*Genome analysis*

# Virtual Footprint and PRODORIC: an integrative framework for regulon prediction in prokaryotes

Richard Münch[1], Karsten Hiller[1], Andreas Grote[1,2], Maurice Scheer[1,3], Johannes Klein[1], Max Schobert[1] and Dieter Jahn[1,*]

[1]Institute of Microbiology, Technical University Braunschweig, Spielmannstrasse 7, D-38106 Braunschweig, Germany, [2]Institute of Biochemical Engineering, Technical University Braunschweig, Gaußstrasse 17, 38106 Braunschweig, Germany and [3]Department of Informatics, University of Applied Sciences Wolfenbüttel, Am Exer 2, 38302 Wolfenbüttel, Germany

## ABSTRACT

**Summary:** A new online framework for the accurate and integrative prediction of transcription factor binding sites (TFBSs) in prokaryotes was developed. The system consists of three interconnected modules: (1) The PRODORIC database as a comprehensive data source and extensive collection of TFBSs with corresponding position weight matrices. (2) The pattern matching tool Virtual Footprint for the prediction of genome based regulons and for the analysis of individual promoter regions. (3) The interactive genome browser GBPro for the visualization of TFBS search results in their genomic context and links to gene and regulator-specific information in PRODORIC. The aim of this service is to provide researchers a free and easy to use collection of interconnected tools in the field of molecular microbiology, infection and systems biology.

**Availability:** http://www.prodoric.de/vfp
**Contact:** d.jahn@tu-bs.de

**Table 1.** Statistics of the PRODORIC database content with respect to TFBSs, regulated genes, regulons, PWMs and promoters

| Organism | TFBSs | Genes | Regulons | PWMs | Promoters |
|---|---|---|---|---|---|
| *Bacillus subtilis* | 662 | 662 | 77 | 65 (47) | 410 |
| *Escherichia coli* | 1608 | 1015 | 81 | 84 (73) | 719 |
| *Helicobacter pylori* | 9 | 11 | 2 | 2 | 9 |
| *Listeria monocytogenes* | 11 | 12 | 3 | 2 | 12 |
| *Pseudomonas aeruginosa* | 176 | 241 | 32 | 21 (18) | 144 |
| *Streptococcus pyogenes* | 13 | 8 | 3 | 3 (2) | 6 |
| *Others* | 38 | 45 | 19 | 5 | 36 |
| *Sum* | 2517 | 1994 | 217 | 182 (149) | 1336 |

Status in August 2005.
In the case of several generated PWMs for one regulator, the non-redundant number is given in parentheses.

## INTRODUCTION

The accurate prediction of transcription factor binding sites (TFBSs) and whole regulons is still a crucial step towards the understanding of complex regulatory networks in systems biology. Current bioinformatic methods of pattern recognition usually suffer from their low specificity. This often results in an accumulation of false-positive matches (Frech *et al.*, 1997; Benitez-Bellon *et al.*, 2002). Therefore, we developed a new framework for the straightforward evaluation and visualization of *in silico* results with focus on bacterial gene regulation. The user can interactively identify putative TFBSs using genome wide searches and immediately obtain detailed information on the promoters, corresponding genes, operons and encoded proteins found. This includes the genomic localization and detailed information about potentially regulated genes via links to the prokaryotic database of gene regulation (PRODORIC) database as well as relevant links to external sources. Moreover, it is possible to evaluate the matches obtained according to their phylogenetic conservation. The software is organized in three major interconnected components which are the PRODORIC

database, the pattern search tool Virtual Footprint and the genome browser GBpro.

### The PRODORIC database

The PRODORIC database is a comprehensive source of prokaryotic genomes and their underlying gene regulatory networks (Münch *et al* 2003). Among many other features it contains a compilation of over 2500 TFBSs from several bacterial species including their interacting transcriptional regulators. The data of PRODORIC are all based on experimental evidence which were manually extracted from the original literature. Using this huge collection of TFBSs over 170 species-specific position weight matrices (PWMs) were generated which serve as a library for the pattern matching tool Virtual Footprint. Besides the exclusive TFBS information, PRODORIC contains data of genes and proteins, promoter details, operon structures and links to relevant databases. Table 1 summarizes the data content of PRODORIC in the field of bacterial gene regulation. More recently the database was supplemented with additional information such as expression data from trans criptomics and proteomics experiments and metabolic

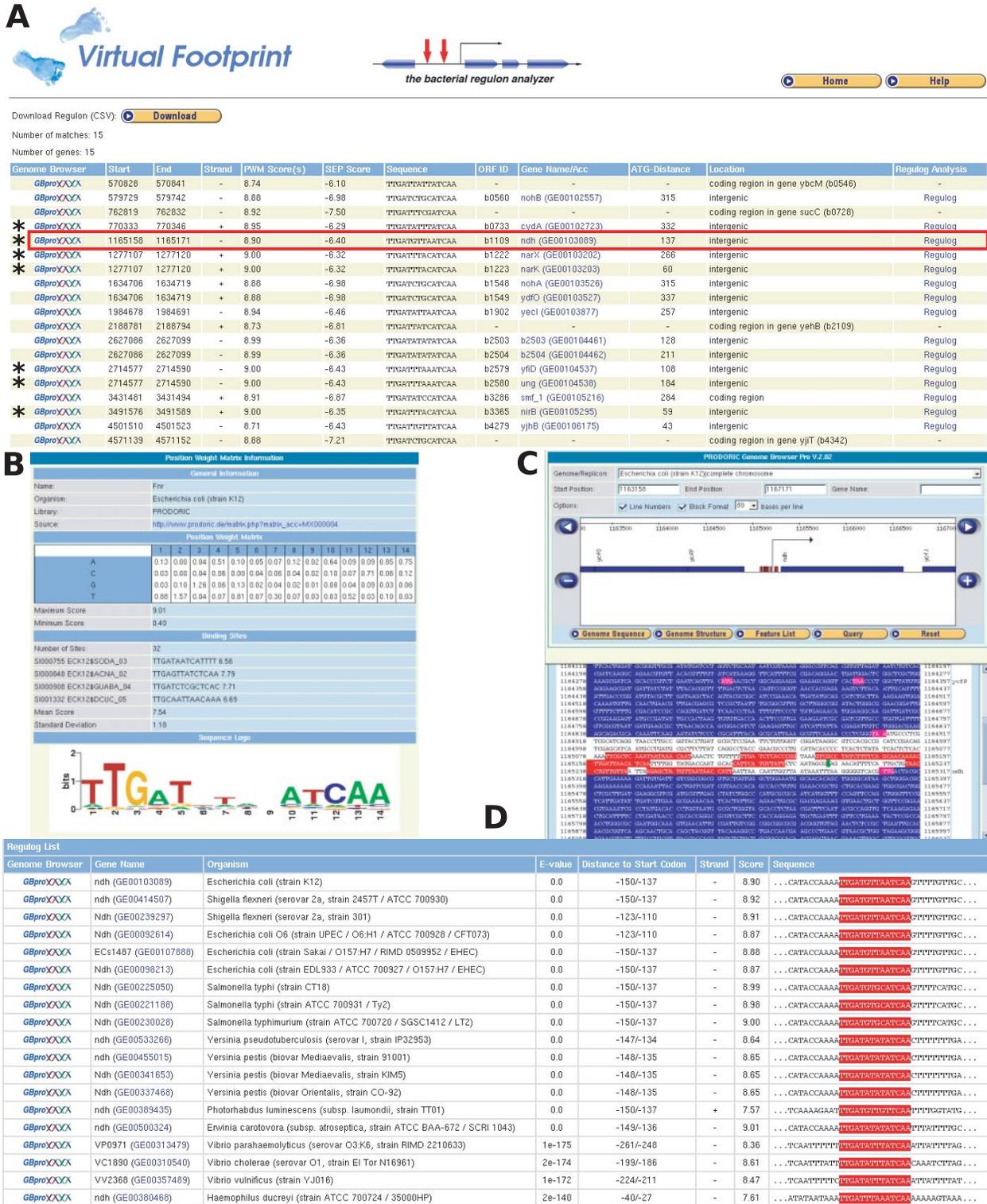*To whom correspondence should be addressed.

**Fig. 1.** Example for the application of Virtual Footprint for a genome wide prediction of Fnr (fumerate nitrate reduction regulatory protein) binding sites in *Escherichia coli*. (**A**) List of high scoring matches including their genomic positions, scores, assigned downstream genes, with links to PRODORIC, the genome browser GBpro and regulog analyses. Matches marked with asterisks are documented in the PRODORIC database. One of the matches upstream of the *ndh* gene (respiratory NADH dehydrogenase) is framed (**B**) Fnr PWM with sequence logo (the number of used regulatory sequences was reduced due to space limitations). (**C**) Genome browser view of the genomic context of *ndh*. (**D**) Regulog analysis of the Fnr binding site upstream of *ndh*. The resulting list provides links to the genome browser and to the gene entries of PRODORIC.

networks. This integrated approach makes PRODORIC well suited as a platform for systems biology in prokaryotes.

## Virtual Footprint—pattern matcher

The new pattern search tool Virtual Footprint offers fast searches of complex DNA patterns in whole bacterial genomes. Usually the search pattern is defined as PWM provided by PRODORIC (Fig. 1B). We also added PWMs from other resources (Robison *et al.*, 1998; Salgado *et al.*, 2004). However, in some cases, e.g. when sufficient sequence data are not available, it is necessary to use other pattern definitions (Stormo, 2000). Therefore, we implemented search algorithms using IUPAC consensi and regular expressions (Betel *et al.*, 2002). Some TFBSs are not only variable in their sequence conservation but also in their sequence lengths. Common examples are the occurrence of two half-sites separated by a variable spacer, the conserved −10 and −35 hexamers found in $\sigma^{70}$ regulated promoters and other so, called composite elements (Kel-Margoulis *et al.*, 2000). Therefore, Virtual Footprint allows the definition of bipartite patterns via the combination of up to two subpatterns separated by a variable spacer. Different pattern types can be freely combined e.g. a PWM with a IUPAC string. This enables a flexible definition of search patterns. The list of obtained matches can be evaluated by several different *in silico* approaches in accordance with their genomic context. Usually matches are directly linked to the downstream genes. Identified genes or operons are linked to the PRODORIC database and the genome browser GBpro (Fig. 1). If it is not possible to assign a gene to a match, the corresponding genomic location is specified. For a search the size of the upstream region (distance to the start codon) can be defined, the pattern orientation can be selected and matches in coding regions can be excluded. Matches can be further evaluated by analyzing upstream regions of orthologous genes for the same pattern (Fig. 1D). In this case orthologous sequence stretches from different genomes are extracted by the use of BLAST (Altschul *et al.*, 1990) and then analyzed via Virtual Footprint. This kind of investigation is also called regulog analysis (Alkema *et al.*, 2004). Furthermore, the GC-content of a promoter region and the resulting stacking energy plot shown by the genome browser can help to evaluate matches since functional targets should be localized in chromosomal regions where the GC-content and stacking energy are expected to be below the average. In addition to the whole genome search, Virtual Footprint allows the analysis of single promoter regions. In this case the upstream sequence of a gene or a pasted user defined sequence is compared with all PWMs provided by PRODORIC. Virtual Footprint has been successfully applied to define the ResD regulon of *Bacillus subtilis* (Härtig *et al.*, 2004) and to detect split tRNA genes in *Nanoarchaeum equitans* (Randau *et al.*, 2005). The program offers many options and settings not described here. A detailed description is available in the online help of the program.

## GBpro—Genome Browser

GBpro is a genome browser for an interactive navigation through all bacterial genomes available in PRODORIC. Genes, promoters and binding sites are displayed in parallel as graphical maps and highlighted sequences. Optionally, the GC-content and stacking energy of a DNA sequence of interest can be visualized. All results of Virtual Footprint are directly linked to GBpro and can thus be visualized in their genomic context (Fig. 1C). Similarly, genes and TFBSs present in PRODORIC are directly linked to this GBpro.

## REFERENCES

Alkema,W.B.L. *et al.* (2004) Regulog analysis: detection of conserved regulatory networks across bacteria: application to *Staphylococcus aureus. Genome Res.*, **14**, 1362–1373.

Altschul,S.F. *et al.* (1990) Basic local alignment search tool. *J. Mol. Biol.*, **215**, 403–410.

Baldi,P. and Baisnee,P.F. (2000) Sequence analysis by additive scales: DNA structure for sequences and repeats of all lengths. *Bioinformatics*, **16**, 865–889.

Benitez-Bellon,E. *et al.* (2002) Evaluation of thresholds for the detection of binding sites for regulatory proteins in *Escherichia coli* K12 DNA. *Genome Biol.*, **3**, 13.

Betel,D. and Hogue,C. (2002) Kangaroo—A pattern-matching program for biological sequences. *BMC Bioinformatics*, **3**, 20.

Frech,K. *et al.* (1997) Finding protein-binding sites in DNA sequences: the next generation. *Trends Biochem. Sci.*, **22**, 103–104.

Härtig,E. *et al.* (2004) *Bacillus subtilis* ResD induces expression of the potential regulatory genes yclJK upon oxygen limitation. *J. Bacteriol.*, **186**, 6477–6484.

Kel-Margoulis,O.V. *et al.* (2000) COMPEL: a database on composite regulatory elements providing combinatorial transcriptional regulation. *Nucleic Acids Res.*, **28**, 311–315.

Münch,R. *et al.* (2003) PRODORIC: prokaryotic database of gene regulation. *Nucleic Acids Res.*, **31**, 266–269.

Randau,L. *et al.* (2005) *Nanoarchaeum equitans* creates functional tRNA from separate genes for their 5′- and 3′-halves. *Nature*, **433**, 537–541.

Robison,K. *et al.* (1998) A comprehensive library of DNA-binding site matrices for 55 proteins applied to the complete *Escherichia coli* K12 genome. *J. Mol. Biol.*, **284**, 241–254.

Salgado,H. *et al.* (2004) RegulonDB (version 4.0): transcriptional regulation, operon organization and growth conditions in *Escherichia coli* K12. *Nucleic Acids Res.*, **32**, 303–306.

Stormo,G.D. (2000) DNA binding sites: representation and discovery. *Bioinformatics*, **16**, 16–23.