

Motif Analysis of Yeast Transcriptional Regulatory Regions

Tetsushi Yada¹ **Yasushi Totoki**² **Kenta Nakai**³
yada@tokyo.jst.go.jp totoki@imslab.co.jp nakai@imcb.osaka-u.ac.jp

¹ Japan Science and Technology Corporation (JST)
5-3 Yonbancho, Chiyoda-ku, Tokyo 102, Japan

² Information and Mathematical Science Laboratory, Inc.
2-43-1 Ikebukuro, Toshima-ku, Tokyo 171, Japan

³ Institute for Molecular and Cellular Biology, Osaka University
1-3 Yamada-oka, Suita 565, Japan

Most of the genetic information can be conveniently classified into two categories: the information on the amino acid sequences and the information on their regulation. With the recent progress of large-scale sequencing efforts, a large number of hypothetical amino acid sequences have been determined. However, the computational study of their regulatory sequences is still in its infant stage. To overcome this situation, we presented a new scheme to predict the function of open reading frames, identified in the genomes of *B. subtilis* [5] and *E. coli* through their sigma-factor dependency. Since eukaryotic cells do not have sigma-factors, something different approach should be explored for the analyses of eukaryotic ORFs. We report here our recent attempt in this direction.

Last year, we introduced the YEBIS system, which can be used to extract a set of motifs from unaligned DNA sequences without any prior knowledge [4]. Two of the strongness of YEBIS are its high speed and its capacity to treat with many sequences at a time. Moreover, the extracted motifs are represented in a hidden Markov model, which allows more flexible representation than other existing methods, such as the regular expression and the weight matrix. In addition, YEBIS outputs a set of motifs in a single run, which is practical because users do not know how many motifs exist in a given set of sequence data. Recently, we further improved its algorithm to detect conserved motifs using recently compiled benchmark data sets [1]. As a result, YEBIS shows comparative or sometimes better sensitivity compared with other programs but it can still process much more sequences quickly (Yada, *et al.*, submitted).

The above-mentioned feature of YEBIS, *i.e.*, its high capacity to process a large number of input sequences, allows us to exploit a new approach to discover any *cis*-elements that characterize a set of genes which are likely to be regulated by some common factors. First, two sets of genes are prepared: one set is related to a specific function, say, amino acid metabolism, while the other is known to be unrelated to its function. We used the functional classification of yeast genes by MIPS [2]. For both data sets, the upstream sequences of length 600bp were extracted and were applied to YEBIS. Then, the obtained two sets of motifs were compared using the group-to-group dynamic programming method, and the motifs which appear in the former set only were selected. Thus, we call this approach ‘*In silico* subtraction assay’. The extracted motifs were tested whether they are previously characterized ones or not by comparing with the known binding sites stored in the TRANSFAC database [3]. In Table 1, some examples of extracted motifs are shown. Further experiments using other datasets are also ongoing. Because of the recent progresses of both the DNA chip technology and the functional genome projects, there will be a number of works that systematically determine the expression profile of various genes. Thus, our approach should be useful to interpret these results from the sequence data.

Table 1: Typical Motifs Extracted from Various Functional Groups of Yeast Genes. Motifs are shown in a simpler representation than original HMM. The last column indicates the name of known transcription factors which have similar sequence specificity.

Group	Examples of Significant Motifs	Simil. Fact.
Metabolism	C ₁₀₀ T ₁₀₀ C ₆₀ s ₉₉ A ₁₀₀ G ₁₀₀ A ₅₆ n n w ₆₉ T ₅₉ y ₇₅ T ₈₇ C ₁₀₀ C ₈₅ G ₇₃ C ₁₀₀ G ₅₈ G ₅₇	CAR1 HSTF, CUP2
Energy	y ₇₀ n y ₆₈ C ₆₂ T ₉₂ T ₉₂ T ₇₀ C ₁₀₀ T ₁₀₀ T ₁₀₀ T ₅₉ n w ₆₈	AP-1
Cell growth	A ₈₄ A ₈₆ A ₁₀₀ C ₁₀₀ G ₁₀₀ C ₁₀₀ G ₁₀₀ T ₆₅	DSC1
Transcription	A ₆₂ G ₈₃ A ₉₇ A ₈₉ G ₇₄ A ₁₀₀ G ₁₀₀ G ₁₀₀	UME6
Protein synth.	T ₅₆ A ₆₂ C ₆₃ A ₈₅ T ₅₆ C ₉₈ C ₈₅ G ₆₃ T ₁₀₀ A ₁₀₀ C ₁₀₀ A ₁₀₀ T ₈₀ y ₇₆ T ₅₆	RAP1
Protein destin.	G ₁₀₀ C ₁₀₀ G ₁₀₀ A ₁₀₀ T ₁₀₀ G ₁₀₀ A ₁₀₀	BUF
Cell. organization	A ₆₅ A ₉₄ T ₁₀₀ A ₁₀₀ A ₈₉ T ₁₀₀ A ₁₀₀ T ₈₀ A ₅₉ G ₁₀₀ C ₁₀₀ T ₁₀₀ G ₁₀₀ C ₈₉ T ₆₇ G ₈₉ w ₇₆ C ₁₀₀ m ₈₈ G ₁₀₀ T ₈₆ G ₁₀₀ G ₁₀₀ A ₅₉ r ₆₅	MAL63 CAR1 DAL82
Protein kinases	A ₇₃ A ₁₀₀ A ₁₀₀ A ₈₇ A ₁₀₀ G ₁₀₀ A ₇₂ A ₈₃	UME6

Acknowledgments

This work was carried out as a part of ALIS (Advanced Lifescience Information System) project for genome analysis by Japan Science and Technology Corporation (JST). KN was supported in part by a Grant-in-Aid for Scientific Research on Priority Areas, "Genome Science", from The Ministry of Education, Science, Sports, and Culture in Japan.

References

- [1] Frech, K., Quandt, K., and Werner, T., "Software for the analysis of DNA sequence elements of transcription," *Comput. App. Biosci.*, 13:89–97, 1997.
- [2] Mewes, H. W., Albermann K, Bähr, M., Frishman, D., Gleissner, A., Hani, J., Heumann, K., Kleine, K., Maieryl, A., Oliver, S. G., Pfeiffer, F., and Zollner, A., "Overview of the yeast genome," *Nature*, 387(6632 Suppl):7–65, 1997.
- [3] Wingender, E., Dietze, P., Karas, H., and Knüppel, R., "TRANSFAC: a database on transcription factors and their DNA binding sites," *Nucl. Acids Res.*, 24:238–241, 1996.
- [4] Yada, T., Totoki, Y., Ishikawa, M., and Asai, K., "Automatic discovery of hidden markov representations for functional sites within DNA sequences," *Genome Informatics 1996*, pp.210–211, Univ. Acad. Press Inc., Tokyo, 1996.
- [5] Yada, T., Totoki, Y., Ishii, T., and Nakai, K., "Functional prediction of *B. subtilis* genes from their regulatory sequences," *Intell. Syst. Mol. Biol.*, 5:354–357, 1997.