# Efficient Online Transcription Factor Binding Site Adjustment by Integrating Transitive Graph Projection with MoRAine 2.0

Tobias Wittkop<sup>\*1,</sup>, Sven Rahmann<sup>2,</sup>, Jan Baumbach<sup>3,4,5,</sup>

<sup>1</sup>Genome Informatics, Bielefeld University, Bielefeld, Germany

<sup>2</sup>Bioinformatics for High-Throughput Technologies, Computer Science XI, TU Dortmund, Dortmund, Germany

<sup>3</sup>International Computer Science Institute, Berkeley, USA

<sup>4</sup>Max Planck Institute for Informatics, Germany

#### <sup>5</sup>Saarland University, Cluster of Excellence for Multimodal Computing and Interaction, Germany

#### Summary

We investigated the problem of imprecisely determined prokaryotic transcription factor (TF) binding sites (TFBSs). We found that the identification and reinvestigation of questionable binding motifs may result in improved models of these motifs. Subsequent model-based predictions of gene regulatory interactions may be performed with increased accuracy when the TFBSs annotation underlying these models has been re-adjusted.

We present MoRAine 2.0, a significantly improved version of MoRAine. It can automatically identify cases of unfavorable TFBS strand annotations and imprecisely determined TFBS positions. With release 2.0, we close the gap between reasonable running time and high accuracy. Furthermore, it requires only minimal input from the user: (1) the input TFBS sequences and (2) the length of the flanking sequences.

Conclusions: MoRAine 2.0 is an easy-to-use, integrated, and publicly available web tool for the re-annotation of questionable TFBSs. It can be used online or downloaded as a stand-alone version from http://moraine.cebitec.uni-bielefeld.de.

## 1 Background

We recently considered the problem of imprecisely determined transcription factor (TF) binding sites (TFBSs), which are stored in public databases that are dedicated to prokaryotic gene regulatory networks [19]. We found that the identification and reinvestigation of questionable binding motifs may result in improved *in silico* models of these motifs. Hence, subsequent model-based bioinformatics predictions of gene regulatory interactions may be performed with increased accuracy [10]. We developed two tools that tackle the problem of TFBS re-adjustment: (1) MotifAdjuster, a stand-alone software and (2) MoRAine 1.0, a fast heuristic. While the MotifAdjuster algorithm is based on expectation maximization and highly accurate,

<sup>\*</sup>Corresponding author: Tobias.Wittkop@CeBiTec.Uni-Bielefeld.DE.

MoRAine is optimized to provide a fast, online-available web solution but asks the user for various parameters to trade between running time and solution quality.

Here, we present a significantly improved version of MoRAine, which now closes the gap between reasonable running time and high accuracy. As release 1.0, MoRAine 2.0 is an online tool but it requires only minimal input from the user: (1) the input TFBS sequences and (2) the length of the flanking sequences, i.e. the expected length of the TF binding motif (TFBM). In our study, we focus on prokaryotic TFBSs, i.e. those of *Escherichia coli*.

In the following, we introduce the biological background, explain why annotation problems may occur, and show why this may result in a poor TFBS prediction performance.

#### Transcription factor binding site annotation - A difficult and error-prone task

Transcription factors are important components of the cell's regulatory machinery. They are DNA-binding proteins that are able to detect intra- and extracellular signals. By binding to so-called transcription factor binding sites they control the expression of their target genes and thereby decisively influence genetic programs like growth, reproduction, and defense [7,1,2,24, 22]. Given a set of known TFBSs for a certain regulator, we can build mathematical models to perform *in silico* predictions of further TFBSs in order to predict regulatory networks. This task is generally complicated by the relatively low level of TFBS conservation. The most widely used model for TFBSs are so-called position frequency matrices (PFMs) [25]. PFMs can be converted to position specific scoring matrices (PSSMs) by calculating log-odds scores. These matrices are used in turn to predict TFBSs in the upstream sequences of putative target genes for a certain TF. Various software tools are available: PoSSuMsearch [11], Virtual Footprint [21], MATCH [20], and P-MATCH [14], to name a few.

Nowadays, TFBS wet lab determination is done by electrophoretic mobility shift assays (EMSA) [17], DNAse footprinting [16], ChIP-chip [26], ChIP-seq [18], or mutations of putative TFBMs and subsequent expression studies. All of these methods lack a precise binding sequence identification that is accurate to one base pair [8]. Another problem occurs: Since TFs bind the double-stranded DNA, it is a matter of interpretation which strand of the TF-binding sequence is annotated. Clearly, both issues directly affect and complicate TFBS modeling as position frequency matrices and hence, the subsequent PSSM-based binding site prediction. This problem occurs when a TFBS from either strand based on approximate knowledge of its position is entered in a reference database and subsequently used blindly for PSSM-based predictions. This does happen in practice for regulatory databases that integrate information from other sources [19, 10], for instance, in CoryneRegNet [3, 4, 5, 8, 9]. Here, the data is added manually to the data repository by curation teams. They scan the literature and transfer the published knowledge into a formal data structure. This task is difficult, error-prone and, hence, further supports putative mistakes with the TFBS annotation process.

For mis-annotated TFBSs, we may observe a poor information content of the subsequently computed PFM, which consequently leads to a decreased binding motif prediction for the PSSM that was constructed from that PFM. We attack this problem by re-annotating the TFBSs by possibly switching their strands and/or shifting them a few positions, in order to maximize the information content of the resulting PFM.

Therefore, MoRAine 1.0 provided a combination of two clustering approaches, a variant of

k-means (km) and cluster growing (cg), applied to two similarity measures, cluster similarity (simC) and seed node similarity (simS). In the following, we describe why and how we replaced these methods in MoRAine 2.0 in order to provide better results within decreased running times. Afterwards, we compare both releases and demonstrate the increased TFBS prediction performance. Subsequently, we show that the PFMs resulting from MoRAine-adjusted TFBSs significantly improve the prediction performance of further binding sites in practice, since the adjusted TFBSs lead to PFMs with higher information content.

### 2 Methods

We need the following definitions to explain how MoRAine works and to compare the readjustment performances of MoRAine 1.0 and 2.0.

Let  $\Sigma := \{A,T,C,G\}$  be the DNA alphabet. In accordance with [10], a position frequency matrix  $F = (f_{\sigma j})$  for a set of *n* TFBSs of length *m* over the alphabet  $\Sigma$  is defined as a  $|\Sigma| \times m$  matrix, where  $f_{\sigma j}$  is the relative frequency of letter  $\sigma$  at position *j*.

Crooks *et al.* introduced in [15] the information content as quality measure for PFMs. The information content  $I_j$  for column j of F is defined as

$$I_j := \log_2 |\Sigma| + \sum_{\sigma \in \Sigma} f_{\sigma j} \cdot \log_2 f_{\sigma j} \quad \text{[bits]}.$$

If all symbols at position j agree,  $I_j$  reaches its maximum with maximal value 2 bits for a 4-letter alphabet  $\Sigma$ . The mean information content I(F) for a given PFM F is defined as the average  $I_j$  over all positions j:

$$I(F) := \frac{1}{m} \sum_{j=1}^{m} I_j.$$

In what follows, we use the mean information content I(F) as a quality measure for a given PFM F and denote it shortly with I if F is fixed. We will use the information content to compare the quality of two different PFMs  $F_1$  and  $F_2$  by comparing  $I(F_1)$  with  $I(F_2)$ . If  $F_2$  is the PFM of the MoRAine-adjusted TFBSs, while  $F_1$  is the PFM computed from the input TFBSs, with  $I(F_1) \leq I(F_2)$ , we can calculate the percentage improvement performance P with  $P = 100 \cdot \frac{I(F_2)}{I(F_1)}$ .

MoRAine now works as follows: The input is a set of n annotated length-m TFBS sequences that extend l bp to the left and r bp to the right. Hence, the length of the given input sequences is  $m^+ := m + l + r$ . First, MoRAine computes the set M of every possible motif of length  $m = m^+ - l - r$  derived by the operations *shift* and *switch* applied to each of the n input sequences. The operation *shift* provides every substring of length m for a given motif of length  $m^+$ , and the operation *switch* its reverse complement sequence. We obtain a set  $S_i$  of M := $|S_i| = 2 \cdot (l + r + 1)$  potential TFBS sequences of length m for each input sequence i, with  $i = 1, \ldots, n$ .

So far MoRAine 1.0 and 2.0 work in a similar way. For both, the goal is to find a set C of TF-BSs that contains exactly one TFBS from each  $S_i$  and maximizes the mean information content of the corresponding PFM  $F_C$ . As mentioned earlier, MoRAine 1.0 offers two heuristic clustering algorithms, (cg) and (km), both working on either of two similarity functions, (simC)

and (simS). Table 1 summarizes the running times and TFBS annotation improvement performance of MoRAine 1.0 for all four combinations. One can see a trade-off between accuracy and running time: (cg/simS) provides best results but (cg/simC) is much faster (see section Results and Discussion for more details).

With MoRAine 2.0, we close this gap and provide a powerful tool that now provides better results than MoRAine 1.0 with (cg/simS) at running times equal to (cg/simC). The goal can be cast as follows: We partition the set of input TFBSs into  $M = 2 \cdot (l + r + 1)$  clusters, where each cluster contains exactly n motifs, one of each  $S_i$  (i = 1, ..., n) and thus is a putative solution. In the following, we describe how we adapted and integrated Transitivity Clustering with MoRAine 2.0 to find such a set C.

Transitivity Clustering is a clustering method based on the weighted transitive graph projection (WTGP) problem. By solving this NP-complete graph modification problem, objects are partitioned into groups of similar elements. We briefly describe the underlying WTGP problem and how it has been modified to fulfill the needs for this specific task. Given a set of objects V and a pairwise similarity function  $s : \binom{V}{2} \to \mathbb{R}$  a similarity graph G = (V, E) is constructed, where V is the set of objects and E is the set of undirected edges between these objects. An edge is present in G if the similarity between the adjacent nodes exceeds a user defined threshold t. The goal is to find a *transitive* graph G' = (V, E') with smallest distance to G. Transitivity means that for all triples  $\{u, v, w\} \in \binom{V}{3}$ :  $\{u, v\} \in E$  and  $\{v, w\} \in E$  implies  $\{u, w\} \in E$ . A transitive graph is a disjoint union of cliques, also called a cluster graph. Therefore a cost function for adding/deleting edges is defined as c(u, v) = |t - s(u, v)| serving as optimization function. Thus, the distance between two edge sets is the cost of transforming one into the other. The resulting disjoint cliques of the minimum cost transitive G' represent the clusters we are looking for.

Transitivity Clustering is flexible and offers the possibility to integrate additional knowledge. As similarity function, instead of the functions from MoRAine 1.0, (cq) and (km), we now use the difference between the motif length  $\ell = |p| = |q|$  and the hamming distance h(p,q)between two TFBSs p, q, hence  $s(p, q) := \ell - h(p, q)$ . To ensure that each cluster of TFBSs contains only one motif from each set  $S_i$ , we set the similarity function s to  $-\infty$  if  $p \in S_i$ and  $q \in S_i$  for some  $S_i$ , i.e. if both potential solutions (the TFBSs p and q) originate from the same input TFBS. The threshold t is set to zero, which guaranties that each cluster contains exactly one TFBS from each set  $S_i$ . Transitivity Clustering has successfully been applied to protein family detection using the layout based heuristic FORCE to solve the NPcomplete WTGP problem [27]. Together with an exact fixed parameter algorithm developed by Böcker et al. [13] and the fast but less accurate heuristic CAST (Cluster Affinity Search Technique) by Ben-Dor et al. [12], this layout-based approach has been integrated into the clustering framework TransClust. The TransClust software combines the different methods to provide very accurate results in reasonable time. Its integration with MoRAine is mainly responsible for the increased performance of MoRAine 2.0, as we will demonstrate in the following section. Further information about Transitivity Clustering is available at the web site http://transclust.cebitec.uni-bielefeld.de.

# 3 Results and Discussion

### 3.1 Implementation

MoRAine 2.0 is an open source JAVA 6 program. It can accessed and downloaded at http://moraine.cebitec.uni-bielefeld.de. As shown for MoRAine 1.0 in [6], release 2.0 of MoRAine may be included into a database back-end as quality assurance tool or to provide a bioinformatics workflow with adjusted position weight matrices for TFBS predictions.

We emphasize that the main advantage of MoRAine is it's easy-to-use web interface. The user may copy and paste binding sequences in FASTA format at the MoRAine web site to calculate the adjusted motifs as well as the corresponding sequence logos by using the Berkeley web logo software [15]. Just as MoRAine 1.0, the second release is an easy-to-use alternative for the computation of sequence logos and the adjustment of transcription factor binding sites, but it now provides increased accuracy at decreased running times and an eased user-interface with less parameters to adjust.

#### 3.2 Increased information content improvement with MoRAine 2.0

In Figure 1 we exemplarily illustrate the output of the MoRAine online service for the binding sites of the transcription factor RamB of *Corynebacterium glutamicum*. The TFBSs have been taken from CoryneRegNet release 5.0. As in most databases, in CoryneRegNet [8], each binding site is annotated in  $5' \rightarrow 3'$  direction relative to the regulated target gene. By using MoRAine 2.0 we improved the average information content from 0.64 (original database TFBSs) to 1.15 (MoRAine-adjusted TFBSs) by switching the strands for 15 of the 38 input sequences. The computation time was less than 2 seconds.



Figure 1: A screenshot from the MoRAine 2.0 web site. A comparison of the sequence logos constructed from the original TFBSs (left side) for the transcription factor RamB of *Corynebacterium glutamicum* and the adjusted TFBSs by using MoRAine 2.0 (right side).

To demonstrate the performance, i.e. decreased running time and increased information content improvement, of MoRAine 2.0, we used the same datasets as in [10]: 1165 binding sites of 85 transcription factors of *Escherichia coli*. We compare the average runtime and the mean information content improvement of MoRAine 2.0 with the four methods implemented in MoRAine 1.0 for different lengths of the flanking sequences (l and r, respectively). As shown in Table 1, with MoRAine 1.0 the combination (cg/simC) had the best runtime, but to gain the best information content improvement, we used the combination (cg/simS) [10]. With the work