# Comparisons Between Different Methods in Measuring Enzyme Similarity for Metabolic Network Alignment

**Wenwei Zhou[1]; Jiangtao Liu[1]; Hiu Xing[1]; Zhengdong Zhang[1]; Xiaoyao Xie[1]; Fengxuan Jing[1]\***
[1] *Guizhou normal University, Key Laboratory of Information and Computing Science of Guizhou Province, Guiyang, Guizhou, China*

## ABSTRACT

*Metabolic network alignments enable comparison of the similarities and differences between pathways in two metabolic networks and help to uncover the conserved sub-blocks therein. Such analysis is important in the understanding of metabolic networks and species evolution. The fundamental parts of metabolic network alignment algorithms all involve comparisons of the similarity between two enzymes as a similarity measure of network nodes. As a result, the study of methods for measuring enzyme similarity becomes highly relevant. Currently, two approaches are mainly used to measure enzyme similarity. One of the methods is based on similarity measures of gene or protein sequences; the other is based on enzyme classification. In this study, multiple metabolic network alignments were performed using both the methods. The results showed that, in general, the sequence similarity method yielded higher accuracy, especially with respect to reflecting evolutionary distances.*

**Key words**: metabolic network alignment, enzyme, similarity measure, classification numbering, sequence similarity

**\***Author for correspondence: 937830175@qq.com

## INTRODUCTION

With the advent of the post-genomic era, a major challenge faced by biologists is to understand how biomolecules such as proteins and enzymes interact with each other to fulfill various functions in organisms, for example, in gene regulatory networks, protein interaction networks, and metabolic networks (Li et al. 2008; Kuchaiev et al. 2010). The differences between these interactive networks of biomolecules lie not only in their composition but also in their topology (Doncheva et al. 2012). Therefore, the study of the differences between various biological networks by using biological network alignment becomes highly relevant. Biological network alignment is a crucial approach for understanding the structures, functions, and evolution of organisms (Sharan et al. 2006). Similarities and differences between network topologies can be determined by performing comparative analysis of the biological networks of different species; this enables further investigations of the conserved regions in a network for uncovering new biological functions and for understanding the relationships between the structures and functions of various molecules (Yamada et al. 2009).

A metabolic network involves multiple in vivo chemical reactions. Individual intracellular biochemical reactions can be studied at the level of the whole network. Therefore, metabolic networks comprise the final output that integrates the interaction of multiple types of data (i.e., data from genes, proteins, and metabolites).

In metabolic network alignments, enzymes that catalyze reactions can be seen as nodes in a pathway diagram, while the substrates and products can be seen as edges connecting these enzymes in the diagram. From this perspective, metabolic network alignment is used to compare the similarity of nodes and edges in pathway diagrams. As a result, matches between enzymes in a network model not only help determine the differences between nodes but also significantly influence the overall network alignment score.

Currently, two approaches are available for measuring enzyme similarity. One is based on the enzyme classification proposed by Tohsato et al. (2000). The method utilizes the Enzyme Commission (EC) classification numbering to quantify the differences between enzymes in two metabolic networks. The other method (Altschul et al. 1997; Francke et al. 2005; Amir-Ghiasvand et al. 2014) refers to the common practice of protein interaction network alignment. Analysis of the genes encoding enzymes in metabolic networks helps to determine the differences between enzymes in two different metabolic networks based on the similarity in their DNA and amino acid sequences.

In this study, we compared the two metabolic network alignment methods described above and obtained the following results: In general, as the sequence similarity method takes into account the topological structure and sequence characteristics of the enzymes, it has higher accuracy, especially with respect to reflecting phylogenetic/evolutionary relationships.

## MATERIAL AND METHODS

### Data Source

All the metabolic network models used in this study were obtained from the Kyoto Encyclopedia of Genes and Genomes (KEGG) Metabolic PATHWAY Database of Kyoto University. The specific steps were as follows: the "Amino sugar and nucleotide sugar metabolism" pathway under KEGG PATHWAY in the KEGG website was selected. Next, under the Organism menu, the following two classes of bacteria were selected: *Gammaproteobacteria - Enterobacteria* and *Betaproteobacteria*; under each class, the desired species and their metabolic networks were selected.

Using the traversal method, all the enzymes in the metabolic network of each species, the sequence similarity measures of each enzyme pair, and the similarity measures of the EC numbers between enzymes were obtained.

### Metabolic Network Alignment Algorithms and Selection of the Enzymes' Similarity Measures

Let a network be represented by a graph $G(V, E)$, where each node $v \in V$ denotes the enzyme node in the network and each edge $(u,v) \in E$ denotes the lines connecting enzyme u and enzyme v. Given two metabolic networks $G1=(V1,E1)$ and $G2=(V2, E2)$, the one-to-one mapping between a subset of V1, denoted as $V1'$, and a subset of V2, denoted as $V2'$, is defined as a network alignment.

In this study, the alignment method used is a modification of the method proposed by Li et al. (Li et al. 2007). This method determines the alignment measure as a constraint-based non-integer quadratic programming problem in order to find the node-matching and edge-matching mappings between two metabolic networks, i.e., the optimal state of matching these two metabolic networks. In this method, $S_{ij}$ denotes the similarity measure of two graphs, where $i \in G1$ and $j \in G2$. As mentioned above, there are two ways to determine $S_{ij}$:

Method 1: This method involves obtaining the EC numbers of two different enzymes from KEGG and finding the lowest class in the hierarchy shared by the two EC numbers. For example, considering Enzyme (1.1.1.1) and Enzyme (1.1.1.2), the lowest class in the hierarchy is (1.1.1.-). The enzyme similarity is then defined as the inverse of all the enzyme numbers below the lowest class in the hierarchy.

Method 2: This method involves obtaining the protein sequences of the enzymes of a given species from KEGG and determining the similarity in the protein sequences of the two enzymes using the common BLASTP algorithm. The ratio of the sequence identity and length is then used to define the enzyme similarity.

## Selection of Metabolic Networks and Plotting Pathway Diagrams

To compare the differences between the two aforementioned metabolic pathway alignment methods, we analyzed a few aspects of the results obtained using both the methods. First, amino sugar and nucleotide sugar metabolic model alignment was conducted for *Dickeya dadantii* Ech703 and *Escherichia coli* K-12 MG1655, and the differences between the accurate matching results obtained from the two methods were compared (Table 1). Then, the performances of these two methods were investigated with respect to aligning multiple metabolic networks and constructing phylogenetic trees. The source of the metabolic networks is shown in Table 1.

The metabolic pathway diagrams (Fig. 1) of these two species (*D. dadantii* Ech703 and *E. coli* K-12 MG1655) were plotted using the software VANTED (Junker B et al. 2006) (version 2.2.1). The resulting relationship diagrams (Fig. 2 and Fig. 3) of the metabolic pathway alignment were preliminarily plotted using the newest version of Cytoscape.js (Ono et al. 2014) together with HTML. The phylogenetic trees (Fig. 4 and Fig. 5) were first constructed using the R package (Ihaka et al. 1996) and then plotted using the EvolView software (Zhang et al. 2012).



**Figure 1-** *Dickeya dadantii* Ech703 and *Escherichia coli* K-12 MG1655 amino sugar and nucleotide sugar metabolic networks from KEGG.

Jing, F et al.



**Figure 2-** Comparison of results obtained with the EC classification method. This result is the alignment between the *Dickeya dadantii* Ech703 and *Escherichia coli* K-12 MG1655 amino sugar and nucleotide sugar metabolic networks.



**Figure 3-** Comparison of results obtained with the sequence similarity method. This result is the alignment between the *Dickeya dadantii* Ech703 and *Escherichia coli* K-12 MG1655 amino sugar and nucleotide sugar metabolic networks.

## RESULTS

### Comparison of the Precise Details of the Two Network Alignment Results

The amino sugar and nucleotide sugar metabolic network models were selected from KEGG for *E.*

*coli* K-12 MG1655 and *D. dadantii* Ech703, each of which possessed 47 and 43 nodes, respectively, in their models. A 43*47-dimension similarity matrix was derived using the two similarity measurement methods, and the detailed results of the comparison are given in Table 1.

**Table 1-** Details of comparison between *D. dadantii* Ech703 and *E. coli* K-12 MG1655

| NO. | *D.dadantii ech703* | *E.coli K-12 MG1655* (Method of EC) | *E.coli K-12 MG1655* (Method of Blastp) |
|---|---|---|---|
| 0-A[a] | 2.7.1.69[b](Dd703_3193[c]) | 2.7.1.69(b0679[d]) | 2.7.1.69(b2417) |
| 1-S | 4.2.1.126(Dd703_2775) | 4.2.1.126(b2428) | 4.2.1.126(b2428 ) |
| 2-A | 2.7.1.69(Dd703_1115) | 2.7.1.69(b2417) | 2.7.1.69(b0679) |
| 3-D | 3.2.1.52(Dd703_1606) | 3.2.1.52(b1107) | 5.1.3.9(b3223) |
| 4-D | 2.7.1.59(Dd703_1621) | 2.7.1.59(b1119) | 3.2.1.52(b1107) |
| 5-S1 | 2.7.7.23(Dd703_3997) | 2.7.7.23(b3730) | 2.3.1.157(b3730 ) |
| 6-S | 3.5.1.25(Dd703_1113) | 3.5.1.25(b0677) | 3.5.1.25(b0677) |
| 7-S1 | 2.3.1.157(Dd703_3997) | 2.3.1.157(b3730) | 2.7.7.23(b3730 ) |
| 8-S | 2.7.1.8(Dd703_0709) | 2.7.1.60(b3222) | 2.7.1.60(b3222) |
| 9-S | 5.4.2.10(Dd703_3350) | 5.4.2.10(b3176) | 5.4.2.10(b3176) |
| 10-S | 5.1.3.14(Dd703_0205) | 5.1.3.14(b3786 ) | 5.1.3.14(b3786 ) |
| 11-S | 1.1.1.336(Dd703_0206) | 1.1.1.336(b3787 ) | 1.1.1.336(b3787 ) |
| 12-S | 3.5.99.6(Dd703_1114) | 3.5.99.6(b0678) | 3.5.99.6(b0678) |
| 13-S | 2.6.1.16(Dd703_3998) | 2.6.1.16(b3729) | 2.6.1.16(b3729) |
| 14-S | 2.7.1.4(Dd703_0696) | 2.7.1.4(b0394) | 2.7.1.4(b0394) |
| 15-S | 2.5.1.7(Dd703_3644) | 2.5.1.7(b3189) | 2.5.1.7(b3189) |
| 16-S | 1.3.1.98(Dd703_3741) | 1.3.1.98(b3972) | 1.3.1.98(b3972 ) |
| 17-S | 5.3.1.9(Dd703_0463) | 5.3.1.9(b4025 ) | 5.3.1.9(b4025 ) |
| 18-D | 4.1.1.35(Dd703_2685) | 4.1.3.3(b3225) | 3.2.1.37(b0271) |
| 19-S1 | 1.1.1.305(Dd703_4017) | 1.1.1.305(b2255) | 2.1.2.13(b2255) |
| 20-S | 2.6.1.87(Dd703_4015) | 2.6.1.87(b2253) | 2.6.1.87(b2253) |
| 21-S1 | 2.1.2.13(Dd703_4017) | 2.1.2.13(b2255) | 1.1.1.305(b2255) |
| 22-S | 2.4.2.53(Dd703_4016) | 2.4.2.53(b2254) | 2.4.2.53(b2254 ) |
| 23-S | 3.5.1.-(Dd703_4018) | 3.5.1.-(b2256) | 3.5.1.-(b2256) |
| 24-A | 2.7.1.69(Dd703_1595) | 2.7.1.69(b1817) | 2.7.1.69(b1101) |
| 25-S | 1.1.1.22(Dd703_1202) | 1.1.1.22(b2028) | 1.1.1.22(b2028) |

| | | | |
|---|---|---|---|
| 26-S | 2.7.7.9(Dd703_1203) | 2.7.7.9(b1236) | 2.7.7.9(b1236) |
| 27-S | 2.7.1.2(Dd703_0024) | 2.7.1.2(b2388 ) | 2.7.1.2(b2388 ) |
| 28-S | 5.4.2.2(Dd703_1123) | 5.4.2.2(b0688) | 5.4.2.2(b0688) |
| 29-S | 5.3.1.8(Dd703_1902 ) | 5.3.1.8(b1613) | 5.3.1.8(b1613) |
| 30-S | 2.7.7.12(Dd703_1172) | 2.7.7.12(b0758) | 2.7.7.12(b0758) |
| 31-S | 5.4.99.9(Dd703_1556) | 5.4.99.9(b2036) | 5.4.99.9(b2036) |
| 32-S | 5.1.3.2(Dd703_1173 ) | 5.1.3.2(b0759) | 5.1.3.2(b0759) |
| 33-D | 5.1.3.6(Dd703_1201) | 5.1.3.9(b3223) | 1.1.1.271(b2052) |
| 34-S | 2.7.1.6(Dd703_1171) | 2.7.1.6(b0757) | 2.7.1.6(b0757) |
| 35-A | 2.7.1.69(Dd703_2021) | 2.7.1.69(b1101) | 2.7.1.69(b1817) |
| 36-S | 5.4.2.8(Dd703_3280) | 5.4.2.8(b2048) | 5.4.2.8(b2048) |
| 37-S | 2.7.7.13(Dd703_3281 ) | 2.7.7.13(b2049) | 2.7.7.13(b2049) |
| 38-S | 4.2.2.1.47(Dd703_3277 ) | 4.2.1.47(b2053) | 4.2.1.47(b2053) |
| 39-S | 2.7.1.69(Dd703_1595) | 2.7.1.69(b1101) | 2.7.1.69(b1101) |
| 40-S | 2.7.1.2(Dd703_0024) | 2.7.1.2(b2388) | 2.7.1.2(b2388) |
| 41-S | 5.4.2.2(Dd703_1123) | 5.4.2.2(b0688) | 5.4.2.2(b0688 ) |
| 42-S | 2.7.7.27(Dd703_0279) | 2.7.7.27(b3430) | 2.7.7.27(b3430) |

a. A indicates that the results of the two methods have the same Enzyme Commission (EC) number, but have different entries; D indicates that the two methods have different EC numbers and different entries; S indicates that the two methods have the same EC number and same entry; S1 indicates that the two methods have different EC numbers and same entries.
b. The EC number of the amino sugar and nucleotide sugar metabolic network of *D. dadantii* Ech703 from KEGG.
c. The enzyme entry in the amino sugar and nucleotide sugar metabolic network of *D dadantii* Ech703 from KEGG.
d. The enzyme entry in the amino sugar and nucleotide sugar metabolic network of *E. coli* K-12 MG1655 from KEGG.

From Table 1, it is apparent that 31 consistent alignments were obtained using the two methods, which account for the majority of the metabolic network. This result demonstrates the strong conservation characteristic of the amino sugar and nucleotide sugar metabolic pathway.

Table S1 also reveals some differences between the two methods. The differences could mainly be categorized into two types: The first type included matches of No. 0, 2, 24, and 35. In match No. 0, PTS system glucose-specific transporter (Dd703_3193) of *D. dadantii* Ech703 was matched to fused N-acetyl glucosamine-specific PTS enzyme IIC, IIB, and IIA components (b0679) of *E. coli* K-12 MG1655

with the EC classification method, while it was matched to the same glucose-specific enzyme IIA component of PTS (b2417) with the sequence similarity method. Evidently, the sequence similarity method provided an accurate result. In match No. 2, PTS system N-acetylglucosamine-specific transporter subunit IIBC of *D. dadantii* Ech703 was matched to glucose-specific enzyme IIA component of PTS in the EC classification method, while in the sequence similarity method, it was matched to fused N-acetyl glucosamine-specific PTS enzyme IIC, IIB, and IIA components. Here, the sequence similarity method again delivered a better result than the EC classification method. Clearly, in this type of

matching, the EC classification method failed to match the correct enzymes, while accurate results were achieved with the sequence similarity method. We defined this type of matching as sequence-first matching. The other type of matches included nos. 3, 4, 18, and 33. In match No. 3, beta-N-acetylhexosaminidase of *D. dadantii* Ech703 was matched to beta N-acetyl-glucosaminidase of *E. coli* K-12 MG1655 in the EC classification method and it was matched to putative N-acetylmannosamine-6-P epimerase of *E. coli K-12 MG1655* in the sequence similarity method. In this case, the EC classification approach clearly produced the more superior result. This type of matching was therefore defined as the EC-first matching.

**Comparison of the Phylogenetic Trees Constructed Using the Two Alignment Methods**
In this study, ten strains of bacteria from seven genera of two classes, *Gammaproteobacteria* and *Betaproteobacteria*, were selected from KEGG (Table 1).

**Table 2-**List of the strains or organisms whose amino sugar and nucleotide sugar metabolic networks were used for alignment.

| Class | genus | species |
|---|---|---|
| *Gamma proteobacteria* | *Escherichia* | *Escherichia coli K-12 MG1655* |
| | *Escherichia* | *Escherichia coli O157H7 EDL933 (EHEC)* |
| | *Dickeya* | *Dicheya dadantii ech703* |
| | *Dickeya* | *Dickeya zeae Ech1591* |
| | *Klebsiella* | *Klebsiella pneumoniae subsp. pneumoniae MGH 78578* |
| | *Serratia* | *Serratia proteamaculans* |
| *Beta proteobacteria* | *Neisseria* | *Neisseria gonorrhoeae FA 1090* |
| | *Neisseria* | *Neisseria meningitidis Z2491 (serogroup A)* |
| | *Snodgrassella* | *Snodgrassella alvi* |
| | *Ralstonia* | *Ralstonia solanacearum GMI1000* |

The amino sugar and nucleotide sugar metabolic network alignment was conducted for each pair of the ten bacterial strains, and the alignment results were converted into distances between each pair of bacterial strains. We then performed clustering analysis based on these distances. The dendrograms (Figs. 4 and 5) derived from both methods were plotted using the EvolView software. It is clear from these dendrograms that both methods successfully explained the correlation between metabolic networks and the proximity in the phylogenetic relationships between each strain. Both methods also illustrated the species attribution of bacteria in *Gammaproteobacteria* and *Betaproteobacteria* and accurately described the phylogenetic relationships between individual bacterial strains in the genera *Escherichia* and *Neisseria*. However, for the two species *D. dadantii* Ech703 and *D. zeae* Ech1591 in the genus *Dickeya*, the phylogenetic tree constructed using the EC classification method did not accurately reflect their phylogenetic relationship (see the second rectangular mark part in Fig. 5).



**Figure 4-**Phylogenetic tree produced using the EC classification method. The rectangular marks represent the strains or species belonging to the same genus. Vertical bar represents the strains or species belonging to the same class.

Jing, F et al.



**Figure 5-** Phylogenetic tree produced using the sequence similarity method. The rectangular marks represent the strains or species belonging to the same genus. Vertical bar represents the strains or species belonging to the same class.

## DISCUSSION

We found that, for the alignment of any two networks, it is possible to obtain inconsistent matching results from these two methods. These inconsistencies mainly included two types of matches, i.e., the sequence-first matches and the EC-first matches. The reason for the emergence of a sequence-first match is mainly due to the existence of other classes of enzymes below the lowest class in the classification hierarchy; i.e., when the current enzyme classification system failed to accurately describe lower enzyme classes, the EC classification method would produce multiple equivalent matches. As a result, the program could not make accurate choices and consequently generated many false matches, while the sequence similarity method was able to make correct choices. On the other hand, the main reason for the appearance of EC-first matches is that the pairwise evolutionary distance between two species or enzymes was too large. This will result in an excessively low sequence similarity measure of the same enzyme from different species, which is possibly even lower than the similarity measures with certain enzymes that work on the same substrates but perform different functions. In this case, it is prone to result in this type of match.

Therefore, the probability of sequence-first or EC-first matches depends on the specific circumstances with regard to the enzymes and sequences of the metabolic networks being compared, as well as the evolutionary distances between the species. Under normal circumstances, the evolutionary status of the enzymes does not differ significantly from each other, and, as a result, the probability of the sequence-first match is higher than that of the EC-first match. In this case, the inaccuracy of the sequence similarity method will be lower than that of the method based on EC classification. Furthermore, because the EC classification method does not contain any evolutionary information, the sequence similarity method will provide superior results when constructing phylogenetic trees. In the comparison of the metabolic networks for ten bacterial strains, the sequence similarity method accurately reflected the phylogenetic relationships between their metabolic networks, while the results of the EC classification method had several mistakes.

In summary, for most pathways, the matching results of the metabolic network alignments from the EC classification method and the enzyme gene sequence similarity method were the same. However, there were some differences in the results, which depend on the context under which the comparison is being made. The similarity measurement method based on the EC classification does not take into account the fact that enzymes show differential expression under specific conditions in certain species and that these differences result in distinct metabolic networks. Therefore, in general, similarity measures based on the sequence similarity method will yield higher accuracy. This is especially the case when reflecting phylogenetic relationships, where the sequence similarity method evidently performs better than the EC

classification method.

## ACKNOWLEDGMENTS

## REFERENCES

Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res*. 1997; 25(17): 3389-3402.

Amir-Ghiasvand F, Nowzari-Dalini A, Momenzadeh V. Pin-Align: A New Dynamic Programming Approach to Align Protein-Protein Interaction Networks.*Comput Math Method M*. 2014;2014(13); Available from: http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4241705/

Doncheva N T, Assenov Y, Domingues F S, Albrecht M. Topological analysis and interactive visualization of biological networks and protein structures[J]. *Nature protocols*, 2012, 7(4): 670-685.

Francke C, Siezen R J, Teusink B. Reconstructing the metabolic network of a bacterium from its genome. *Trends Microbiol*. 2005;13(11): 550-558

Ihaka R, Gentleman R. R: a language for data analysis and graphics. *J Comput Graph Stat*. 1996;5(3): 299-314.

Junker B H, Klukas C, Schreiber F. VANTED: a system for advanced data analysis and visualization in the context of biological networks. *BMC Bioinformatics*. 2006;7:109(13). Available from:http://www.biomedcentral.com/1471-2105/7/109/

Kuchaiev O, Milenkovic T, Memisevic V, Hayes W, Przulj N. Topological network alignment uncovers biological function and phylogeny. *J R Soc Interface*. 2010;7(50): 1341-1354.

Li Y, Ridder D, Groot M,Reingers M. Metabolic pathway alignment (M-Pal) reveals diversity and alternatives in conserved networks. *BMC Syst Biol*, 2008;2:111(13) Available from: http://siplab.tudelft.nl/sites/default/files/apbc040a.pdf

Li Z, Zhang S, Wang Y, Zhang XS, Chen L. Alignment of molecular networks by integer quadratic programming. *Bioinformatics*. 2007;23(13): 1631-1639.

Ono K, Demchak B, Ideker T. Cytoscape tools for the web age: D3. js and Cytoscape. js exporters. *F1000Research*[Internet].2014.Available from:http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4264639/

Sharan R, Ideker T. Modeling cellular machinery through biological network comparison[J]. *Nature Biotechnology*. 2006, 24(4): 427-433.

Tohsato Y, Matsuda H, Hashimoto A. A multiple alignment algorithm for metabolic pathway analysis using enzymehierarchy. In: the Eighth International Conference on Intelligent Systems for Molecular Biology, 2000. *Syst Mol Biol*. 2000;376–383. Available from: ftp://genbank.sdsc.edu/pub/sdsc/biology/ISMB00/026.pdf

Yamada T, Bork P. Evolution of biomolecular networks—lessons from metabolic and protein interactions[J]. *Nature Reviews Molecular Cell Biology*. 2009, 10(11): 791-803.

Zhang H, Gao S, Lercher MJ, Hu S, Chen WH. EvolView, an online tool for visualizing, annotating and managing phylogenetic trees. *Nucleic Acids Res*. 2012; 40(W1): W569-W572.