epcALGO: a home-grown algorithm for entire proteome comparison

1 2

4

5

6

7

8

9

10

11 12

13

14

15

16

17

Abstract 3

Due to the advancement of bioinformatics and genome sequencing project, entire genome and proteome sequences of different organisms are available in the public domain. These vast data are repeatedly compared and explored to find out identical and similar sequence patterns. In this paper we employed NCBI's Standalone BLAST program for entire proteome comparison of any two strains / species and illustrate a simple algorithm for the same. The implementation of this epcALGO algorithm is to identify systematically conserved proteins that are missing in a given proteome and also identify proteins unique to a particular species. This algorithm is simple and quick to apply for revealing the species / strain variation among any two closely related species / strains by identifying identical and non-identical proteins in their proteomes and also identifying where there is mutation in the protein sequence. We implemented this algorithm for proteome comparison of two strains of Mycobacterium tuberculosis H₃₇Rv and H₃₇Ra and elucidated the methodology for finding out their proteomic variation.

1	Satish Kumar ^{1*} , Lingaraja Jena ^{1,2}		
2 3 4	 Bioinformatics Centre & Biochemistry, Mahatma Gandhi Institute of Medical Sciences, Sevagram, Maharashtra, India Department of Bioinformatics, Shri JJT University, Jhunjhunu, Rajasthan 		
5			
6	Email address:		
7	SK: satishangral@gmail.com		
8	LJ: lingaraj.jena@gmail.com		
9			
10			
11	*Corresponding Author:		
12	Dr. Satish Kumar		
13	Professor, Biochemistry &		
14	Dy Coordinator, Bioinformatics Centre		
15	Mahatma Gandhi Institute of Medical Sciences		
16	Sevagram-442 102 (Wardha) Maharashtra, India		
17	Tel +91 7152 284679		
18	Fax +91 7152 284038		
19	Email: satishangral@gmail.com		
20			
21			
22			

Introduction

- 2 Due to the rapid development of new sequencing technologies (Li and Homer, 2010) and
- hasty progress in bioinformatics, the complete genome and proteome sequences of numerous 3
- organisms have become available in the public databases. As of September 2012, Genomes 4
- Online Database (GOLD) version 4.0 (Pagani et al., 2012), contains information on 3699 5
- complete sequencing projects of different organisms (http://www.genomesonline.org). In 6
- 7 order to retrieve and analyze those huge amount of sequence data, development of novel
- algorithms and techniques are becoming increasingly important. 8
- Comparison is an essential feature of all biological research and in early days, the 9
- 10 comparisons were revolved around morphological and physiological level of research due to
- 11 unavailability of DNA and proteins sequences (Bachhawat, 2006). In recent era of Genomics
- & Proteomics the comparison paradigm has been shifted to entire genome & proteome level 12
- since complete genome & proteome sequences of several organisms have become freely 13
- 14 available in the public databases. Comparative genomics / proteomics among different
- pathogenic organisms can not only reveal evolutionary history among them but also identify 15
- 16 the similarity and variation in their DNA / protein sequences, which accounts for their
- morphological & physiological changes as well. Besides, such analysis is of great value in 17
- 18 phylogenetic reconstruction, drug discovery and functional annotation of hypothetical
- proteins (Bachhawat, 2006). 19
- Many tools have been developed for complete determination of genome sequence of a huge 20
- 21 number of bacteria, but still, their proteomes remain relatively poorly defined. In the post
- genomic era, proteomics is a rapidly growing field of research for studying proteins involved 22
- 23 in carcinogenesis as well as novel biomarker discovery for clinical use such as screening,
- diagnosis, prognosis, detection of recurrent disease etc (Cho, 2007). Since proteins are 24
- specifically directly involved in both normal and disease related biochemical processes, a 25
- more comprehensive understanding of disease may be achieved by looking directly into the 26
- proteins within a disease cell or tissue (Cho, 2007). Proteomics has much promise in novel 27
- drug discovery by targeting proteins of pathogenic organisms causing different diseases in 28
- host, whereas comparative proteomics is very significant in studying the proteomic variations 29
- 30 among different pathogens.
- The simultaneous development of rapid sequence comparison algorithms has revolutionized 31
- 32 the role of biological sequence comparison in molecular biology. The Basic Local Alignment

- 1 Search Tool (BLAST) (Altschul et al., 1990) is now prevails over as the fastest and most
- 2 widely-used tool for sequence similarity searches (McGinnis and Madden, 2004). The stand-
- 3 alone executable BLAST from the NCBI BLAST site
- 4 (ftp://ftp.ncbi.nlm.nih.gov/blast/executables/release/LATEST/) provide easy ways for a user
- 5 to perform blast searches via command line or a local web server. However, the standalone
- 6 version requires each user to install and configure the program and customize the databases
- 7 for performing specialized research (Deng et al., 2007).
- 8 In this study we have developed a simple algorithm for entire proteome comparison of any
- 9 two organisms (strains / species) using NCBI standalone BLAST and discussed on how to
- 10 retrieve and analyze those comparison data. We also implemented the algorithm for
- 11 comparison of virulent (MTB H₃₇Rv) and avirulent (MTB H₃₇Ra) strains of *Mycobacterium*
- 12 tuberculosis (MTB). Our observations provide a unique platform for discovery of proteomic
- variation in different strains / species.

Materials and Methods

15 The methodology for entire proteome comparison and analysis involves following simple

steps using EpcALGO algorithm (Fig. 1).

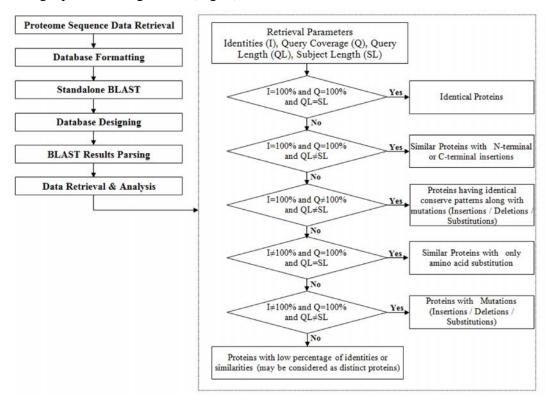


Figure 1: epcALGO Algorithm.

18

Data Retrieval 1

- Dataset was prepared by retrieving the entire proteome of MTB H₃₇Rv (NCBI RefSeq: 2
- 3 NC_000962.3) & MTB H₃₇Ra (NCBI RefSeq: NC_009525.1) from NCBI FTP site
- (ftp://ftp.ncbi.nih.gov/genomes/). 4

Database formatting 5

- Protein sequence data of H₃₇Rv (NC_000962.faa) & H₃₇Ra (NC_009525.faa) were formatted 6
- using formatdb application of NCBI Standalone BLAST-2.2.26. 7

Perform Standalone BLAST 8

- Blastall application was used to perform protein BLAST (Altschul et al., 1990) between the
- proteome of MTB H₃₇Rv against MTB H₃₇Ra to carry out proteomic comparison. 10

11 **Database Designing**

- Database table was created using Micosoft SQL (MS SQL) Server for storing proteome 12
- 13 comparison information of MTB H₃₇Rv vs MTB H₃₇Ra. For database designing other
- database server such as MySQL can also be used. 14

Parsing BLAST results 15

- 16 The output of the BLAST result was parsed and stored in MS SQL relational database tables
- using in-house developed PERL script using Bio::SearchIO module of Bioperl (Stajich et al., 17
- 2002). While parsing BLAST output results, percentage identities, positivities, number of 18
- gaps, identical residues, bits, bits score, e-value, query length, subject length, query sequence, 19
- subject sequence, consensus sequence etc of the first hit obtained were taken into 20
- consideration for each protein comparison. 21

22 Data retrieval & analysis

- Different SQL queries were written to retrieve and analyze comparison data from MS SQL 23
- database tables. 24

25

26 **Results and Discussion**

- In spite of several studies in the past the potential causes for variation in virulence between 27
- 28 MTB H₃₇Rv and MTB H₃₇Ra have remained unclear. A single amino acid mutation in protein
- sequence may cause alteration in protein structure and function that may account for 29

18

- 1 virulence and drug resistance properties of pathogenic organisms. Therefore, the development
- 2 of an *in silico* technology to study the proteomic variations of different strains of genetically
- 3 intractable pathogens such as MTB will enhance the analysis of virulence and drug resistance
- 4 properties and significantly advance the understanding of the mechanisms of disease.
- 5 In our previous study we have implemented this algorithm and developed *Mycobacterium*
- 6 tuberculosis proteome comparison database (MTB-PCDB) which provides integrated access
- 7 to protein sequence comparison with identical and non identical protein data for five strains
- 8 of MTB (H₃₇Rv, H₃₇Ra, CDC 1551, F11 and KZN 1435) (Jena et al., 2011). In this study, we
- 9 implemented epcALGO algorithm to performed comparative proteomic analysis of MTB
- 10 H₃₇Rv and MTB H₃₇Ra. While comparison, protein sequence of MTB H₃₇Rv was taken as
- query and sequences of MTB H_{37} Ra were taken as database sequences (subject).
- 12 A total of 4018 protein-coding sequences (CDS) are identified amongst 4111 genes in the
- MTB H_{37} Rv genome while there are 4084 genes with 4034 protein coding sequences in the
- 14 genome of MTB H₃₇Ra. There were seven categories (Table 1) obtained depending on the
- 15 percentage identities, query coverage, query and subject length upon entire proteome
- 16 comparison between these two strains.

Table 1. Comparison of proteomic variations between MTB H₃₇Rv and MTB H₃₇Ra

Category	Features	Total
		Number
1	Identical Proteins in both the strains	3804
2	Proteins having 100% identities and query coverage but with varying sequence length	20
3	Proteins having 100% identities but with variation in query coverage	21
4	Proteins having same sequence length and 100% query coverage but with variation in identities	36
5	Proteins having 100% query coverage but variation in length and identities	31
6	Proteins having variation in length, identities and query coverage	101
7	Proteins with no significant similarities	5

- 1 Category 1: Identical Proteins in both the strains (MTB H₃₇Rv and MTB H₃₇Ra)
- 2 Proteomic comparison of M. tuberculosis H₃₇Rv and M. tuberculosis H₃₇Ra in our study
- 3 revealed 3804 identical proteins between these two strains. A protein of MTB H₃₇Rv is said
- 4 to be identical compared to the corresponding protein of MTB H₃₇Ra when both identities
- 5 and query coverage are 100% and the query length is equal to the subject length. So, 214
- 6 proteins were identified as non identical proteins between these two strains.
- 7 Category 2: Proteins having 100% identities and query coverage but with varying
- 8 sequence length
- 9 There were 20 proteins identified in this category. In this case the length of subject (MTB
- H_{37} Ra) sequences was found to be greater in comparison to sequences of MTB H_{37} Rv. This
- observation revealed that there were insertions in the respective proteins of MTB H_{37} Ra.
- 12 Category 3: Proteins having 100% identities but with variation in query coverage
- 13 21 proteins of MTB H_{37} Rv were found where identities was 100% but with variation in query
- 14 coverage (less than 100%) when compared with proteins of MTB H₃₇Ra. This revealed that
- there were mutations (insertions / deletions / substitutions) in these sequences.
- 16 Category 4: Proteins having same sequence length and 100% query coverage but with
- 17 variation in identities
- Out of 214 non identical proteins, 36 proteins of MTB H_{37} Ra were identified with only amino
- 19 acid substitution compared to corresponding proteins of MTB H₃₇Rv, as there was only
- variation in identities with same sequence length and 100% query coverage.
- 21 Category 5: Proteins having 100% query coverage but variation in length and identities
- 22 31 proteins were identified in this category. In this case BLAST results showing 100% query
- coverage with varying length and identities revealed that there were insertions / deletions /
- 24 substitutions in these sequences.
- 25 Category 6: Proteins having variation in length, identities and query coverage
- There were 101 proteins observed with variation in length, identities and query coverage. So,
- 27 most of the dissimilar proteins were found in this category. As many proteins in this group
- 28 were found to have lower percentage of similarities and identities, these proteins may be
- 29 considered as distinct proteins in MTB H₃₇Ra, which needs further study.
- 30 Category 7: Proteins with no significant similarities
- 31 There were five proteins of MTB H₃₇Rv identified with no significant similarities compared
- 32 to the proteins of MTB H_{37} Ra. So, these proteins were unique to MTB H_{37} Rv.

22

- 1 Overall comparative analysis provide proteomic differences between MTB H₃₇Rv and H₃₇Ra,
- 2 which may be useful for better understanding of the basis of pathogenesis of *Mycobacterium*
- 3 tuberculosis and virulence attenuation in MTB H₃₇Ra (Jena et al., 2013). Further studies on
- 4 functional characterization of non-identical proteins identified in MTB H₃₇Ra and MTB
- $_{5}$ $H_{37}Rv$ may help in understanding the important role of variation among these two strains.

6 Conclusion

- 7 epcALGO algorithm has been successfully implemented in designing MTB-PCDB (Jena et
- 8 al., 2011). This algorithm is designed especially for protein sequence comparison using
- 9 'blastall' application of standalone BLAST, however the same method can also be applicable
- 10 for DNA sequence comparison using nucleotide blast. As genome and proteome sequences of
- 11 different organisms are available and easily accessible to researchers, this simple and
- proficient algorithm developed in this study, would be helpful in finding proteomic variations
- in different strains / species and subsequently reveals the morphological dissimilarities
- 14 amongst them.

Acknowledgments

- 16 Authors express gratitude to Dr. B.C. Harinath, Director, JBTDRC & Coordinator,
- 17 Bioinformatics Centre for providing Bioinformatics Laboratory wherein this study has been
- carried out. Grateful thanks to Shri D.S. Mehta, President, Kasturba Health Society, Dr.
- 19 (Mrs.) P. Narang, Secretary, Kasturba Health Society, Dr. B.S. Garg, Dean, MGIMS & Dr.
- 20 S.P. Kalantri, Medical Superintendent, Kasturba Hospital, MGIMS, Sevagram for their
- 21 encouragement & unconditional support.

References

1

- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *Journal of Molecular Biology*, 215(3):403-410.
- 4 Bachhawat AK. 2006. Comparative genomics. *Resonance*, 11(8):22-40.
- 5 Cho WC. 2007. Contribution of oncoproteomics to cancer biomarker discovery. *Molecular* 6 *Cancer*, 6:25.
- 7 Cho WC. 2007. Proteomics technologies and challenges. Genomics Proteomics. 8 *Bioinformatics*, 5(2):77-85.
- Deng W, Nickle DC, Learn GH, Maust B, Mullins JI. 2007. ViroBLAST: a stand-alone BLAST web server for flexible queries of multiple databases and user's datasets. *Bioinformatics*, 23(17):2334-2336.
- Jena L, Kashikar S, Kumar S, Harinath BC. 2013. Comparative proteomic analysis of *Mycobacterium tuberculosis* strain H₃₇Rv versus H₃₇Ra. *International Journal of Mycobacteriology*, 2(4):220-226.
- Jena L, Wankhade G, Kumar S, Harinath BC. 2011. MTB-PCDB: *Mycobacterium* tuberculosis proteome comparison database. *Bioinformation*, 6(3):131-133.
- Li H, Homer N. 2010. A survey of sequence alignment algorithms for next-generation sequencing. *Briefings in Bioinformatics*, 11(5):473-483.
- McGinnis S, Madden TL. 2004. BLAST: at the core of a powerful and diverse set of sequence analysis tools. *Nucleic Acids Research*, 32:W20-W25.
- Pagani I, Liolios K, Jansson J, Chen IM, Smirnova T, Nosrat B, Markowitz VM, Kyrpides NC. 2012. The Genomes OnLine Database (GOLD) v.4: status of genomic and metagenomic projects and their associated metadata. *Nucleic Acids Research*, 40:D571-579.
- Stajich JE, Block D, Boulez K, Brenner SE, Chervitz SA, Dagdigian C, Fuellen G, Gilbert
 JG, Korf I, Lapp H, Lehväslaiho H, Matsalla C, Mungall CJ, Osborne BI, Pocock
 MR, Schattner P, Senger M, Stein LD, Stupka E, Wilkinson MD, Birney E. 2002.
 The Bioperl toolkit: Perl modules for the life sciences. *Genome Research*,

29 12(10):1611-1618.