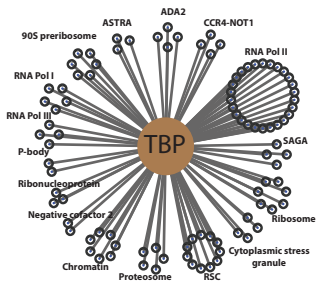


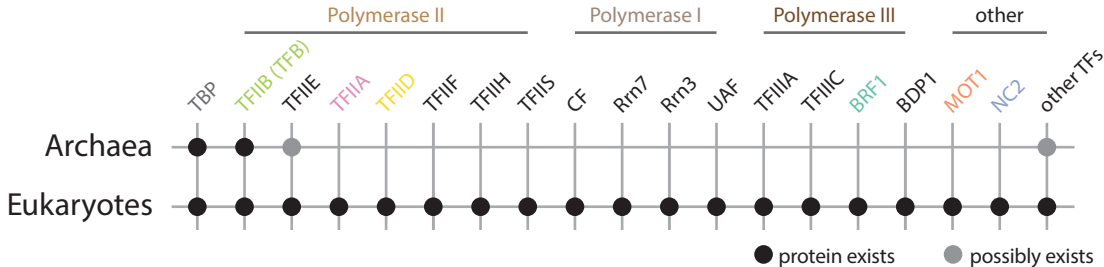
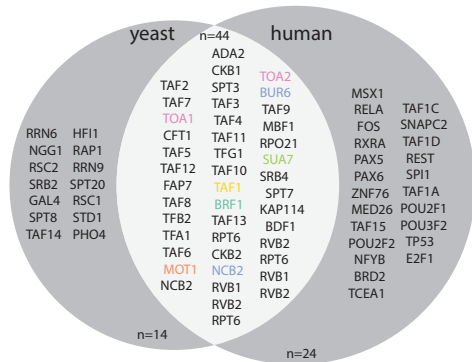
Supplementary Information

**Molecular determinants underlying
functional innovations of TBP and their
impact on transcription initiation**

Ravarani et al.

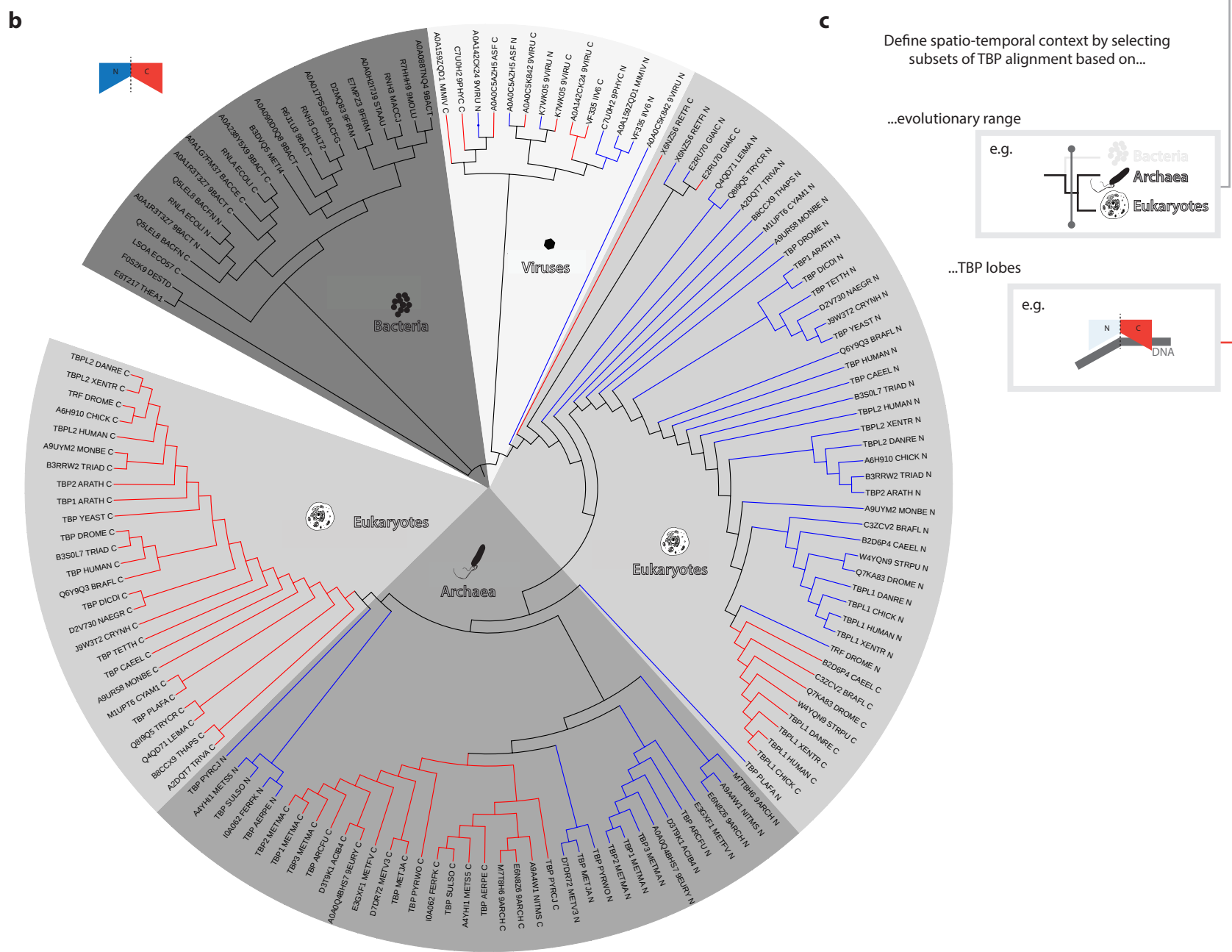
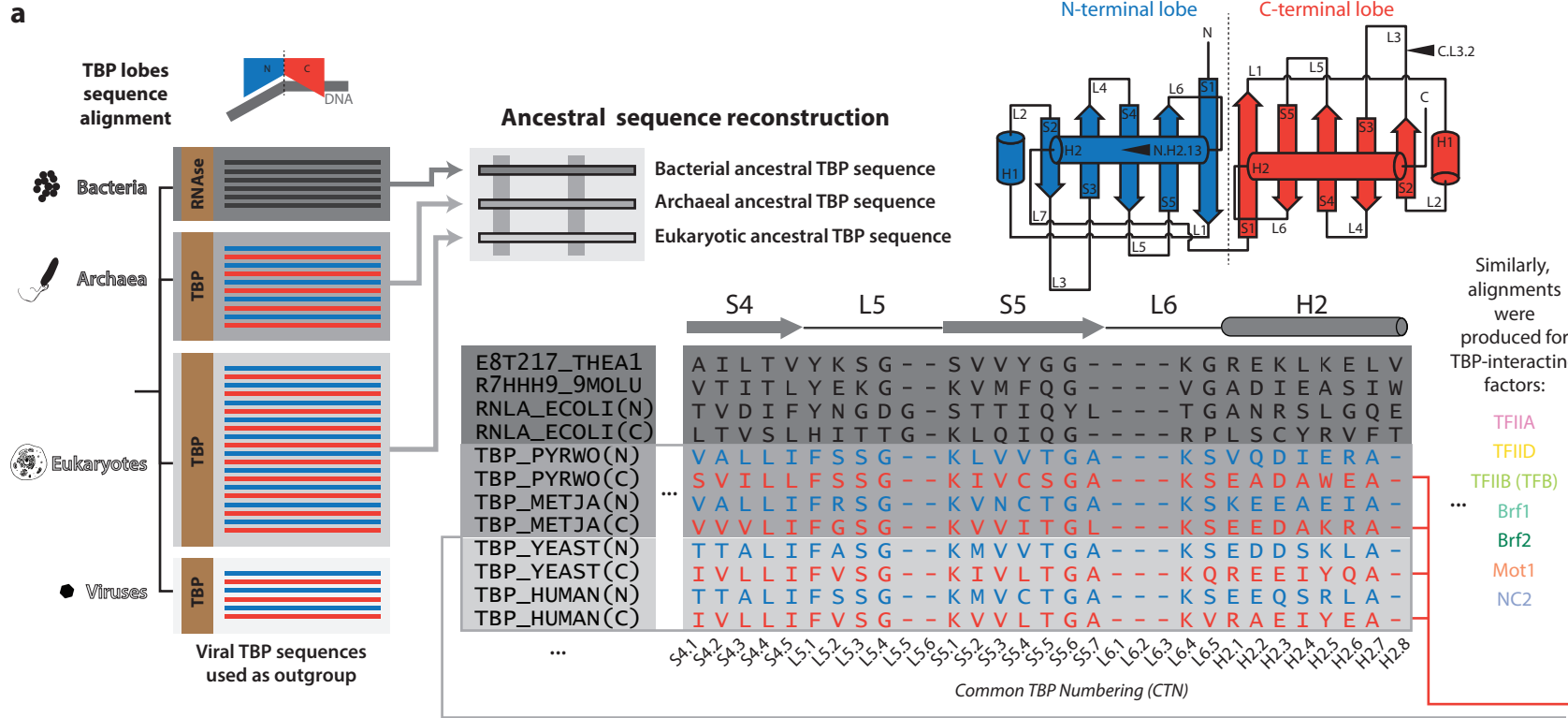


Conservation of interaction partners across evolution



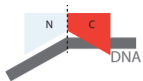
Supplementary Figure 1. Physical interaction partners of TBP and their conservation between organisms

- a) TBP functions as a hub of physical interactions for various transcriptional protein complexes. Some of these protein complexes with TBP include several chromatin-remodeling complexes such as ASTRA, ADA2 and SAGA complexes, basal transcription apparatus such as RNA Pol I, II and III holoenzymes and RNA processing machineries such as CCR4-NOT1.
- b) Numerous components of such transcriptional complexes are conserved between organisms and also physically interact with TBP. For instance, about 65% of interaction partners are conserved and physically interact with their respective TBPs between yeast and human TBPs.
- c) Basal transcription apparatus and its components that are common between two superkingdoms, i.e. archaea and eukaryotes, are shown. The presence of numerous additional components in eukaryotic basal transcription machinery is possibly due to the existence of parallel transcription systems, i.e. PolII, PolIII and PolIII and complex chromatin dynamics during the transcription process.

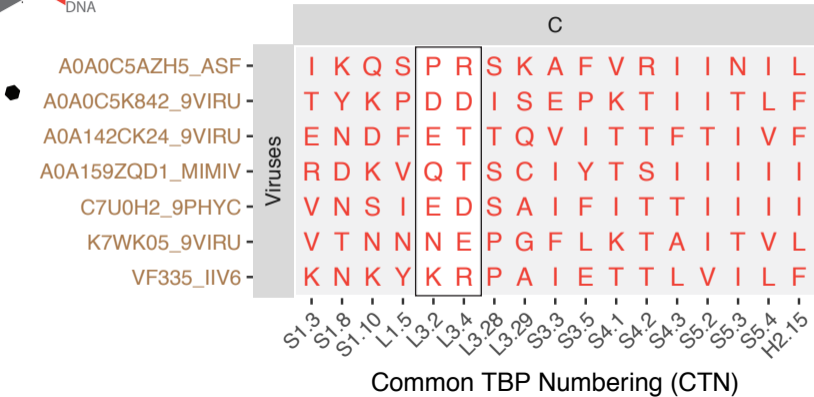


Supplementary Figure 2. Ancestral sequence reconstruction, common TBP-lobe numbering system and dendrogram of TBP-lobe like regions

- a) Alignment of TBP-lobe like regions in RefMSA was used to construct ancestral sequences for each of the superkingdoms i.e. bacteria, archaea and eukaryotes (see **Methods** and **Main text**). While the horizontal bars indicate TBP-lobe like sequences, blue and red colours respectively denote N-terminal and C-terminal lobes for archaea, eukaryotes and viruses. The most conserved residues in the ancestral sequences between three superkingdoms were considered as universally conserved residues (see **Methods**). The ancestral sequences were integrated with the common TBP-lobe numbering (CTN) system for further analyses. A snapshot of actual CTN assignment for a part of RefMSA is shown. In the CTN system, each residue position in RefMSA is referred as <Lobe>.<Secondary structure type and secondary structure number>.<Alignment position> (e.g. N.L5.1 refers to Phe 116 in yeast; See **Methods and above**).
- b) Dendrogram of TBP-lobe like regions: TBP-lobe RefMSA was used to construct this dendrogram (see **Methods**). The dendrogram indicates that TBP-lobe like regions of bacteria and viruses form distinct clusters, suggesting their significant divergence from TBP lobes of eukaryotes and archaea. For the construction of the dendrogram, we considered only sequences from representative organisms in the RefMSA alignment. We made sure to select sequences from the RefMSA that span the whole diversity of lineages represented in the RefMSA. Poxvirus sequences were not considered for the generation of the dendrogram as they are the most divergent versions of TBP-lobe like sequences (see **Methods**).
- c) Approach to uncover molecular signatures for functional innovations of TBP involved utilization of two independent contexts i.e., spatial and temporal. To identify molecular signatures of TBP, evolutionary range and TBP-lobe regions from the RefMSA were used to define the temporal and spatial contexts, respectively. Here, two insets together define a particular spatio-temporal context, wherein temporal context is confining the analyses to evolutionary range of archaea and eukaryotes. The spatial context is the C-terminal TBP-lobe centric investigations.



*Viral residues corresponding to conserved
C-terminal residue signatures*

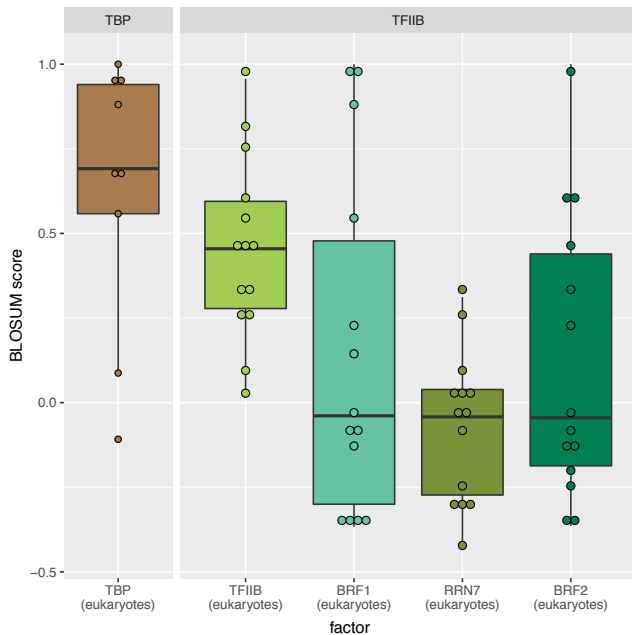


Supplementary Figure 3. C-terminal lobe specific molecular signature in viral TBPs

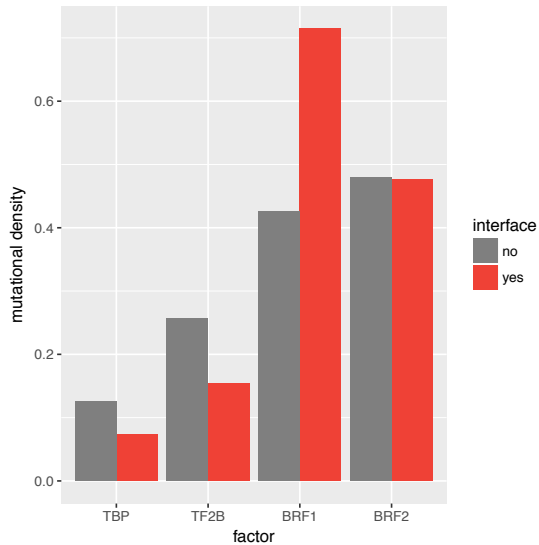
C-terminal lobe specific molecular signatures from archaea and eukaryotes (**Figure 3a**) mapped onto C-terminal lobes of representative viral TBPs from the RefMSA. This mapping indicates there is not high conservation of Glu or Asp at L3.2 or L3.4 in the C-terminal lobe of viral TBP-lobe sequences.

a

*Conservation between eukaryotic
paralogous TBP-TF2B/BRF2 interactions*

**b**

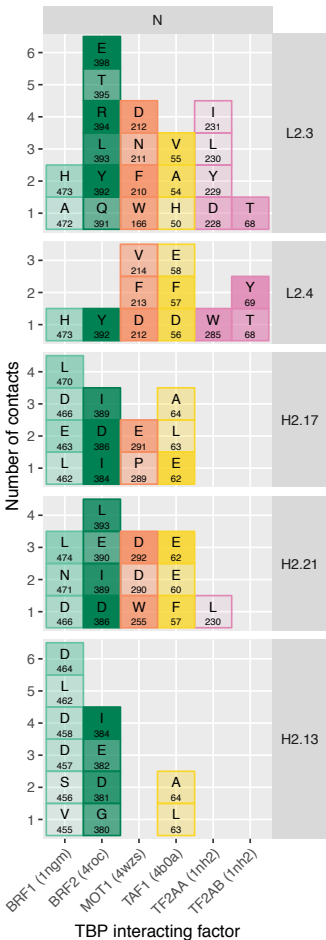
*Natural variation of TBP-interacting
residues with factors in the C-terminal lobe*



Supplementary Figure 4. Evolutionary conservation and natural variation of interaction interface residues in the C-terminal lobe of TBP

- a) The box-plots display the distribution of evolutionary conservation at the residue-level, measured in terms of normalized BLOSUM scores, of interaction mediating interface residues for eukaryotic: (i) orthologs of TBP and TFIIB and (ii) paralogs of TFIIB. The distribution of the evolutionary conservation was derived based on co-complex crystal structure (PDB: 1c9b) data of human TBP-TFIIB interactions as the reference for the comparison. The structure of human TBP-TFIIB was used to identify interaction mediating residues in both TBP and TFIIB. These residues were mapped on to equivalent residue positions in the following subsets of sequence alignments to evaluate their respective normalized BLOSUM scores across the alignment: (i) eukaryotic C-terminal TBP-lobes from RefMSA (for TBP only), (ii) TFIIB orthologs alignment in eukaryotes, (iii) alignment of TFIIB/BRF1 orthologs, (iv) alignment of BRF2 orthologs and (v) alignment of TAF1B/Rrn7 orthologs (see **Supplementary data**). Horizontal line within each of the box plots indicates median score of conservation. TBP displays highest conservation, relative to the other interacting factors, at these interface residue positions.
- b) The bar plot depicts mutation densities for natural variation of residues at the interaction interfaces of C-terminal lobe of TBP with TFIIB and its paralogs. Mutation density (rate of mutations normalized to the sequence length) bar plots has been made based on data from known co-complex structures of TBP C-terminal lobe with various factors such as TFIIB, BRF1, BRF2 (PDB: 1c9b, 6f40 and 4roc). “Red bar” indicate mutation density for the interface residues, while “gray bar” indicates the same measure for non-interface residues. The bar plots indicate that TBP (C-terminal lobe) has the least natural variation at the interface as compared to the other factors (**Methods**). Hence similar to the pattern of evolutionary conservation within eukaryotes, natural variations indicate that interface residues of TBP interacting factors does display a greater mutational tendency compared to TBP interface residues for the C-terminal lobe interactions.

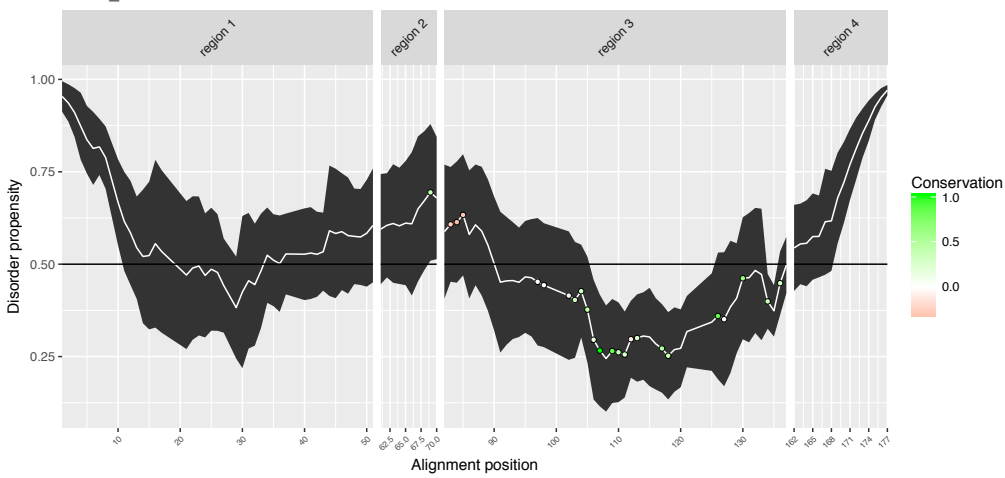
Interactions between TBP and factors
(“periodic table view”)



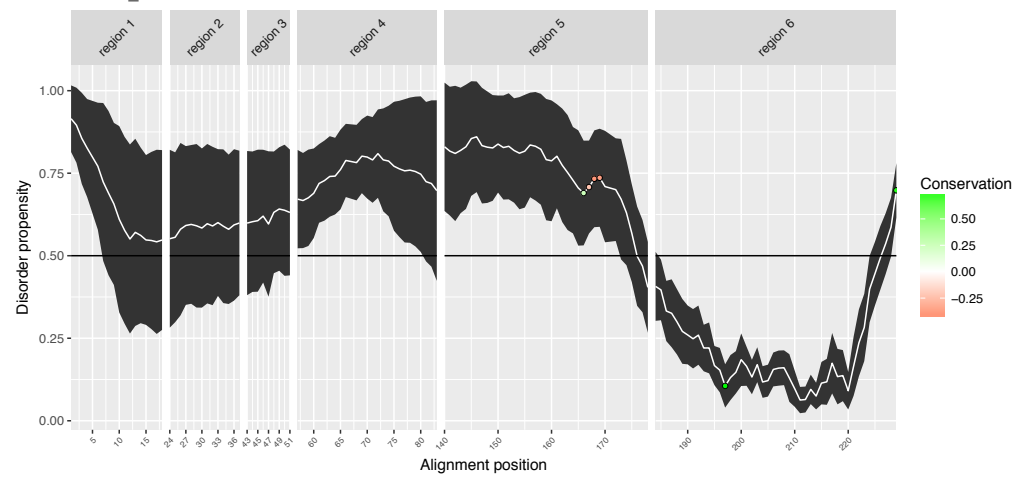
Supplementary Figure 5. Interaction residues of various factors that mediate interactions with the molecular signature positions in the N-terminal lobe of TBP.

The plot displays detailed view of number of interaction contacts of the 5 conserved molecular signature positions (shown in the right y-axis; see **Figure 4**) in the N-terminal lobe TBP and their interacting residues in various factors (shown in the x-axis). The darker the color of boxes higher the evolutionary conservation of that residue. These interactions are deduced based on available co-complex structures of these factors with TBP (**Figure 4** and **Methods**). The majority of the residues displayed in the plot are either acidic or aromatic residues suggesting the existence of a significant number of electrostatic interactions between N-terminal lobe of TBP and these factors (see **Figure 4**). These acidic or aromatic residues in some cases do not display a strong evolutionary conservation as these residues fall in intrinsically unstructured or disordered regions. However, while not conserving individual acidic or aromatic residues, overall they preserve negatively charged or aromatic characteristics in these residue sequence neighborhoods (see **Supplementary Figure 6** and **Figure 4**).

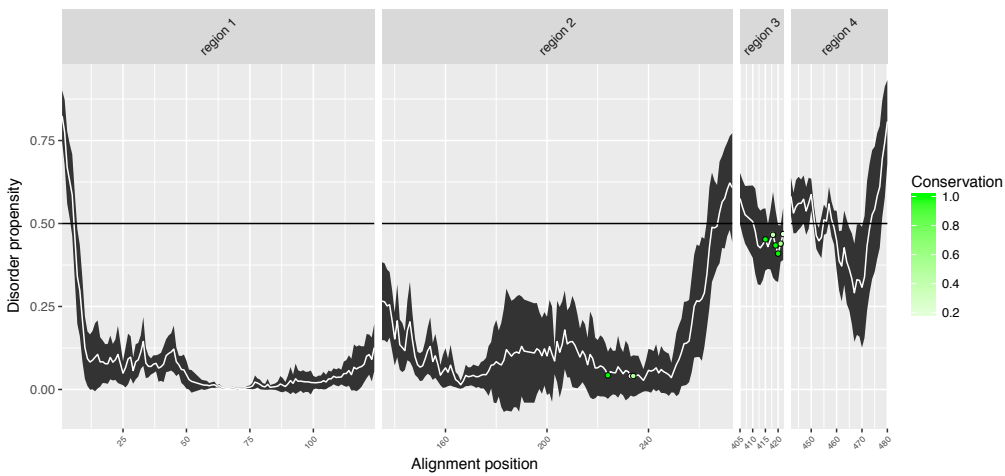
BRF1_all.fasta



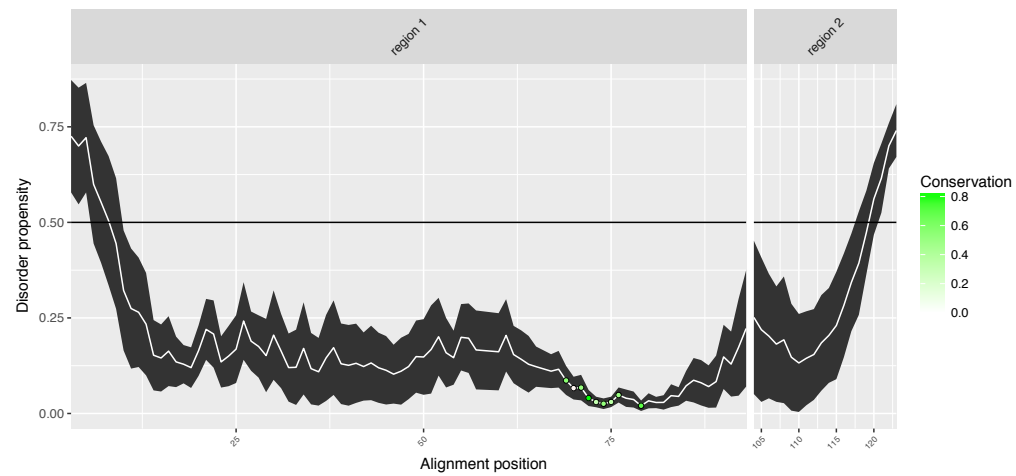
TF2AA_all.fasta



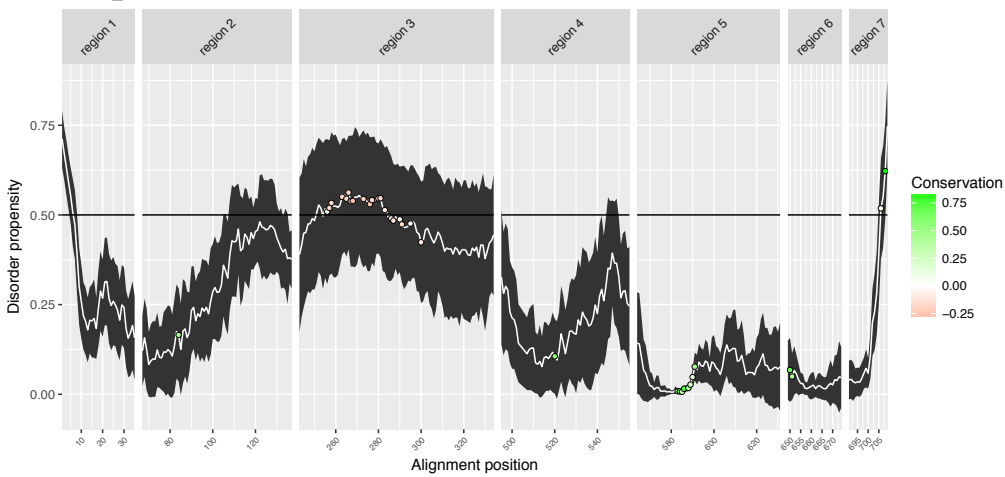
BRF2_all.fasta



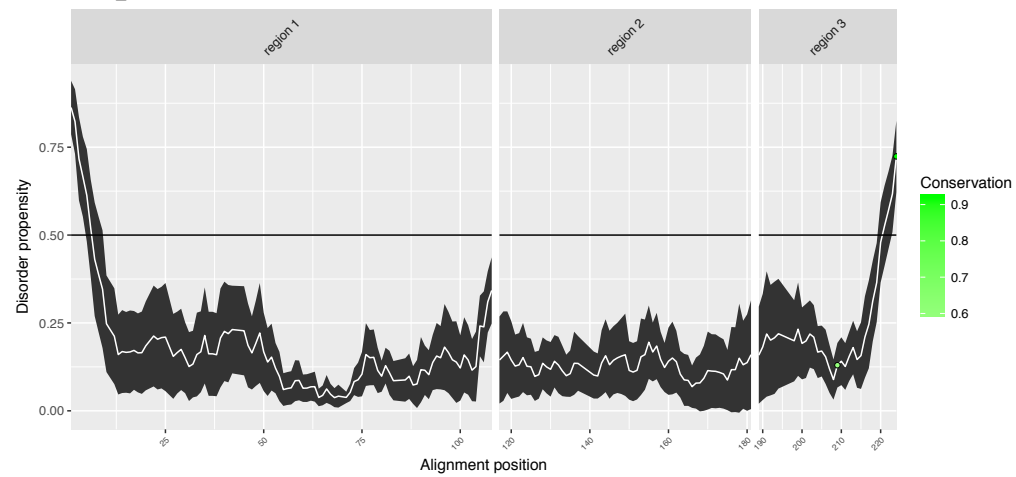
TF2AB_all.fasta



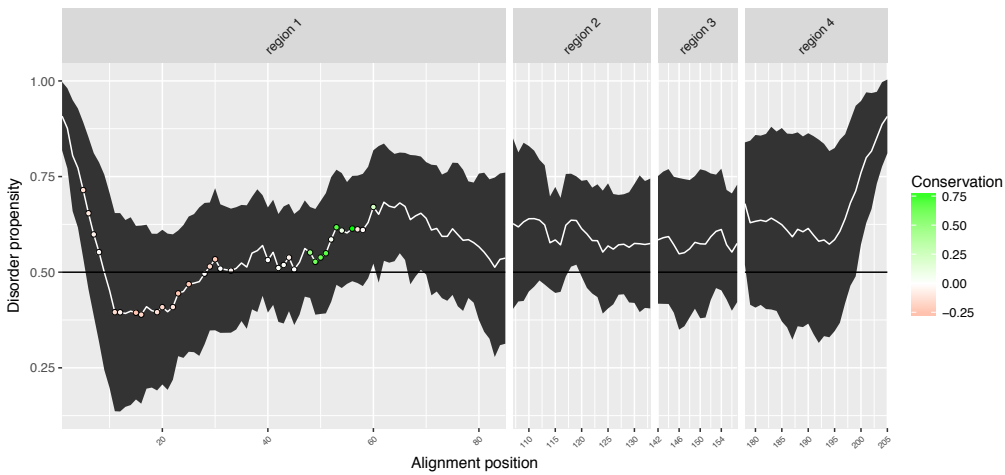
MOT1_all.fasta



TFIIB_all.fasta



TAF1_all.fasta

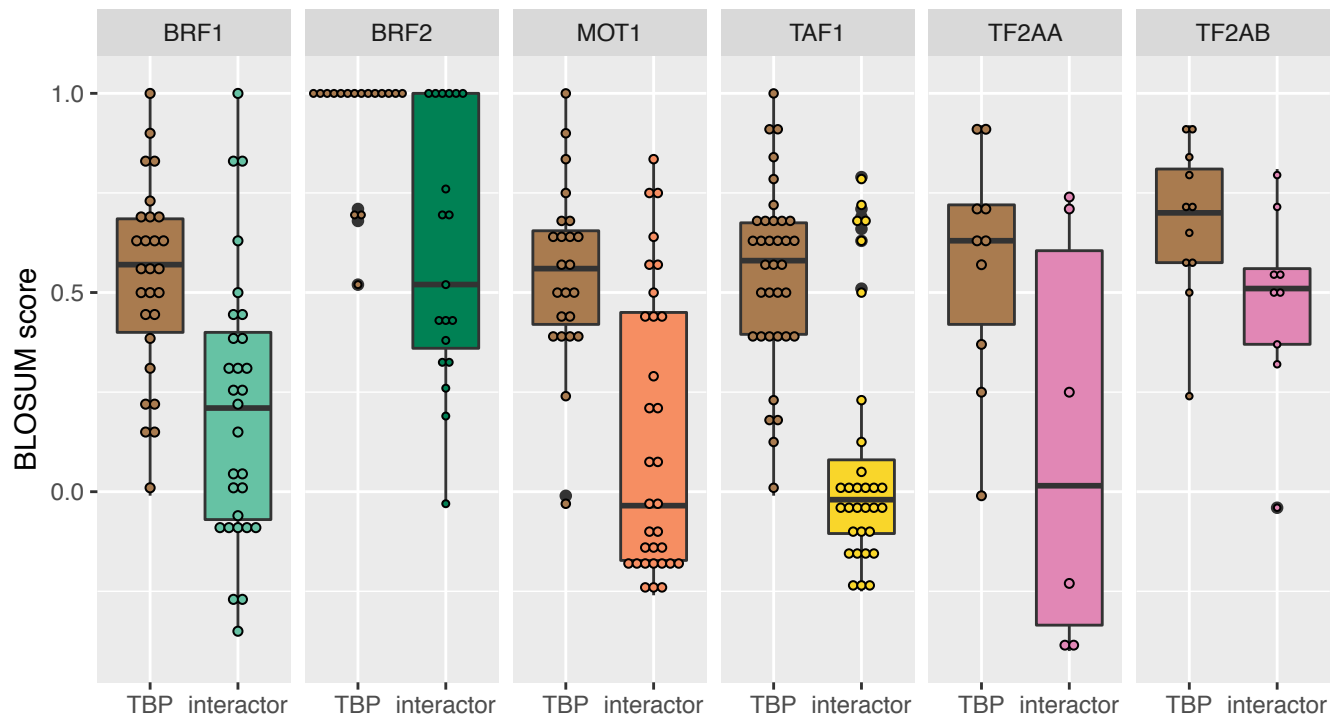
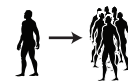


Supplementary Figure 6. Disorder propensities of the interacting segments of various factors that interact with the N-terminal lobe of TBP in eukaryotes.

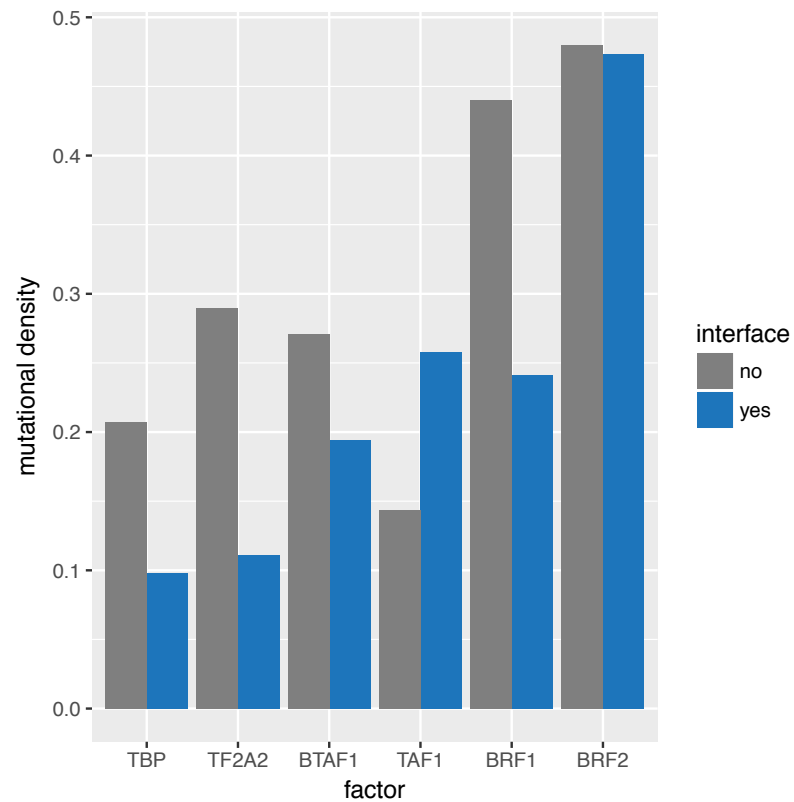
Each plot indicates the disorder propensity, a measure of the unstructured nature of a given residue, for interacting residues of every TBP interacting factor with N-terminal lobe of TBP. The disorder propensity is calculated using IUPRED (see **Methods**) and higher values in y-axis indicate greater tendency for being intrinsically unstructured or disordered for that given residue. Typically, values in the y-axis greater than 0.5 signify assignment of disorder tendency to a given residue. Disorder propensities across organisms, as in the alignment of TBP interacting factors (see **Supplementary data**), are shown as black shaded regions in each plot and central white lines indicate mean value of disorder propensity for the respective positions (across organisms). The residue conservation is indicated by BLOSUM score is shown as green for high conservation to red for poor conservation.

a

Conservation of interaction mediating residues between TBP and various factors

**b**

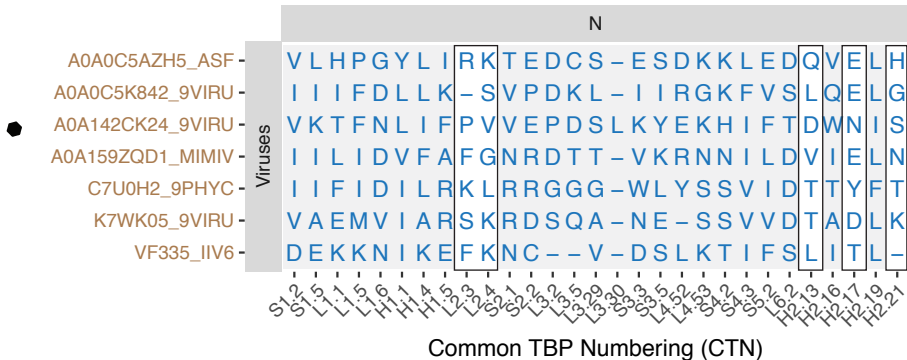
SNP based conservation comparison between TBP and its factors



Supplementary Figure 7. Evolutionary conservation and natural variation of interaction interface residues at TBP N-terminal lobe

- a) Evolutionary conservation, measured in terms of BLOSUM score, of interaction mediating interfaces residues for TBP and its various interacting factors at the N-terminal lobe of TBP. This is identified based on co-complex structure data of TBP with various factors; in particular, the extent of evolutionary conservation or divergence for various intermediating residues in various factors are evaluated based on sequence alignment (see **Supplementary data**). It should be noted that conservation plots are made for TBP and its interacting factors with an identical phylogenetic range in order to have a meaningful comparison. Horizontal lines within the box plot indicate median score of conservation for each factor's interface residues. In the majority of cases, TBP displays better evolutionary conservation than its interaction partners at the interaction mediating interface residue positions.
- b) The bar plots of natural variation of residues measured as mutational densities (see **Methods and above**) at the interface for various factors that interact with the N-terminal lobe of TBP. "Blue bars" indicate mutation density for the residue at the interface mediating interactions, while "gray bar" represents mutation density for non-interface residues. The barplots indicate TBP has the least natural variation at the interface as compared to other factors. Hence, similar to the pattern of evolutionary conservation and natural variation of interaction interfaces at the C-terminal lobe of TBP, these plots indicate that the interface residues of TBP interacting factors does display a greater mutational frequency compared to TBP interface residues.

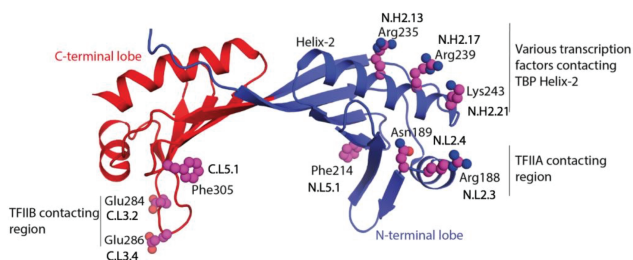
*Viral residues corresponding to conserved
N-terminal residue signatures*



Supplementary Figure 8. Molecular signature residue positions mapped to the N-terminal lobe of TBP of dsDNA viruses.

Highly conserved residue positions in the N-terminal lobe of TBP (see **Figure 4**) mapped on their dsDNA viruses orthologs. The five residue positions that are either positively charged positions or contains asparagine are indicated within rectangular boxes. It is clear from these indicated positions that, in the viruses, these positions, by and large, are devoid of positively charged residues. Given the positively charged residues are critical for mediating TBP interactions (**Figure 4** and **Supplementary Figure 5**) with various factors, which in turn function as regulatory controls, this is suggestive that viruses potentially avoid the host regulatory controls. These regulatory controls could otherwise place restrictions on functions of viral TBPs in their respective hosts.

a



b

| Mutation | Organism | Phenotype | Lobe | CTN | Signature Type | Pubmed ID |
|-----------------|----------|--|------|---------|-----------------|-----------|
| F116Y | Yeast | Electrophoretic mobility shift assay with promoter DNA probe | N | N.L5.1 | Universal | 2015629 |
| F116A | | | | | | 24865972 |
| F207A | | | | | | 24865972 |
| F214A | Human | Electrophoretic mobility shift assay with promoter DNA probe | C | C.L5.1 | | 16179647 |
| F305A and F305K | | | | | | 16179647 |
| F288K and F288A | Human | Significant reduction in transcription activity (<i>in vitro</i>) at one or more promoters | C | C.L3.27 | | 12535529 |
| F305A and F305K | | | | | | 12535529 |
| F305A | Human | Significant reduction in basal transcription (<i>in vitro</i>) | C | C.L5.1 | | 8843200 |
| S307F | | | | | | 8843200 |
| K316E | | | | | | 8843200 |
| L165E | | | | | 8843200 | |
| D179R | | | | | 8843200 | |
| A184E | | | | | 8843200 | |
| R188E | | | | | 8843200 | |
| R188E | | | | | 12535529 | |
| N189E | Human | Significant partial reduction in transcription activity (<i>in vitro</i>) at one or more promoters | N | N.L2.3 | 12535529 | |
| N189E | | | | | 8843200 | |
| A190E | Human | Significant reduction in basal transcription (<i>in vitro</i>) | N | N.L2.4 | 8843200 | |
| A190E | | | | | 8843200 | |
| E191K | Human | Significant reduction in transcription activity (<i>in vitro</i>) at one or more promoters | N | N.S2.2 | N-lobe specific | 12535529 |
| E191R | Human | Partial reduction in basal transcription (<i>in vitro</i>) | N | N.S2.2 | | 8843200 |
| E93R | Yeast | Gene expression changes in yeast system and perturbation to interaction with TFIIA | N | N.S2.2 | 17407552 | |
| N193R | Human | Significant reduction in basal transcription (<i>in vitro</i>) | N | N.L3.2 | 8843200 | |
| E206R | Human | Marginal reduction in basal transcription (<i>in vitro</i>) | N | N.L4.52 | 8843200 | |
| R235E | Human | Potential increase in basal transcription (<i>in vitro</i>) | N | N.H2.13 | 8843200 | |
| R239S | | | | | 8843200 | |
| K145E | Yeast | Gene expression changes in yeast system and perturbation to interaction with Mot1p | N | N.H2.21 | 17407552 | |
| S261E | Human | Significant reduction in basal transcription (<i>in vitro</i>) | C | C.S1.8 | C-lobe specific | 8843200 |
| D263R | | | | | | 8843200 |
| E284R | | | | | | 8843200 |
| E284A and E284R | Human | Significant reduction in transcription activity (<i>in vitro</i>) at one or more promoters | C | C.L3.2 | | 12535529 |
| E286R | | | | | | 12535529 |
| E286R | Human | Significant reduction in basal transcription (<i>in vitro</i>) | C | C.L3.4 | | 8843200 |

Supplementary Figure 9. Mapping of Mutational data.

- a) Molecular signature residues of TBP and experimental data. Highlighted residues that are prominent molecular signatures for interactions that have support from experimental data. For more comprehensive data please see the table in b) below.
- b) Mutational data that provide support for molecular signatures.

Supplementary Figure 10. Conserved regions, protein-protein interactions and tissue expression in TBP and its paralogs.

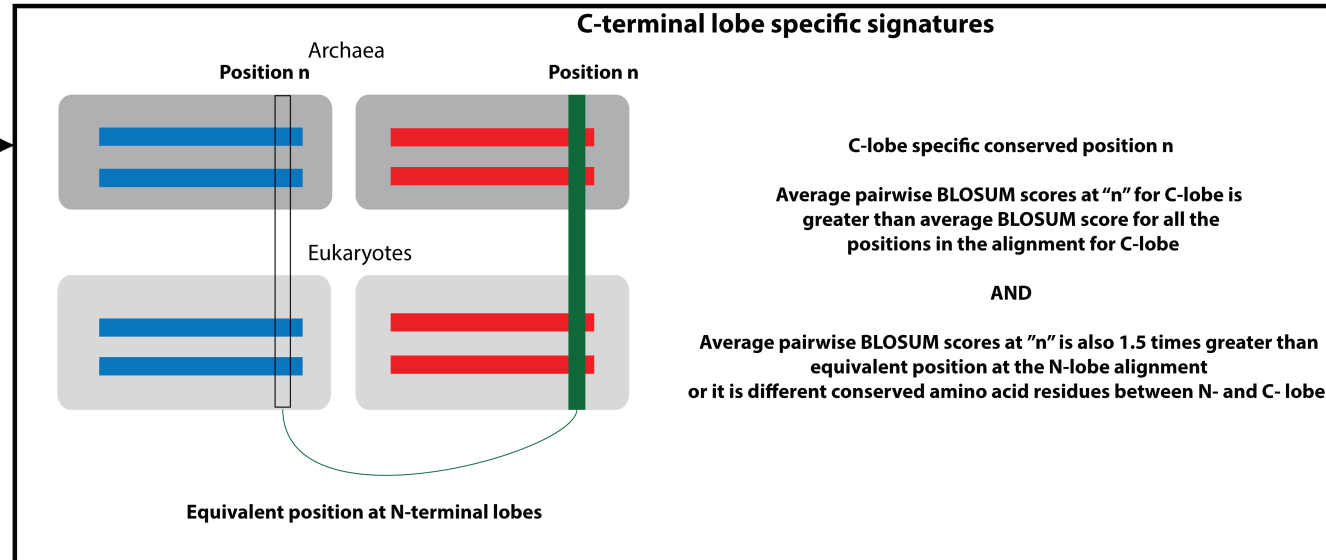
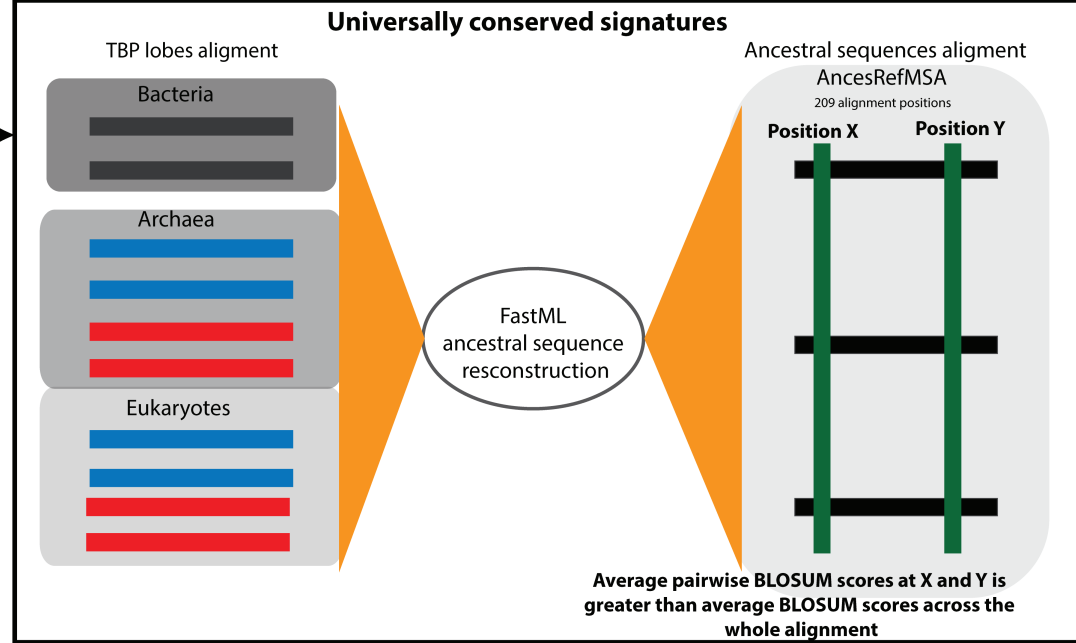
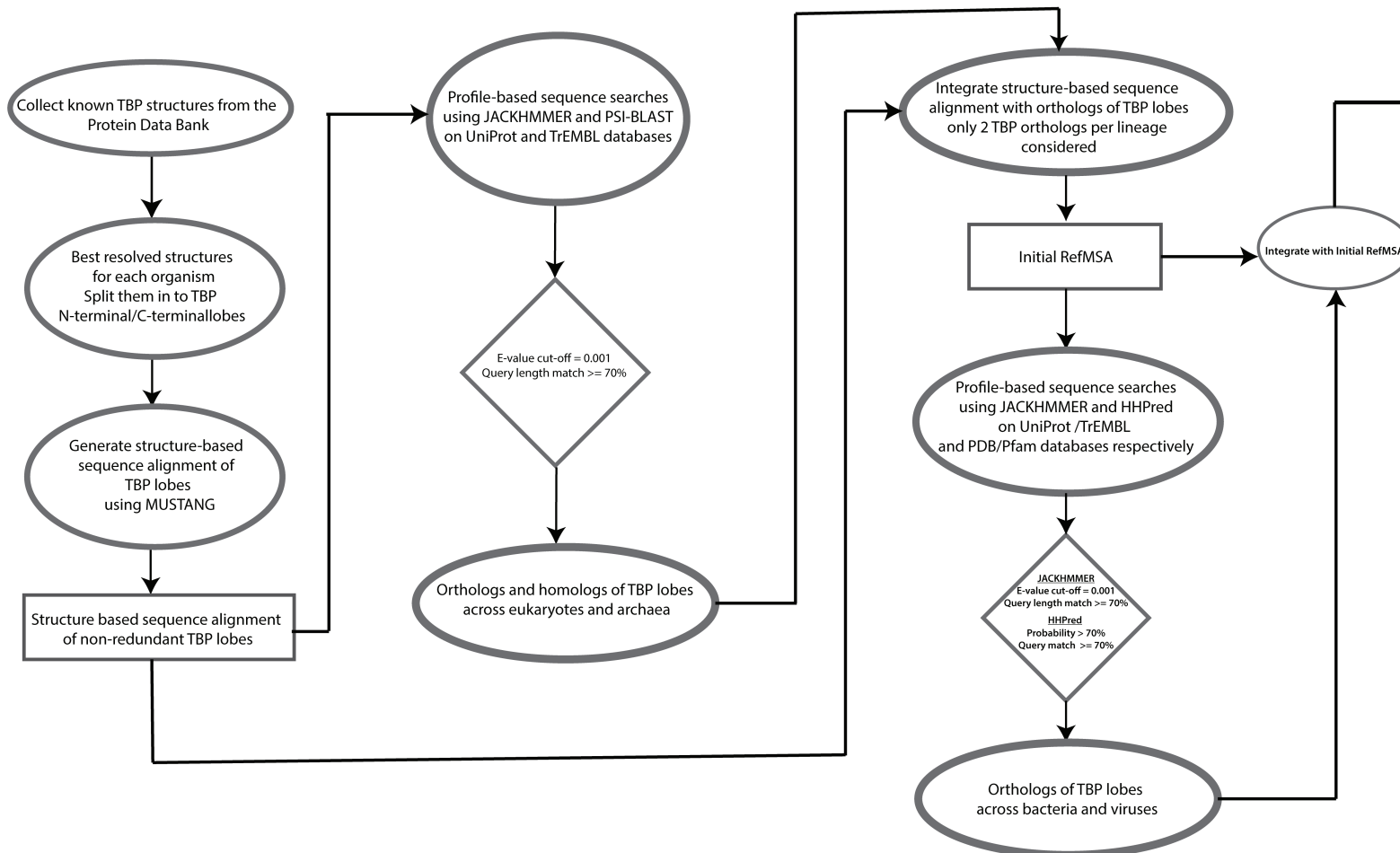
- a) Dendrograms depicting relationship between TBP and its paralogs (**Supplementary Data**). These dendrograms have been constructed using representative sequences that capture diversity of animal lineages in which these proteins were detected (see **Methods**). TBP and TBPL2 are found in early or primitive animals, while TBPL1 is found only in invertebrates and vertebrates suggesting that TBPL1 is the most recent paralog of TBP as well as the most divergent.
- b) Multiple sequence alignment of N-terminal PolyGln stretches containing regions of TBP. Only a part of the alignment of this evolutionary conserved region is shown (see **Supplementary Data** for full alignment of this domain). This region is present only in vertebrates. It is clear that there is an approximate linear increase in sequence length of PolyGln (PolyQ) stretches in TBP from zebra fish to humans.
- c) Multiple sequence alignment of N-terminal proline rich regions of TBPL2. Only a part of the alignment of this evolutionary conserved region is shown (see **Supplementary data** for full alignment of this region). This region is present only in vertebrates and particularly prominent in mammals.
- d) Physical interactions of human TBP, TBPL1 and TBPL2. While the data is sparse for TBPL1 and TBPL2, their shared interactions with TBP appears statistically significant. However, albeit with limited protein-protein interaction data, we observe the existence of unique interaction partners of TBPL1 and TBPL2. This is suggestive of potential distinct function contexts for TBPL1 and TBPL2. Given the higher sequence divergence and distinct pattern of protein expression across human tissues of TBPL1, interaction divergence could be additional contributing factor for functional divergence between TBP and TBPL1 (see **Supplementary Figure 10e**). The physical interaction data has been obtained from BIOGRID and Intact databases (Oughtred et al, 2019; Kerrien et al, 2012).
- e) Protein level expression of human TBP, TBPL1 and TBPL2 across 17 different tissues (Kim et al, 2014). TBPL1 is almost ubiquitously expressed in majority of these tissues, while TBP and TBPL2 are more restricted. This may be due to their low-level of expression, which might be hard to detect by high throughput methods. Nevertheless, TBPL1 appears to have a divergent expression pattern indicated by Pearson correlation of 0.2 with TBP. This along with the earlier observation that TBPL1 is most divergent of the paralogs of TBP is suggestive of distinct functional niche of TBPL1. TBPL2 has a robust co-expression with TBP with a Pearson correlation of 0.7 and this might suggest some level of functional overlap between TBP and TBPL2 for e.g. DNA sequence recognition, protein-protein interaction, etc.

Analysis Overview

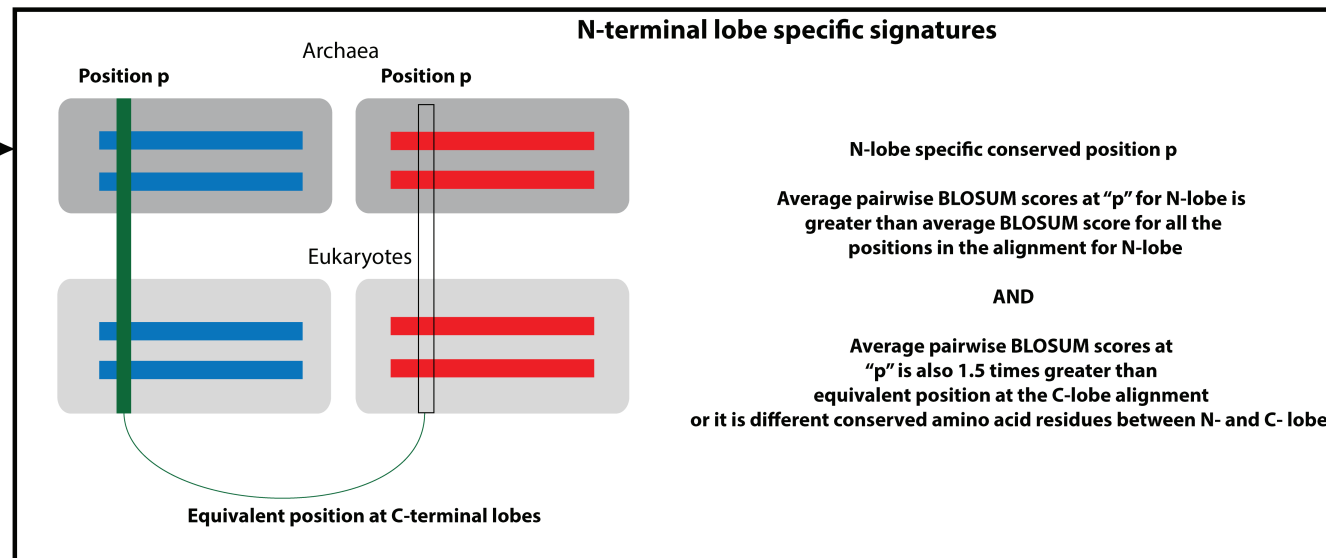
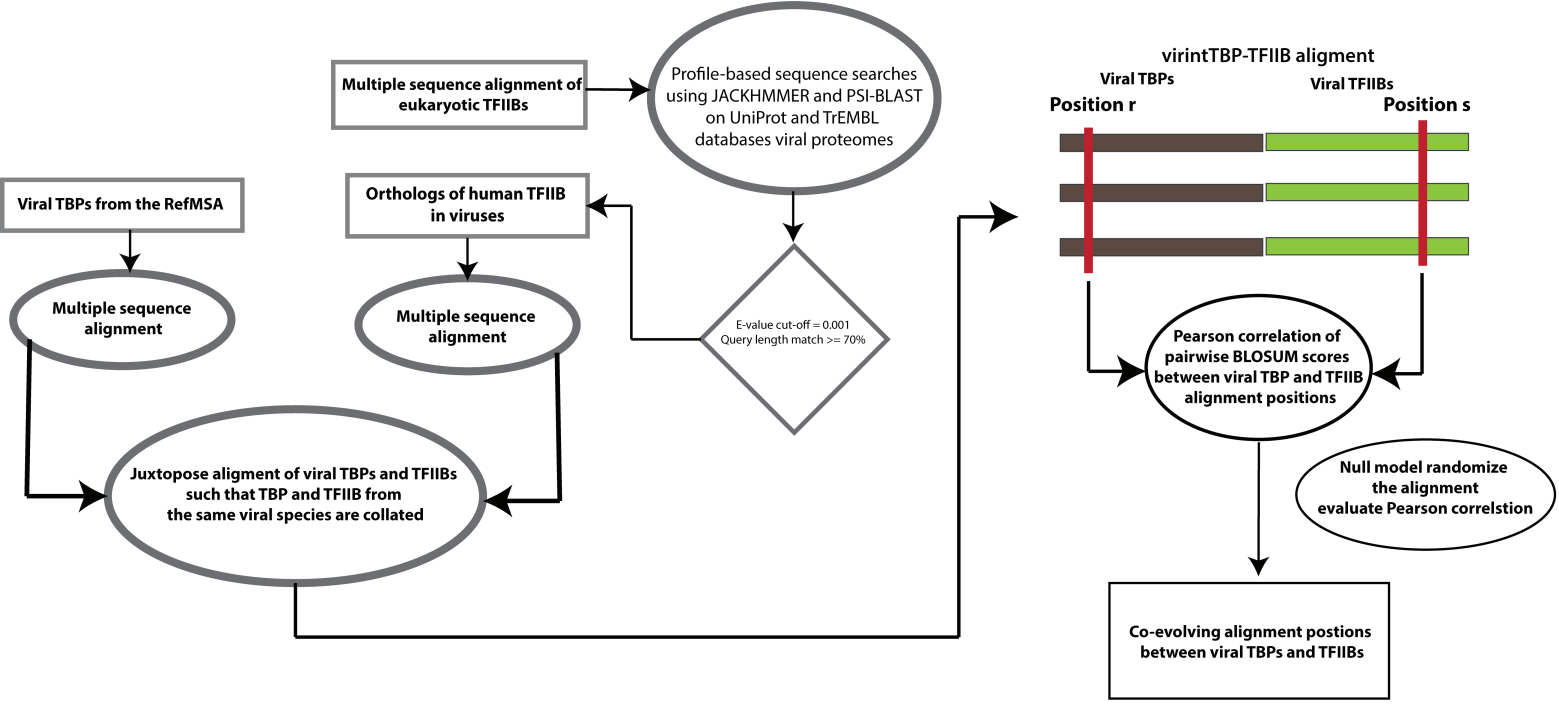
| | |
|---|--|
| Comparison on N- and C-terminal lobes of TBP | ✓ |
| Evolution of interactions | ✓ |
| Region of main focus in TBP | All regions in TBP were considered |
| Number of adaptors studied | 8 adaptors |
| Number of sequences analyzed | ~400 sequence domains |
| Evolutionary range covered | Bacteria, Archaea, Eukaryotes and Viruses were considered |
| Reduction of potential bias | Representatives from every lineage were carefully considered to minimise bias in discovery of signature residues |
| Number of structures analyzed | 10 co-complex structures |
| Number of molecular signature residues identified | 5 residues for ds nucleic acid binding 2 residues for C-lobe interactions 5 residues for N-lobe interactions |
| Integration of additional datasets | Genome-scale data such as transcriptomic and proteomics data were integrated |
| Support from mutational data (new analysis) | ✓ |
| Distinguishing orthologs and paralogs | ✓ |
| Common numbering system | ✓ |
| Comprehensive resource of alignments of TBP, TBP-like proteins and its interaction partners | ✓ |
| Approach applicable to other protein families | ✓ |

Supplementary Figure 11: Overview of datasets used and analyses performed in Ravarani *et al.*

Overall strategy for the construction of RefMSA



Co-evolution of Viral TBP and Viral TFIIIB interactions



Supplementary Figure 12: Overview of Methods.