

Supporting Material

Genome-wide organization of eukaryotic pre-initiation complex is influenced by nonconsensus protein-DNA binding

Ariel Afek and David B. Lukatsky*

Department of Chemistry, Ben-Gurion University of the Negev, Beer-Sheva 84015, Israel

**Corresponding author:*
Email: lukatsky@bgu.ac.il

Supporting Figure Legends

Figure S1. This figure demonstrates the robustness of the computed free energy of nonconsensus TF-DNA binding with respect to the global variability of the nucleotide content along the yeast genome (**A**); the robustness with respect to the width of the sliding window, L (**B**); and the robustness with respect to the number of contacts, M , that the TF makes with DNA (**C**). (**A**) The average free energy of nonconsensus TF-DNA binding per bp, $\langle f \rangle = \langle \langle F \rangle_{\text{TF}} \rangle_{\text{seq}} / M$, (red), as compared with the corresponding normalized free energy per bp, $\langle \delta f \rangle = \langle \langle \delta F \rangle_{\text{TF}} \rangle_{\text{seq}} / M$, (blue), where $\delta F = F - F_{\text{rand}}$. For a given TF, F is computed as described in the main text, and F_{rand} is the free energy computed for a randomized sequence (in the same sliding window as F), and averaged over 25 random realizations. The described procedure removes the bias in the free energy, stemming from the global variability of the nucleotide content. (**B**) The normalized, average free energy, $\langle \delta f \rangle$, computed using different values of the width of the sliding window, $L = 30$ (red), $L = 50$ (black), and $L = 80$ (blue). We used $L = 50$ for all the calculations described in the main text. (**C**) The normalized, average free energy, $\langle \delta f \rangle$, computed using different values of the TF size, $M = 6$ (red), $M = 8$ (black), and $M = 10$ (blue). We used $M = 8$ for all the calculations described in the main text. In all plots, (**A**), (**B**), and (**C**), the average free energies are computed using 6,045 yeast transcripts.

Figure S2. This figure further demonstrates the statistical significance of the correlation between the free energy of nonconsensus TF-DNA binding and the PIC occupancy in the vicinity of the TSS. (**A**) Correlation between the minimal value of the free energy of nonconsensus TF-DNA binding, $f_{\text{min}} = \min(f)$, where $f = \langle F \rangle_{\text{TF}} / M$, and the combined occupancy of all GTFs, computed for individual genes in non-overlapping windows of 80 bp within the entire interval $(-400, 400)$ around the TSS for each of these 3,945 genes. The data are binned into 50 bins. The notation $\langle \text{PIC} \rangle$ describes the average, combined occupancy profile of all nine GTFs. (**B**) Correlation between the minimal value of the free energy of nonconsensus TF-DNA binding, $f_{\text{min}} = \min(f)$, and the maximal combined occupancy of eight GTFs (all GTFs less the Pol II occupancy), computed for individual genes within the entire interval $(-150, 0)$ around the TSS for each of these 3,945 genes. The data are binned into 25 bins.

Figure S3. (**A**) Correlation between the minimal value of the free energy of nonconsensus TF-DNA binding, $f_{\text{min}} = \min(f)$, where $f = \langle F \rangle_{\text{TF}} / M$, and the average, combined GTF occupancy computed for individual genes in non-overlapping windows of 80 bp within the entire interval $(-2990, 2070)$ around the open reading frame (ORF) ends for 2,903 mRNA genes. The data are binned into 50 bins. (**B**) Similarly computed correlation of f_{min} with the nucleosome occupancy. (**C**) The heat maps represent the free energy of nonconsensus TF-DNA binding, f , the combined occupancy of the GTFs, and the nucleosome occupancy, respectively, for 1,860 tandem mRNA genes aligned with respect to the ORF ends. The genes are sorted by intergenic length.

Figure S4. Correlation between the minimal value of the free energy of nonconsensus TF-DNA binding, $f_{\text{min}} = \min(f)$, where $f = \langle F \rangle_{\text{TF}} / M$, and the average GTF occupancy

computed for individual genes, in non-overlapping windows of 80 bp within the entire interval $(-2990, 2070)$ around the ORF ends for 2,903 mRNA genes. The data are binned into 50 bins.

Figure S5. (A) Correlation between the minimal value of the free energy of nonconsensus TF-DNA binding, $f_{\min} = \min(f)$, where $f = \langle F \rangle_{\text{TF}} / M$, and the average GTF occupancy of 676 TATA-containing genes. The correlation is computed for individual genes in non-overlapping windows of 80 bp within the entire interval $(-990, 990)$ around the TSS. The data are then binned into 50 bins. (B) Similar to (A), but now f_{\min} is correlated with the nucleosome occupancy of these TATA-containing genes. (C) Correlation between f_{\min} , and the average GTF occupancy of 3,269 TATA-less genes. The correlation is computed for individual genes in non-overlapping windows of 80 bp within the entire interval $(-990, 990)$ around the TSS. The data are binned into 50 bins. (D) Similar to (C), but now f_{\min} is correlated with the nucleosome occupancy of these TATA-less genes. (E) The average free energy of nonconsensus TF-DNA binding per bp, $\langle f \rangle$, for a larger set of 5,034 TATA-less genes (blue), and 1,011 TATA-containing genes (red), around the TSSs.

Figure S6. This figure is supplementary to **Figure 7** of the main text. It demonstrates the robustness of our conclusions with respect to an alternative definition of the TATA-like box occupancy score. (A) The TATA-like box occupancy scores were computed based on the PWM taken from V.X. Jin et al., BMC Bioinformatics 7:114 (2006). We used 1,432 genes with double-peak in the nucleosome occupancy (red) and 2,513 genes with single-peak in the nucleosome occupancy (blue), as described in **Figure 7** of the main text. The obtained TATA-like box occupancy profiles are similar to those presented in **Figure 7** of the main text, and these profiles are statistically indistinguishable between those two groups of genes. (B) We also computed the initiator (INR) element PWM occupancy score for the selected two groups of genes, as described in (A). The obtained occupancy profiles are statistically indistinguishable between these two groups of genes.

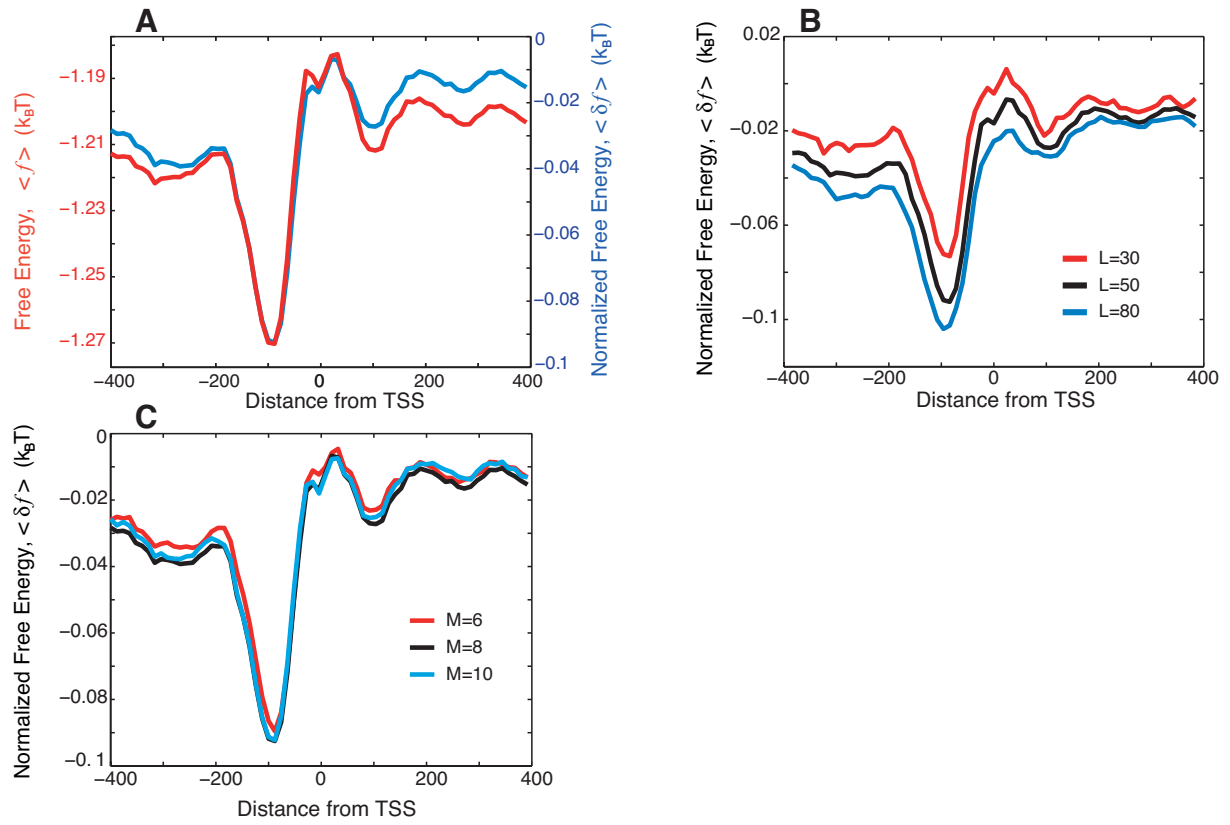


Figure S1

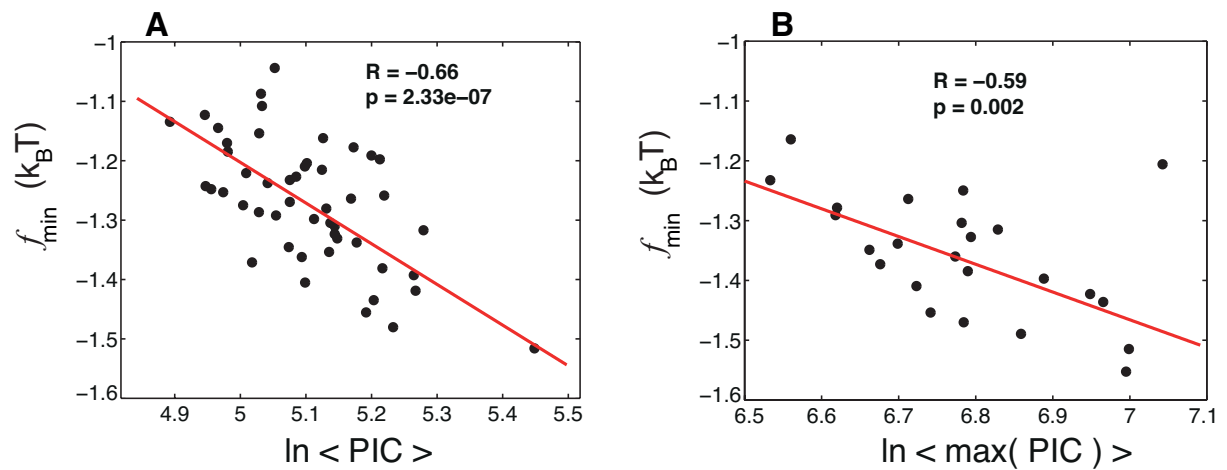


Figure S2

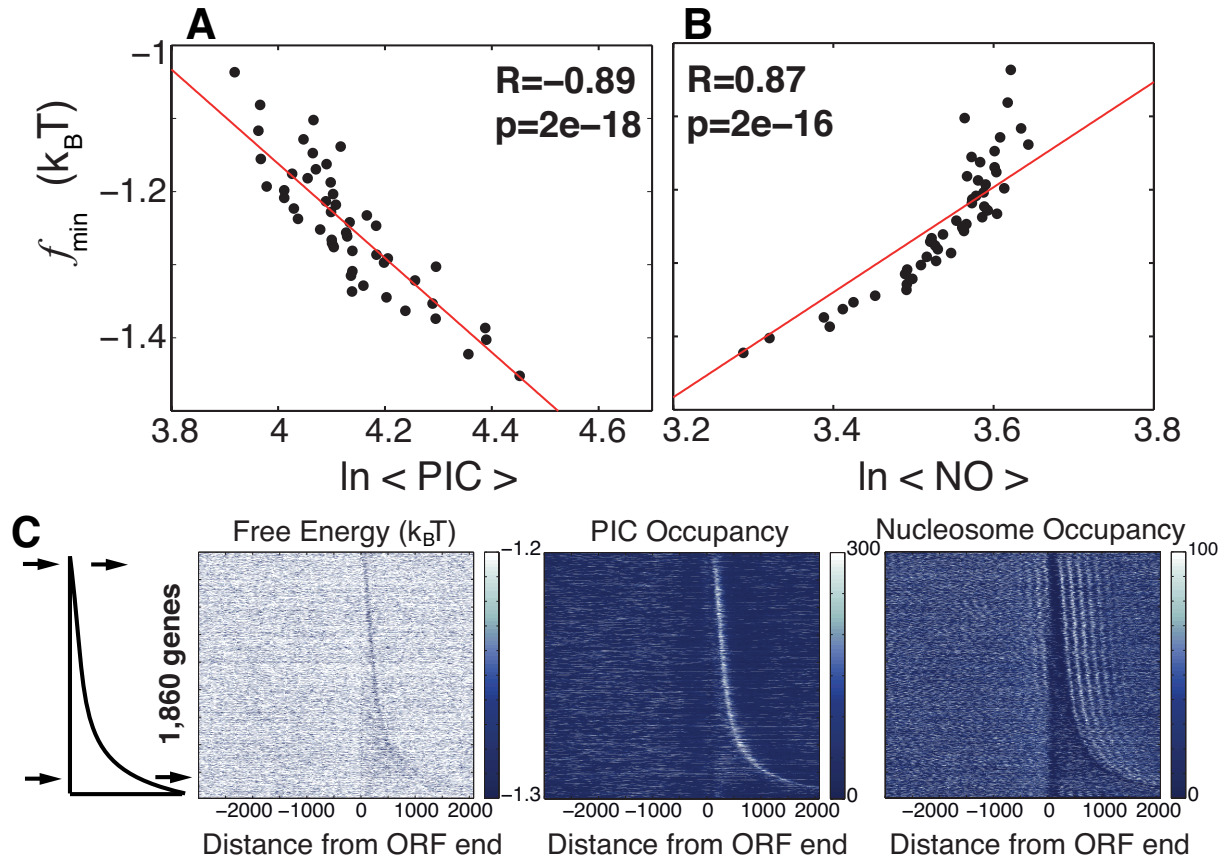


Figure S3

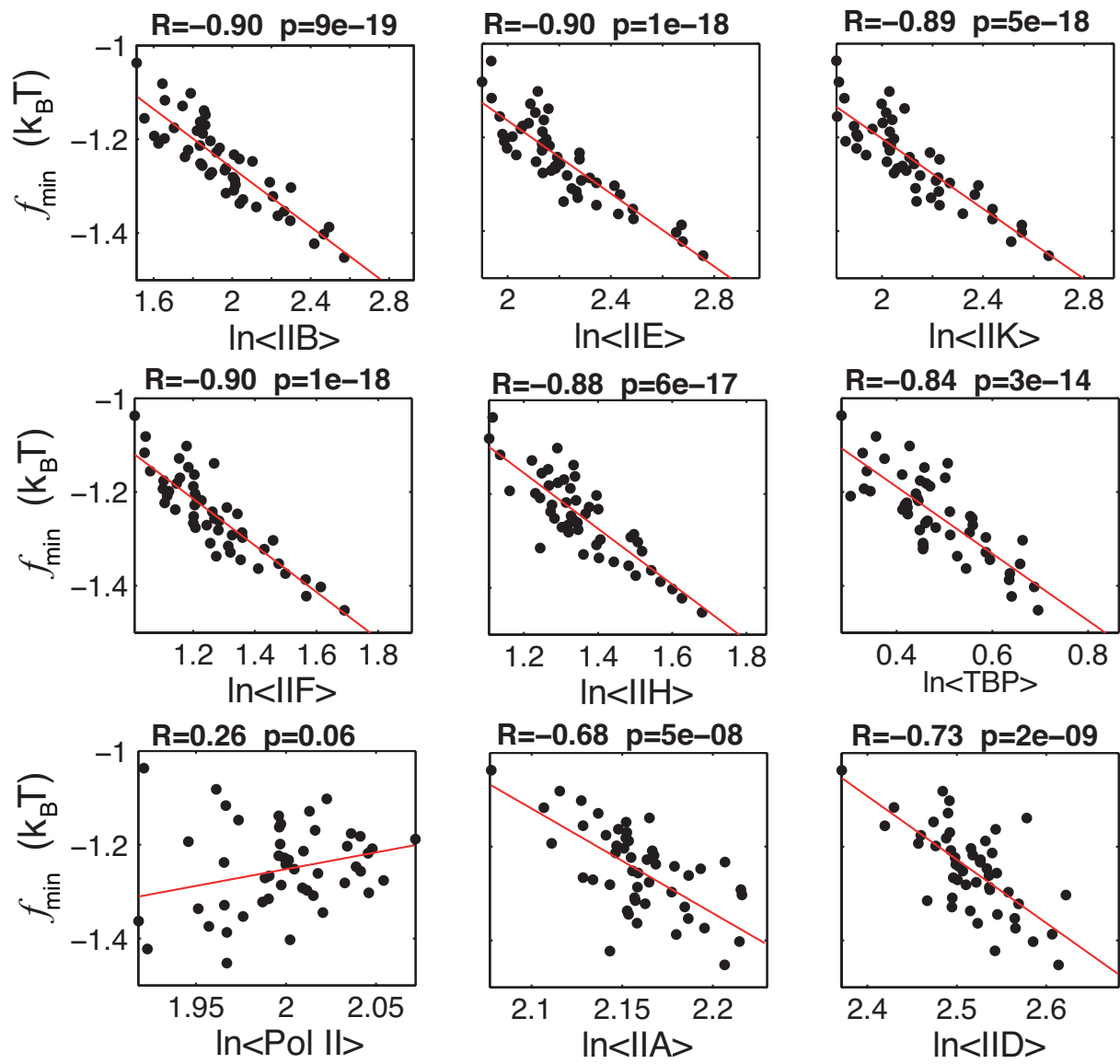


Figure S4

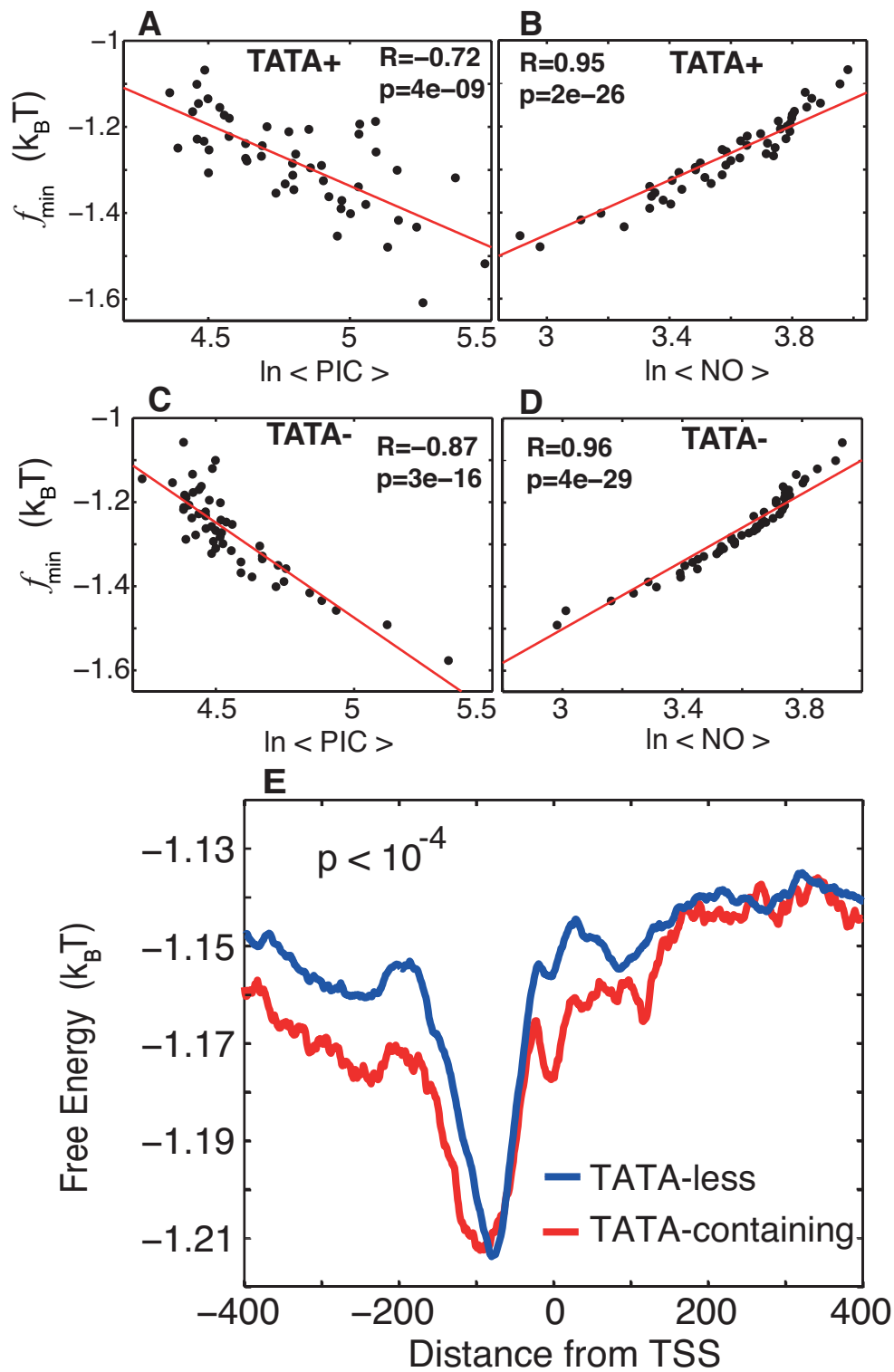


Figure S5

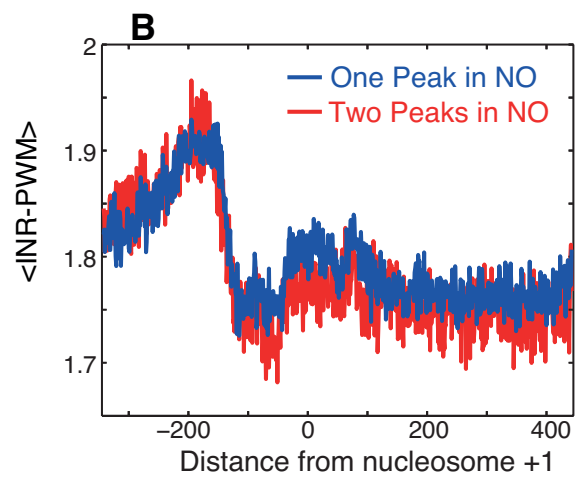
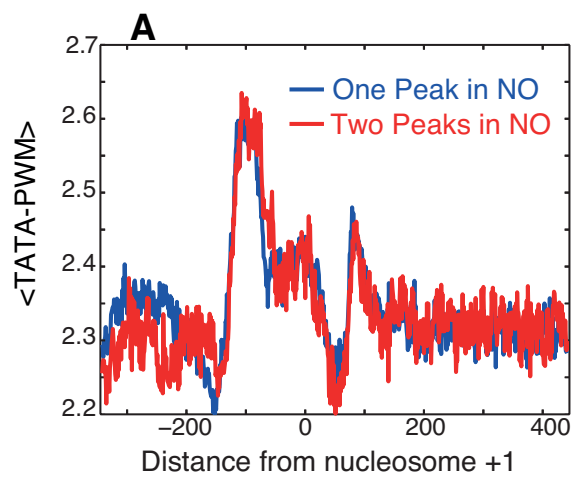


Figure S6