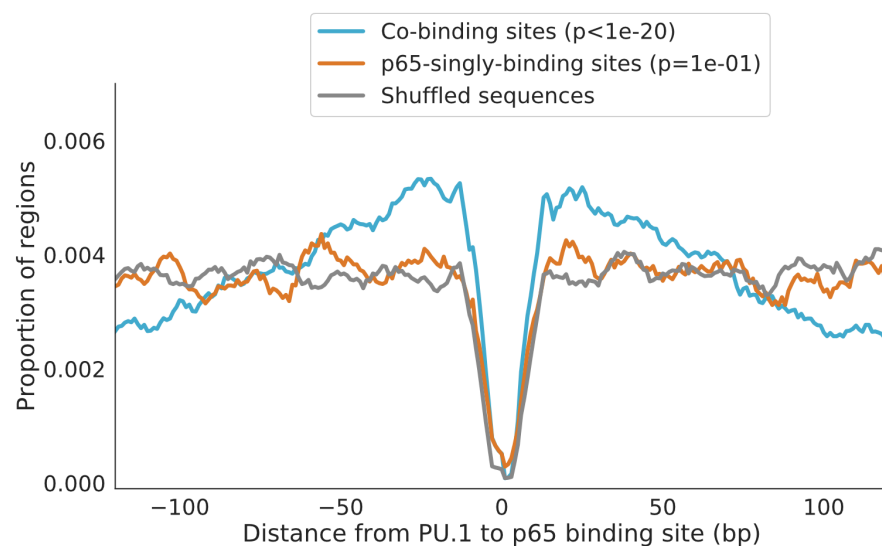


Figure 3—figure supplement 6

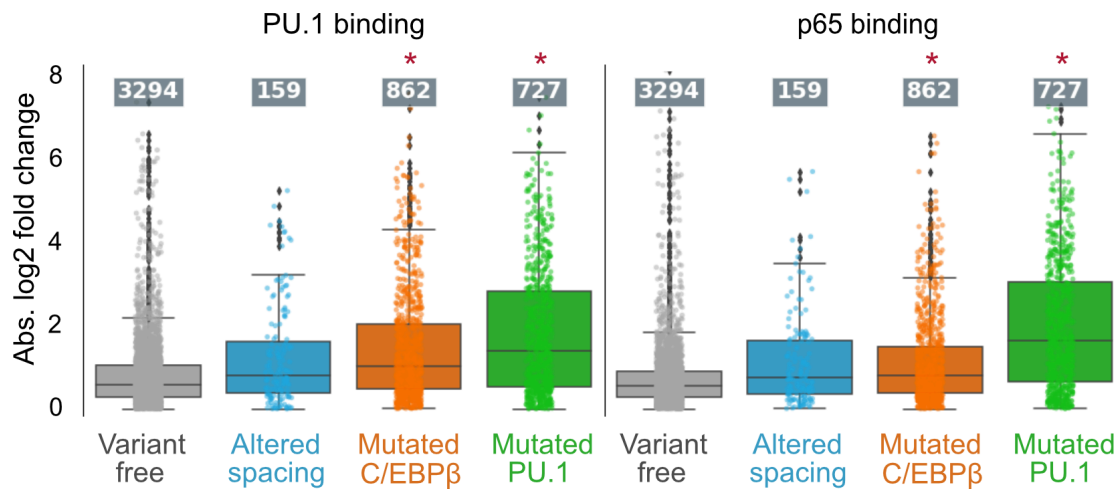


Figure 4—figure supplement 1

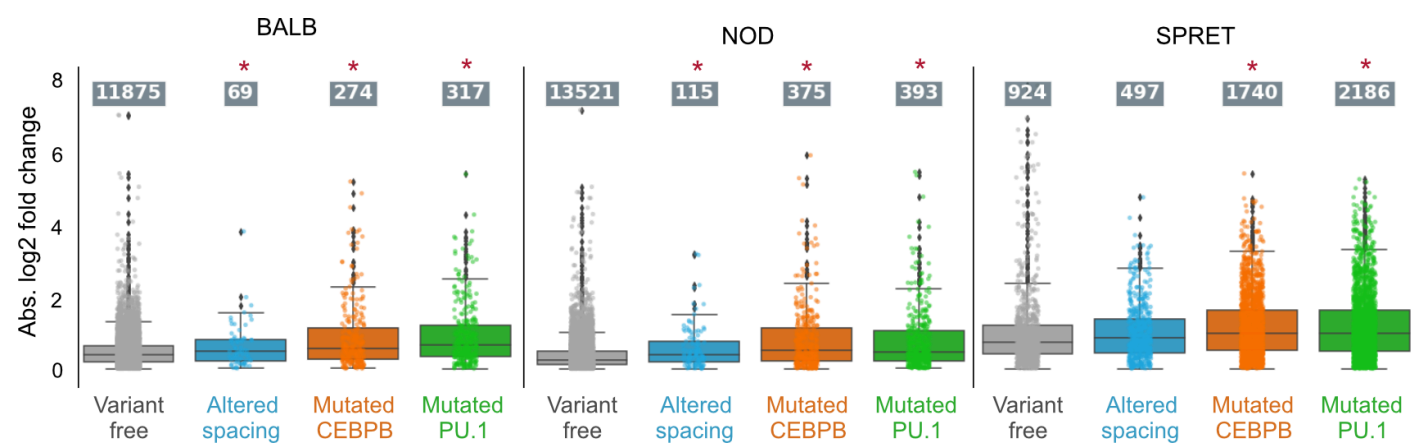


Figure 4—figure supplement 2

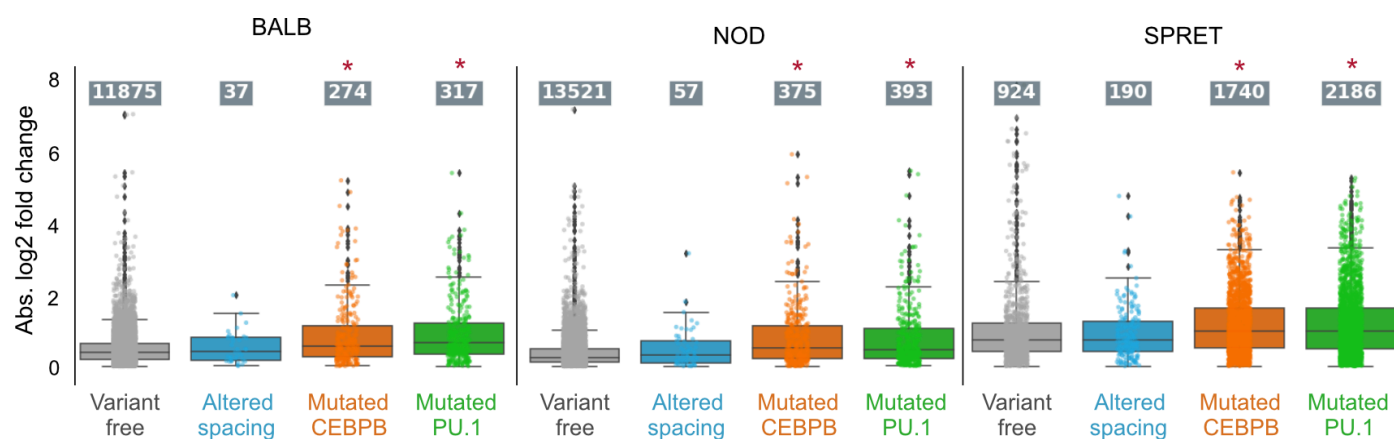


Figure 4—figure supplement 3

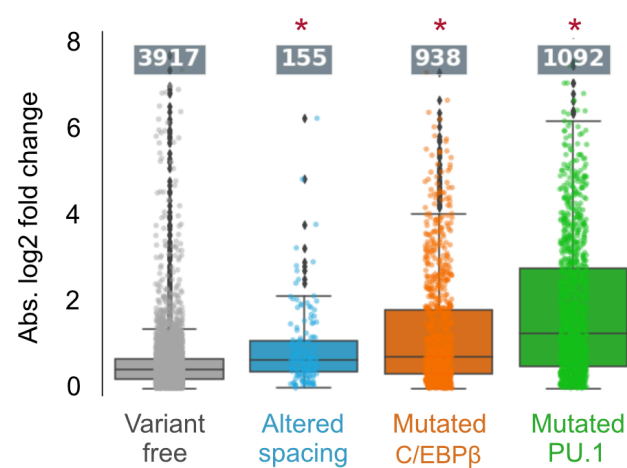


Figure 4—figure supplement 4

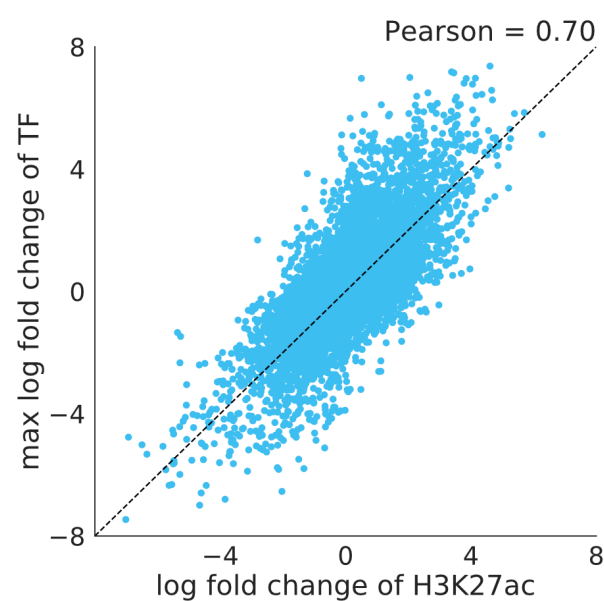


Figure 5—figure supplement 1

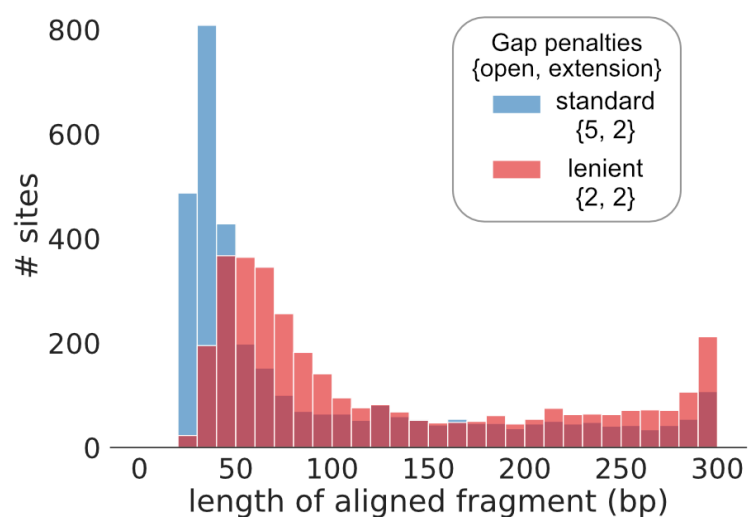


Figure 5—figure supplement 2

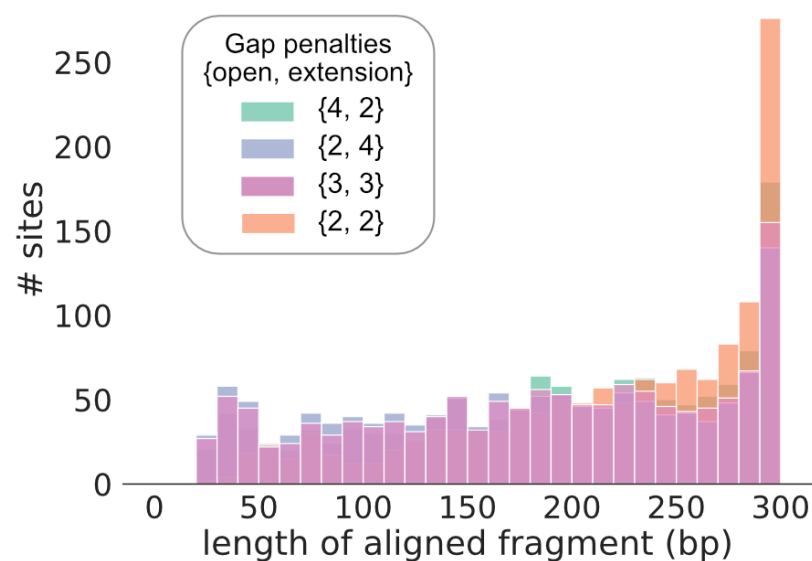


Figure 5—figure supplement 3



Figure 5—figure supplement 4

