

VFDB 2008 release: an enhanced web-based resource for comparative pathogenomics

Jian Yang¹, Lihong Chen^{1,2}, Lilian Sun¹, Jun Yu³ and Qi Jin^{1,2,*}

¹State Key Laboratory for Molecular Virology and Genetic Engineering, National Institute for Viral Disease Control and Prevention, Chinese Center for Disease Control and Prevention, Beijing 100176, ²Institute of Pathogen Biology, Chinese Academy of Medical Sciences, Beijing 100730, China and ³Strathclyde Institute of Pharmacy and Biomedical Sciences, University of Strathclyde, Royal College, Glasgow, G1 1XW, UK

Received September 12, 2007; Revised October 11, 2007; Accepted October 15, 2007

ABSTRACT

Virulence factor database (VFDB) was set up in 2004 dedicated for providing current knowledge of virulence factors (VFs) from various medical significant bacterial pathogens to facilitate pathogenic research. Nowadays, complete genome sequences of almost all the major pathogenic microbes have been determined, which makes comparative genomics a powerful approach for uncovering novel virulence determinants and hidden aspects of pathogenesis. VFDB was therefore upgraded to present the enormous diversity of bacterial genomes in terms of virulence genes and their organization. The VFDB 2008 release includes the following new features; (i) detailed tabular comparison of virulence composition of a given genome with other genomes of the same genus, (ii) multiple alignments and statistical analysis of homologous VFs and (iii) graphical comparison of genomic organizations of virulence genes. Comparative analysis of the numerous VFs will improve our understanding of the nature and evolution of virulence, as well as the development of new therapeutic and preventive strategies. VFDB 2008 release offers more user-friendly tools for comparative pathogenomics and it is publicly accessible at <http://www.mgc.ac.cn/VFs/>.

INTRODUCTION

Infectious diseases remain to be one of the biggest threats to public health despite the advance of modern medicine in post-genome era (1). Virulence factors (VFs) refer to the traits encoded by 'virulence genes' that pathogenic microbes are equipped to cause infection.

To combat infectious diseases, a better understanding of VFs is absolutely necessary to decipher the mechanisms pathogenic microbes employ. VFDB was built to meet the challenge of providing up-to-date knowledge about VFs from various medically important bacterial pathogens (2).

The term pathogenomics is given to describe genomic approaches in studying microbial pathogens as to how they interact with their hosts, and in other words, pathogenomics is the study of pathogenic microbes and the entities they infect on the genomic level. The availability of complete genome sequences of different microbial species enables comparative studies to identify the common as well as species- or strain-specific VFs. Pathogenic bacteria have acquired various VFs that allow them to colonize diverse niches, cause infection and to survive in the hosts. Commonly shared VFs indicate universal requirement to cause infection by related pathogens, whereas narrowly distributed VFs often determine species- and/or strain-specific characteristics.

As a consequence, comparative genomic approaches were introduced into VFDB to explore VFs within completely deciphered bacterial genomes. The VFDB 2008 release has not only collected up-to-date knowledge about VFs from over 200 complete genomes of pathogenic bacteria, but also has incorporated a set of analytical tools to meet the desire of comparative pathogenic studies.

DATABASE UPDATES

Data source and construction for comparative analysis

Information of publicly available bacterial genomes was retrieved from the summary page of 'Complete Microbial Genomes' at NCBI (<http://www.ncbi.nlm.nih.gov/genomes/lproks.cgi>). RefSeq is a curated non-redundant collection of sequences with uniform format (3).

*To whom correspondence should be addressed. Tel: +86 10 6787 7732; Fax: +86 10 6787 7736; Email: zdsys@sina.com

The authors wish it to be known that, in their opinion, the first two authors should be regarded as joint First Authors.

© 2007 The Author(s)

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/2.0/uk/>) which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.

For convenience of later data processing only genomes that are available from RefSeq database were included for further comparative analysis (both pathogenic and non-pathogenic isolates). The complete genome sequences and annotations were batch downloaded from the FTP server of RefSeq (<ftp://ftp.ncbi.nih.gov/genomes/Bacteria/>).

The VF loci in each genome were obtained from the original literatures and subsequent reviews. Each of the VF genes was verified by sequence-similarity search against the genomes of related bacterial pathogens. Data collected were manually inspected and each homolog group was further validated by multi-alignments (see below).

The NCBI BLAST software was used for local sequence-similarity search (4). A series of BioPerl scripts were designed to extract features for all desired loci from the downloaded genome files in a semi-automated fashion. An enhanced multiple genome map viewer (5) was employed for graphical comparison of the pathogenic organizations (see below). ClustalW (6) was run in batch by an in-house Perl script to generate multi-alignments for each homolog group.

Full tabular comparison of the pathogenic composition

Tabular style was commonly used in scientific literatures for comparative analysis. For each genus a full comparison of pathogenic composition is given as a spreadsheet to integrate information about VFs and genomes (see Figure 1A for example). The far left column organizes all known VFs in functional groups (toxins, lipase, etc.) and the next column lists all VF genes, and each row gives gene IDs (i.e. 'locus_tag' in annotation files) of the respective genomes, and pseudogenes are highlighted by star marks. Each gene ID in the table is a hyperlink that connects its individual page for full DNA and protein sequences. All tables can be downloaded as Excel files by terminal users.

For a quick glance of overall information, the above full-detailed tables can be converted to simplified tables, where each VF occupies a single row and gene details are replaced by symbols; '+' for presence, '-' for absence, and '±' for partial or non-functional genes. The simplified table can be viewed in text mode, symbolic mode or schematic mode.

In some genus there are many complete genome sequences, and an extreme case is *Streptococcus* with 25 complete genomes up to date. Taking the consideration that a selective subset of genomes might be more interesting to certain studies, a special filter is designed to allow generating tables that contain only selected genomes. For example, a table can be generated by expressing data of only 12 pyogenic genomes within the genus of *Streptococcus*.

Multiple alignments of homologous virulence genes

Analysis of homologous genes is a powerful approach for elucidating gene structure, function and evolution.

The diversity of nucleotide sequences of bacterial genes often reflects particular niches a microbe colonizes in vivo and in the environment (7). From the full-comparison table described above, a single click on gene name will return a page with multi-alignment of both nucleotide and amino acid sequences if homologous gene(s) existing. A summary table on top of the alignment gives statistics about unmatched overhang, length of the alignment and percentage of polymorphic sites. A configurable phylogenetic tree constructed by ClustalW is displayed beneath the alignment when more than two sequences are involved.

A filter is also designed to perform multi-alignment on selected sequences only. In this case however there is no pre-computed result by default, and an online ClustalW must be run which constructs an alignment in a few seconds.

Concise graphical comparison of pathogenic organization

Bacterial genome evolution has been driven by nucleotide substitutions and indels, as well as the changes of the genome architecture by genetic rearrangements including translocations and inversions (8). Recent comparative genomic studies have revealed that the dynamic changes of genome structures contribute greatly to the adaptive evolution of certain bacterial pathogens, such as *Shigella* (9). To unambiguously display the dynamic features of the genomes and to compare VFs' genomic organization among related pathogens, an enhanced multiple genome map viewer was implemented, which depicts all VF genes as clickable arrows (or bars) and color-coded by functional classifications. Since details about genes unrelated to virulence are hidden, the map becomes concise, although not to scale, and suitable for quick examinations of the genome organization of VFs among related genomes.

The viewer page provides three different representing styles: (i) complete mode which exhibits full scale pathogenic map that is informative but usually large in size; (ii) compact mode that provides details of all virulence loci but omits flanking genes/regions; (iii) overview mode that scales the map to fit full screen without giving details. To facilitate interpretation of pathogenic synteny under the overview mode, there are lines to connect homologous VF genes of the adjacent genomes when only one replicon is available (or selected by the users) for each genome. Terminal users can also run the viewer with select genomes of their interest. The usefulness of displaying synteny is highlighted in the case of *Listeria* species (Figure 1B); their genomes exhibit a high synteny in virulence gene organization. It is in agreement with the recent listerial pangenome studies, which revealed the lack of inversions or shifting of large genome segments in the sequenced *Listeria* genomes. The possible reason may be the low occurrence of transposons and insertion sequence elements in those genomes (10).

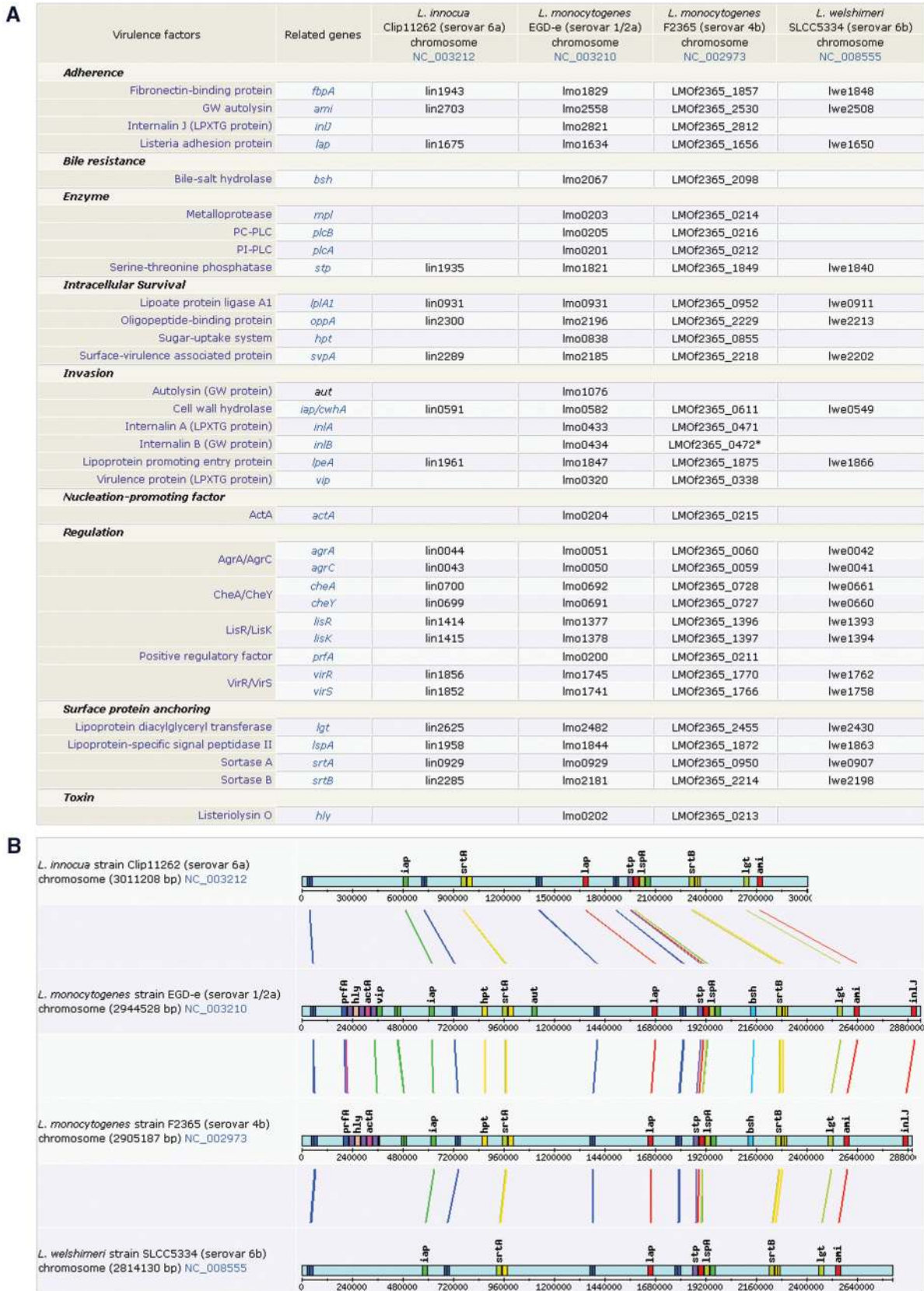


Figure 1. Comparative pathogenomic results of four sequenced *Listeria* genomes. (A) Full tabular comparison of pathogenomic composition. (B) Graphical overview for the comparison of pathogenomic organization. VF genes are color-coded by their functional classifications. Homologues between each adjacent pair of genomes are indicated by connecting lines for convenience of further interpretation.

DISCUSSION

Virulence involves a wide spectrum of biological activities, which is reflected by the diverse VFs employed by pathogenic microbes to colonize the particular niches in the hosts. A fuller investigation of VFs is highly desirable for pathogenomic research. VFDB 2008 release attempts to meet such a challenge by providing all creditable information up to date and by providing more analytical tools to the terminal users. The comparative pathogenomic results indicate that most pathogens have a flexible gene pool encoding VFs. Different combinations of VFs or organizations on microbial genomes or different expression patterns of VFs may in consequence be responsible for the diverse clinical signs of pathogen infections.

VFDB 2008 release has expanded with additional eight pathogens, which are *Brucella*, *Bartonella*, *Campylobacter*, *Clostridium*, *Corynebacterium* and *Enterococcus*, as well as *Chlamydia* and *Mycoplasma*. VFDB will continue to expand by including more medical significant pathogens, and provide up-to-date information by regular updates. For the convenience of local use, full dataset of VFDB is available for batch download in several forms, including FASTA sequences and tabular (Excel) files. Furthermore, new features and analytical tools are under development which we anticipate to make VFDB a useful pathogenomic resource to the scientific community.

ACKNOWLEDGEMENTS

This work is supported by the National Basic Research Program from the Ministry of Science and Technology of China (grant No 2005CB522904), and National Natural Science Fund from the National Natural Science Foundation of China (grant No 30600022).

Funding to pay the Open Access publication charges for this article was provided by the National Natural Science Foundation of China after partially waived by Oxford University Press.

Conflict of interest statement. None declared.

REFERENCES

1. Becker, K., Hu, Y. and Biller-Andorno, N. (2006) Infectious diseases—a global challenge. *Int. J. Med. Microbiol.*, **296**, 179–185.
2. Chen, L., Yang, J., Yu, J., Yao, Z., Sun, L., Shen, Y. and Jin, Q. (2005) VFDB: a reference database for bacterial virulence factors. *Nucleic Acids Res.*, **33**, D325–D328.
3. Pruitt, K.D., Tatusova, T. and Maglott, D.R. (2007) NCBI reference sequences (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res.*, **35**, D61–D65.
4. Altschul, S.F., Madden, T.L., Schaffer, A.A., Zhang, J., Zhang, Z., Miller, W. and Lipman, D.J. (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.*, **25**, 3389–3402.
5. Yang, J., Chen, L., Yu, J., Sun, L. and Jin, Q. (2006) ShiBASE: an integrated database for comparative genomics of *Shigella*. *Nucleic Acids Res.*, **34**, D398–D401.
6. Thompson, J.D., Higgins, D.G. and Gibson, T.J. (1994) CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res.*, **22**, 4673–4680.
7. Gilsdorf, J.R., Marrs, C.F. and Foxman, B. (2004) *Haemophilus influenzae*: genetic variability and natural selection to identify virulence factors. *Infect. Immun.*, **72**, 2457–2461.
8. Mira, A., Klasson, L. and Andersson, S.G. (2002) Microbial genome evolution: sources of variability. *Curr. Opin. Microbiol.*, **5**, 506–512.
9. Yang, F., Yang, J., Zhang, X., Chen, L., Jiang, Y., Yan, Y., Tang, X., Wang, J., Xiong, Z. *et al.* (2005) Genome dynamics and diversity of *Shigella* species, the etiologic agents of bacillary dysentery. *Nucleic Acids Res.*, **33**, 6445–6458.
10. Hain, T., Chatterjee, S.S., Ghai, R., Kuenne, C.T., Billion, A., Steinweg, C., Domann, E., Karst, U., Jansch, L. *et al.* (2007) Pathogenomics of *Listeria* spp. *Int. J. Med. Microbiol.*, **297**, 541–557.