

SubMAP: Aligning metabolic pathways with subnetwork mappings

Ferhat Ay^{*†‡} Manolis Kellis[‡] Tamer Kahveci[†]

Abstract

We consider the problem of aligning two metabolic pathways. Unlike traditional approaches, we do not restrict the alignment to one-to-one mappings between the molecules (nodes) of the input pathways (graphs). We follow the observation that in nature different organisms can perform the same or similar functions through different sets of reactions and molecules. The number and the topology of the molecules in these alternative sets often vary from one organism to another. With the motivation that an accurate biological alignment should be able to reveal these functionally similar molecule sets across different species, we develop an algorithm that first measures the similarities between different nodes using a mixture of homology and topological similarity. We combine the two metrics by employing an eigenvalue formulation. We then search for an alignment between the two input pathways that maximizes a similarity score, evaluated as the sum of the similarities of the mapped subsets of size at most a given integer k , and also does not contain any conflicting mappings. Here we prove that this maximization is NP-hard by a reduction from *Maximum Weight Independent Set (MWIS)* problem. We then convert our problem into an instance of MWIS and use an efficient vertex-selection strategy to extract the mappings that constitute our alignment. We name our algorithm **SubMAP** (**S**ubnetwork **M**appings in **A**lignment of **P**athways). We evaluate its accuracy and performance on real datasets. Our empirical results demonstrate that SubMAP can identify biologically-relevant mappings that are missed by traditional alignment methods and, that it is scalable for metabolic pathways of arbitrary topology, including searching for a query pathway of size 70 against the complete KEGG database of 1,842 pathways.

Keywords. *metabolic pathway alignment, subnetwork mappings, one-to-many mappings, alternative reaction sets, maximum weight independent set*

[†]Computer and Information Science and Engineering, University of Florida, Gainesville, FL

[‡]MIT Computer Science and Artificial Intelligence Laboratory, Cambridge, MA

*Correspondece author. E-mail: fay@cise.ufl.edu

**Implementation in C++ available at <http://bioinformatics.cise.ufl.edu/SubMAP.html>

1 Introduction

Biological pathways show how different molecules interact with each other to perform vital functions. In the literature, the terms *network* and *pathway* are used interchangeably for different types of interaction data. Metabolic pathways, an important class of biological pathways, represent how different compounds are transformed through various reactions. Analyzing these pathways is essential in understanding the machinery of living organisms.

The efforts on analyzing pathways can be classified into two types. The first type takes one pathway into account at a time and explores its important properties such as its robustness [1], steady states [2], modular structure [3, 4], network motifs [5, 6] as well as its representation [7, 8]. The second type is the comparative approach which considers multiple pathways to identify their frequent subgraphs [9, 10] and their alignments [11–21]. Alignment is a fundamental type of comparative analysis which aims to identify similar parts between pathways. For metabolic pathways, these similarities provide insights for drug target identification [22, 23], metabolic reconstruction of newly sequenced genome [24], phylogenetic reconstruction [25, 26] and enzyme cluster and missing enzyme identification [27, 28].

In the literature, alignment is often considered as finding one-to-one mappings between the molecules of two pathways. In this case, the global/local pathway alignment problems are GI/NP complete as the graph/subgraph isomorphism problems can be reduced to them in polynomial time [29]. A number of studies have been done to systematically align different types of biological networks. For metabolic pathways, Pinter *et al.* [14] devised an algorithm that aligns query pathways with specific topologies by using a graph theoretic approach. Tohsato *et al.* proposed two algorithms for metabolic pathway alignment one relying solely on Enzyme Commission (EC [30]) numbers of enzymes and the other considering only the chemical structures of compounds of the query pathways [19, 20]. Latterly, Cheng *et al.* developed a tool, *MetNetAligner*, for metabolic pathway alignment that allows a certain number of insertions and deletions of enzymes [21]. However, these methods focus on a single type of molecule and the alignment is driven only by the similarities of these molecules (e.g., enzyme similarity, compound similarity). Furthermore, some of these methods limit the query pathways to certain topologies, such as trees, non-branching paths or limited cycles, which degrades their applicability to complex pathways. Recently, Singh *et al.* [15] and Ay *et al.* [17, 18] combined both topological features and homological



Figure 1: A portion of Lysine biosynthesis pathway. Each reaction is represented by the Enzyme Commission (EC) number of the enzyme that catalyze it. Circles represent compounds (intermediate compounds are not shown). *E.coli* and *H.sapiens* (human) use the path colored by grey with three reactions, whereas plants and *Chlamydia* achieve this transformation directly through the path with a single reaction shown in white.

similarity of pairwise molecules to find the alignments of protein interaction networks and metabolic pathways respectively. These algorithms showed that this combination increases the accuracy of the alignment. Another advantage of these two methods is that they do not restrict the topologies of query pathways and hence are applicable to arbitrarily complex pathways.

All the methods discussed above limit the possible molecule mappings to only one-to-one mappings. As also pointed out by Deutscher *et al.* [31] considering each molecule one by one fails to reveal its function(s) in complex pathways. This restriction prevents all the above methods from identifying biologically relevant mappings when different organisms perform the same function through varying number of steps. As an example, there are alternative paths for LL-2,6-Diaminopimelate production in different organisms [23, 32]. LL-2,6-Diaminopimelate is a key intermediate compound since it lies at the intersection of different paths on the synthesis of L-Lysine. Figure 1 illustrates two paths both producing LL-2,6-Diaminopimelate starting from 2,3,4,5-Tetrahydrodipicolinate. The upper path represents the shortcut used by plants and *Chlamydia* to synthesize L-Lysine. This shortcut is not an option, for example, for *E.coli* or *H.sapiens* due to the lack of the gene encoding LL-DAP aminotransferase (2.6.1.83). *E.coli* and *H.sapiens* have to use a three step process shown with grey path in Figure 1 to do this transformation. Thus, a meaningful alignment should map the two paths when, let say, the lysine biosynthesis pathways of human and a plant are aligned. However, since these two paths have different number of reactions traditional

alignment methods, limited to one-to-one mappings, fail to identify this mapping.

Our aim in this paper is to design an algorithm that can accurately identify such biologically relevant mappings by allowing one-to-many mappings of molecules. Note that, in Figure 1 both reaction sets form linear paths. It is possible to have reaction sets with different topologies performing a certain function (See Figure 4 for examples). Therefore, we use the term *subnetwork* to include all types of topologies. Also, since we only consider the sets of reactions that are connected, we will simply use the term *subnetwork* instead of *connected subnetwork*.

Problem definition: We consider the problem of aligning two metabolic pathways. Unlike traditional alignment approaches, we allow aligning a molecule of one pathway to a connected subnetwork of the other. More formally, let \mathcal{P} and $\bar{\mathcal{P}}$ be two query pathways and k be a positive integer. We want to find the mapping between the molecules of \mathcal{P} and $\bar{\mathcal{P}}$ with the largest alignment score, such that (1) each molecule in \mathcal{P} ($\bar{\mathcal{P}}$) can map to a subnetwork of $\bar{\mathcal{P}}$ (\mathcal{P}) with at most k molecules and (2) each molecule can appear in at most one mapping.

The first condition above allows one-to-many mappings. The reason why we want to have one-to-many mappings in our alignment is not only that they capture functionally similar parts, but also they enable us to construct many-to-many mappings of arbitrary sizes. Identifying many-to-many mappings of molecule sets of different sizes is essential and is not possible with only one-to-one mappings as their combinations enforce both sides of a many-to-many mapping to be of the same size. The second condition enforces *consistency*. That is, if a molecule is already mapped alone or as a part of a subnetwork, it cannot map to another molecule. We elaborate on consistency and the problem definition later in Section 2. *Note that, allowing one-to-many mappings in alignment introduces new computational challenges (e.g., exponential increase in the problem size, conflicting mappings) that cannot be addressed using existing methods and hence novel methods are needed to tackle this problem.*

Contributions: In this paper, we propose a novel algorithm that finds subnetwork mappings in alignment of pathways. SubMAP accounts for both the effect of pairwise similarities (homology) and the organization of pathways (topology). This combination is motivated by its successful applications on pathway alignment by Singh *et al.* [15] and Ay *et al.* [17, 18]. However, al-

lowing one-to-many mappings makes it impossible to trivially extend these methods to our problem. To address this challenge, we map our problem to an eigenvalue problem. We solve this eigenvalue problem using an iterative technique called power method. The result of the power method converges to a principal eigenvector. This eigenvector defines a weighted bipartite graph where each node corresponds to a molecule or a subnetwork. The edges are only between two nodes from different pathways and their weights define the similarity of these nodes. Unlike the case of only one-to-one molecule mappings, the nodes on the same side of the bipartite graph can be intersecting as the same molecule can appear in more than one subnetwork. If two such intersecting nodes on one side are mapped to two different nodes of the other side, it will create inconsistent mappings for the elements of the intersection. We term such pair of mappings as *conflicting mappings*. Our aim is to extract a set of mappings that has no conflicting pairs and maximizes the sum of the similarity scores of the mappings. We prove that this maximization is NP-hard (See Theorem 1). We then construct a vertex weighted *conflict graph* with nodes representing a mapping of two subnetworks one from each pathway and edges representing a *conflict* between two mappings. The weights of the nodes in this graph are the edge weights of the bipartite graph from the earlier step. Our alignment problem is now equivalent to finding *maximum weight independent subset (MWIS)* of the conflict graph. To extract an independent set from the conflict graph, we use a vertex-selection strategy proposed for MWIS problem. We report the mappings that correspond to the selected set of nodes as the *alignment* of the query pathways. *Our experiments on the metabolic pathways from KEGG [33] database suggest that SubMAP finds biologically meaningful alignments efficiently. Also, SubMAP is scalable for subnetwork sizes up to three or four as it can align two real size metabolic pathways in a about a minute.*

The rest of the paper is organized as follows. Section 2 describes our algorithm. Section 3 presents experimental results. Section 4 concludes the paper.

2 Our Algorithm: SubMAP

In this section, we present our algorithm for pairwise metabolic pathway alignment that allows one-to-many molecule mappings. We begin by introducing some notation that we use throughout this section. Then, we formally

Table 1: Commonly used symbols in this paper

$\mathcal{P}, \bar{\mathcal{P}}$	Query metabolic pathways
k	Parameter for the largest subnetwork size
r_i, \bar{r}_j	Reactions of query pathways
R_i, \bar{R}_j	Subnetworks of query pathways
$\mathcal{R}_k, \bar{\mathcal{R}}_k$	Sets of all subnetworks with size at most k
n, m	Numbers of the reactions in query pathways
N_k, M_k	Numbers of all subnetworks of size at most k
φ_k	Set of all possible one-to-many mappings for a given k
$ \varphi_k $	Number of all possible one-to-many mappings for a given k
S	Support matrix
H	Homological similarity vector
G_c	Conflict graph
α	Parameter adjusting relative weights of homology and topology

state the problem and describe the SubMAP algorithm in detail.

Let, \mathcal{P} be a pathway which is represented by a directed unweighed graph $G = (V, E)$. Here, we only use the reactions of the pathway in graph representation. Hence, the vertex set $V = \{r_1, r_2, \dots, r_n\}$ is the set of all reactions of \mathcal{P} . We include a directed edge e_{ij} from r_i to r_j in E if and only if at least one output compound of r_i is an input compound of r_j . We call r_i a *backward neighbor* of r_j and r_j a *forward neighbor* of r_i if $e_{ij} \in E$. Note that reactions can be reversible (bi-directional) and hence both e_{ij} and e_{ji} can exist.

A *subnetwork* of a pathway is a subset of its reaction set such that the induced undirected graph of the elements of this subset forms a connected graph. Let $R_i \subseteq V$ be such a subnetwork of \mathcal{P} . We define \mathcal{R}_k , the set of all subnetworks of \mathcal{P} that have at most k reactions, as $\mathcal{R}_k = \{R_1, R_2, \dots, R_{N_k}\}$ where $|R_i| \leq k$ for all $i \in [1, N_k]$. Here, $|R_i|$ denotes the cardinality of the reaction set R_i . Using this notation, we define a binary relation that maps a reaction of a query pathway to a subnetwork of the other as follows:

Definition 1 Let \mathcal{P} and $\bar{\mathcal{P}}$ be two pathways and k be a positive integer. Also, let $\mathcal{R}_k = \{R_1, R_2, \dots, R_{N_k}\}$ and $\bar{\mathcal{R}}_k = \{\bar{R}_1, \bar{R}_2, \dots, \bar{R}_{M_k}\}$ be the sets of subnetworks with size at most k of \mathcal{P} and $\bar{\mathcal{P}}$ respectively. We define a binary relation φ between \mathcal{R}_k and $\bar{\mathcal{R}}_k$ that allows one-to-many reaction mappings as $\varphi : \varphi \subseteq \varphi_k = (\mathcal{R}_1 \times \bar{\mathcal{R}}_k) \cup (\mathcal{R}_k \times \bar{\mathcal{R}}_1)$. ■

The number of all possible one-to-many mappings is:

$$|\varphi_k| = nM_k + mN_k - nm \quad (1)$$

where n, m are the number of reactions of \mathcal{P} and $\bar{\mathcal{P}}$ respectively. The *alignment* of \mathcal{P} and $\bar{\mathcal{P}}$ is a binary relation φ that is a subset of all these possible mappings and satisfies certain criteria that we describe next.

Recall that for a mapping $(R_i, \bar{R}_j) \in \varphi$ one of the R_i or \bar{R}_j can contain more than one reaction. Reporting this mapping as a part of our alignment implies that all the reactions of the subnetwork with multiple reactions are aligned to a single reaction of the other. To have a *consistent alignment* none of the reactions of these subnetworks can be included in any other mapping. Next, we formally define the term *conflict* to characterize this property.

Definition 2 Let φ be a binary relation and $R_i, R_u \in \mathcal{R}_k$ and $\bar{R}_j, \bar{R}_v \in \bar{\mathcal{R}}_k$. The distinct pairs $(R_i, \bar{R}_j) \in \varphi$ and $(R_u, \bar{R}_v) \in \varphi$ **conflict** if and only if $(R_i \cap R_u) \cup (\bar{R}_j \cap \bar{R}_v) \neq \emptyset$. ■

Conflicts can cause inconsistencies about which reaction subset of one pathway should be aligned to the one of the other pathway. If φ has a conflicting pair of elements, we say φ is *inconsistent*. Since this is not a desirable property, we limit our alignment to the consistent relations only.

As well as discarding conflicting mappings we also need to use a meaningful scoring score in order to gather biologically relevant alignments. One standard scoring scheme for this purpose incorporates the homology of the aligned molecules with their topologies [15, 17, 18]. Here, we generalize this scheme to one-to-many mappings. We will elaborate on this similarity score later in Section 2.4. Next, we state our problem formally.

Problem formulation: Given k and two pathways \mathcal{P} and $\bar{\mathcal{P}}$, let \mathcal{R}_k and $\bar{\mathcal{R}}_k$ be the sets of subnetworks with size at most k of \mathcal{P} and $\bar{\mathcal{P}}$ respectively. We want to find the *consistent* binary relation $\varphi \subseteq (\mathcal{R}_1 \times \bar{\mathcal{R}}_k) \cup (\bar{\mathcal{R}}_k \times \mathcal{R}_1)$ that *maximizes* the summation of the similarity scores of the aligned subnetworks.

In the following, we present our algorithm SubMAP. Section 2.1 explains the enumeration of the subnetworks of query pathways. Section 2.2 and 2.3 discuss homological and topological similarities respectively. Section 2.4 describes the eigenvalue formulation that combines these similarities and explains the extraction of subnetwork mappings.

2.1 Enumeration of Connected Subnetworks

The first step of SubMAP is to create the sets of all connected subnetworks of size at most k for each query pathway. Here, we describe the enumeration process for a single query pathway. Let $G = (V, E)$ represent a pathway and k be a positive integer. We construct the set of subnetworks \mathcal{R}_k as follows. For $k = 0$, $\mathcal{R}_k = \mathcal{R}_0 = \emptyset$ and for $k = 1$, $\mathcal{R}_k = \mathcal{R}_1 = V$. For $k > 1$ we define \mathcal{R}_k recursively by using \mathcal{R}_{k-1} . At each recursive step we check for each reaction in V if it can be added to already enumerated subnetworks of size $k - 1$ to create a new connected subnetwork of size k . This way the k th recursive step takes $O(|V| \cdot (|\mathcal{R}_{k-1}| - |\mathcal{R}_{k-2}|))$ time.

The size of the set \mathcal{R}_k can be exponential in k when G is dense. However, metabolic pathways are usually sparse (on the average there are 2.5 forward neighbors per reaction). We observe that the number of subnetworks of the metabolic pathways in our dataset for $k = 3$ is around $5|V|$ and for $k = 4$ it is $10|V|$ on the average. In Section 3.2, we provide a detailed discussion of how $|\mathcal{R}_k|$ changes with different pathway sizes and different k values.

2.2 Homological Similarity of Subnetworks

Recall that the relation φ maps a reaction to a subnetwork that can contain multiple reactions. This necessitates computing the similarity between reaction sets. Since reactions are defined by their input and output compounds (i.e., substrates and products) and the enzymes that catalyze them, we measure the homological similarity between reactions using the similarities of these components.

In the literature, there are alternative pairwise similarity scores for compounds, enzymes and reactions. Particularly, two well known measure are information content similarity for enzyme pairs [14] and SIMCOMP [34] for compound pairs. We denote these measures by $SimE$ and $SimC$ respectively. We defer the readers to Ay *et al.* [18] for details on computing these similarities. Here, we utilize these similarity measures to compute the homological similarity between two reaction sets. To calculate this, we first construct the sets of the unions of input compounds (I_i), output compounds (O_i) and enzymes (E_i) of the reactions in each subnetwork R_i . For instance, in Figure 1 if we take upper path as the subnetwork R_i , then $E_i = \{2.3.1.117, 2.6.1.17, 3.5.1.8\}$. Let $\gamma_e, \gamma_i, \gamma_o$ denote the relative weights of

the similarities of enzymes, input and output compounds respectively. Then,

$$\begin{aligned} \text{SimRSet}(R_i, \bar{R}_j) = \\ \gamma_e W(E_i, \bar{E}_j, \text{SimE}) + \gamma_i W(I_i, \bar{I}_j, \text{SimC}) + \gamma_o W(O_i, \bar{O}_j, \text{SimC}) \quad (2) \end{aligned}$$

Here W denotes the sum of edge weights of the pairs returned by the maximum weight bipartite matching (MWBM) of the two sets. MWBM finds an assignment between the nodes of two sets such that maximizes the sum of the weights of these assignments specified by the similarity score. We use $\gamma_i = \gamma_o = 0.3$ and $\gamma_e = 0.4$ as they provide a good balance between enzymes and compounds.

We calculate SimRSet for all possible one-to-many mappings between the subnetworks of two pathways. Therefore, we do this calculation $|\varphi_k|$ times in total. This way, we assess the homological similarities between all possible subnetwork mappings. Even though this scoring can be considered a good measure of similarity, relying solely on this score ignores the topology similarity which we explain next.

2.3 Topological Similarity of Subnetworks

The motivation for utilizing topological similarity is that the induced topologies of two aligned subnetworks should also be similar. In other words, if R_i is mapped to \bar{R}_j , then their neighbors in the corresponding pathways should also be similar. Motivated by this, we first extend the *neighborhood* definition of reactions to reaction subnetworks. Then, we introduce the notion of *support* between two mappings. Using these definitions, we formally describe how we calculate *the support matrix* and show it on an illustrative example.

Definition 3 Let $R_i, R_u \in \mathcal{R}_k$. Then, R_u is a **forward neighbor** of R_i ($R_u \in \text{FN}(R_i)$) if and only if there exists $r_a \in R_i$ and $r_b \in R_u$ such that r_b is a forward neighbor of r_a or $R_i \cap R_u \neq \emptyset$. R_i is a **backward neighbor** of R_u ($R_i \in \text{BN}(R_u)$) if and only if R_u is a forward neighbor of R_i . ■

Definition 4 Let $R_i, R_u \in \mathcal{R}_k$ and $\bar{R}_j, \bar{R}_v \in \bar{\mathcal{R}}_k$. The mapping (R_i, \bar{R}_j) **supports** the mapping (R_u, \bar{R}_v) if and only if both $R_j \in \text{FN}(R_i)$ and $\bar{R}_v \in \text{FN}(\bar{R}_u)$ or both $R_j \in \text{BN}(R_i)$ and $\bar{R}_v \in \text{BN}(\bar{R}_u)$. ■

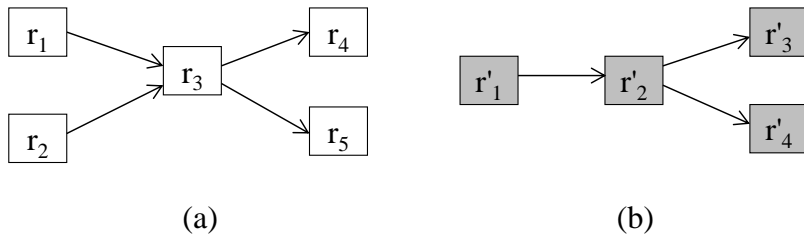


Figure 2: Reaction based representation of two hypothetical metabolic pathways. An arrow from r_i to r_j represents that an output of r_i is used as an input in r_j .

Definition 5 Let $\mathcal{P}, \bar{\mathcal{P}}$ be two pathways and $\mathcal{R}_k, \bar{\mathcal{R}}_k$ be the sets of their subnetworks that have at most k reactions. **The support matrix** S is a $|\varphi_k| \times |\varphi_k|$ matrix with each entry $S[i, j][u, v]$ identifying the fraction of the total support provided by (R_u, \bar{R}_v) mapping to (R_i, \bar{R}_j) mapping. Let $N(u, v) = |BN(R_u)||BN(\bar{R}_v)| + |FN(R_u)||FN(\bar{R}_v)|$ denote the number of all possible mappings between backward neighbors of R_u and \bar{R}_v plus the ones between their forward neighbors. Then, each entry of S is computed as:

$$S[i, j][u, v] = \begin{cases} \frac{1}{N(u, v)} & \text{if } (R_u, \bar{R}_v) \text{ supports } (R_i, \bar{R}_j) \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

■

Definition 4 states that the mapping of R_i to \bar{R}_j favors all possible mappings of forward (backward) neighbors of R_i to those of \bar{R}_j . To see it on an example, let us consider the case when $k = 1$ and focus on the mapping $(\{r_3\}, \{r'_2\})$ in Figure 2. We see that $FN(\{r_3\}) = 2$, $FN(\{r'_2\}) = 1$, $BN(\{r_3\}) = 2$ and $BN(\{r'_2\}) = 2$. Then, by Definition 5 we distribute the support of the mapping $(\{r_3\}, \{r'_2\})$ to $2 \times 1 + 2 \times 2 = 6$ other mappings by placing $\frac{1}{6}$ in the corresponding entries of S .

There can be cases when one mapping does not provide support to any others. In such cases, we simply distribute its support equally to all possible mappings ($|\varphi_k|$). Notice that, by construction, the entries in each column of S sums up to 1. This is important as it ensures the stability and convergence of our algorithm as we explain in Section 2.4. Interested reader can find detailed description of the properties of support matrix in a previous work of ours [18].

Before moving into describing how we use the support matrix, it is important to explain how we create it in our implementation. Trivial but costly way of doing this is to check each mapping against all the others to calculate the support values. However, such an exhaustive strategy will require computing a huge matrix S of size $|\varphi_k| \times |\varphi_k|$. Since the creation of S will incur prohibitive computational costs, we do not construct this matrix literally. Instead, for each mapping (R_i, \bar{R}_j) , we take the sets $FN(R_i)$, $FN(\bar{R}_j)$ and $BN(R_i)$, $BN(\bar{R}_j)$ to generate only the pairs supported by (R_i, \bar{R}_j) . In other words, we use the sparse matrix form of the support matrix S .

2.4 Aligning Two Pathways

2.4.1 Combining Homology and Topology

Both the homological similarities of subnetworks and their topological organization provide us significant information for the alignment of metabolic pathways. A good alignment algorithm needs to combine these two factors in an efficient and accurate way. Here, we describe how we achieve this combination in SubMAP by using an iterative technique called *power method*.

Let k be a given parameter and \mathcal{P} , $\bar{\mathcal{P}}$ be two pathways with connected subnetwork sets $\mathcal{R}_k = \{R_1, R_2, \dots, R_{N_k}\}$ and $\bar{\mathcal{R}}_k = \{\bar{R}_1, \bar{R}_2, \dots, \bar{R}_{M_k}\}$ respectively. We represent the homological similarity of all subnetwork pairs by the column vector H of size $|\varphi_k|$. Each entry of H denotes the homological similarity between two subnetworks one from each pathway which corresponds to a mapping.

Let S be the $|\varphi_k| \times |\varphi_k|$ support matrix as described in Section 2.3. Given a parameter $\alpha \in [0, 1]$ to adjust the relative weights of homology and topology, we combine homology and topology through power method iterations as follows:

$$H^{i+1} = \alpha S H^i + (1 - \alpha) H^0 \quad (4)$$

In this equation, $H^0 = \frac{H}{\|H\|}$. We iterate this equation till $H^{i+1} = H^i$ (i.e., it converges). The resulting vector H^i is the principal eigenvector of the matrix $\alpha S + (1 - \alpha) H^0 e$ where e is a row vector of size $|\varphi_k|$ with all entries equal to 1. This system converges to a unique principal eigenvector when the matrices S and $H^0 e$ are both column stochastic. We ensure this since each column of S adds up to one and H^0 is normalized. Each entry of H^i gives us a combination of homological and topological similarities for the corresponding mapping. We use $\alpha = 0.6$ in this paper since in our previous

work we observed that this value provides a good combination of the two similarities [17, 18].

2.4.2 Extracting Subnetwork Mappings

Recall that our aim is to find an alignment that *maximizes* the summation of the similarity scores defined by H^i while preserving the consistency between different mappings. Here we first show that this maximization is NP-hard by a reduction from *maximum weight independent set (MWIS)* problem in bounded degree graphs. MWIS problem, even for graphs with largest degree 3, is NP-hard [35] and there is no constant factor approximation to the optimal solution in polynomial time [36, 37]. We then describe how we construct a conflict graph from our alignment problem and utilize a vertex-selection heuristic in order to extract the set of mappings that generates the alignment.

Theorem 1 (FINDING THE BEST ALIGNMENT IS NP-HARD) *Let \mathcal{P} and $\bar{\mathcal{P}}$ be two pathways with reaction sets $\{r_1, \dots, r_n\}$ and $\{\bar{r}_1, \dots, \bar{r}_m\}$ respectively. Let $\mathcal{R} = \{R_1, R_2, \dots, R_N\}$ and $\bar{\mathcal{R}} = \{\bar{R}_1, \bar{R}_2, \dots, \bar{R}_M\}$ be all possible reaction subsets of \mathcal{P} and $\bar{\mathcal{P}}$ with size at most a given positive integer k . Also, let $w : (\mathcal{R}, \bar{\mathcal{R}}) \rightarrow [0, 1] \cup \{-\infty\}$ be a similarity function defining the score for each mapping and the conflicts between mappings are defined according to Definition 2. Then, finding a set of mappings (i.e., alignment) that maximizes the sum of mapping scores (i.e., alignment score) and has no conflicting pairs is NP-hard.*

Proof:

We show this by a reduction from MWIS problem in bounded degree graphs. Let $G = (V, E, w)$ be a vertex weighted undirected graph with largest degree $k - 1$ (i.e., $k = \max_{i=1, \dots, |V|} \deg(v_i) + 1$). Let us set $n = |V| + |E|$, $m = |V|$. Then, $N = \sum_{i=1}^k \binom{n}{i}$ and $M = \sum_{i=1}^k \binom{m}{i}$ are the numbers of all possible subsets of size at most k for the sets with sizes n and m respectively. Since k is bounded N and M are polynomial in n and m . We want to construct a new graph $G' = (V', E', w')$ that is of polynomial size in the size of G through a polynomial time reduction such that the solution for MWIS problem on G is equivalent to the set of mappings extracted from G' that gives us the best alignment of \mathcal{P} and $\bar{\mathcal{P}}$.

First, let us create a new vertex $v'_i \in V'$ for each $v_i \in V$ such that $v'_i = (R_i, \bar{R}_i)$ where initially each $R_i = \{r_i\}$, each $\bar{R}_i = \{\bar{r}_i\}$ for all i from 1

to $|V|$. We ensure such r_i and \bar{r}_i exist since $n = |V| + |E|$ and $m = |V|$. So far, we have not created any conflicts (i.e., edges) between any node pair (v'_i, v'_j) since $R_i \cap R_j = \bar{R}_i \cap \bar{R}_j = \emptyset$ for all $i \neq j$. Also, $R_i \in \mathcal{R}$ and $\bar{R}_i \in \bar{\mathcal{R}}$ since $|R_i| = |\bar{R}_i| = 1 \leq k$ for all i from 1 to $|V|$.

Next, for each $e_{ij} \in E$ we add $r_{h(i,j)}$ to both R_i ($R_i = R_i \cup r_{h(i,j)}$) and R_j ($R_j = R_j \cup r_{h(i,j)}$) so that we get a conflict between v'_i and v'_j . In order to ensure that we are only creating a conflict between the desired node pair (v'_i, v'_j) and no other pairs, $h(i, j)$ has to be a perfect hash function mapping each (i, j) to a distinct index between $|V| + 1$ and $n = |V| + |E|$. For this purpose, it suffices to simply choose $h(i, j) = |V| + (i - 1)|V| + j$. Also, for each vertex $v_i \in V$ we set $w'(v'_i) = w(v_i)$.

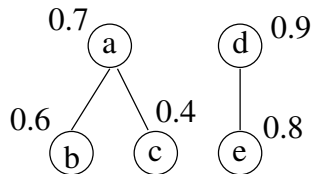
At this point, we have our graph G' isomorphic to G with their vertex weights having a one-to-one correspondence. Also, we have $\bar{R}_i \in \bar{\mathcal{R}}$ as $|\bar{R}_i| = 1 \leq k$ for each i from 1 to $|V|$. To have each $R_i \in \mathcal{R}$ as well, we need to make sure that $|R_i| \leq k$ for all i . Our construction starts with $|R_i| = 1$ for each i and adds one reaction to it for each undirected edge of v_i . Since G is a bounded degree graph with largest degree $k - 1$, we have at most $k - 1$ such additions and hence $|R_i| \leq k$.

Finally, we check for all possible mappings of the subsets, one from \mathcal{R} and one from $\bar{\mathcal{R}}$, whether (R_i, \bar{R}_j) exists as a node of G' or not. If this node does not exist, then we add it as an isolated node in G' with its weight set to $-\infty$ (i.e., $w'(R_i, \bar{R}_j) = -\infty$). After adding all such nodes, G' has $|V'| = NM$ vertices and $|E'| = |E|$ edges. Since k is bounded and N, M are polynomial in $n = |V| + |E|$ and $m = |V|$ respectively, the size of G' as well as the overall reduction is polynomial in the size of G . Furthermore, it can be seen from the construction that the maximum weight independent subset of G corresponds to the maximum weight independent subset of G' , hence to the set of non-conflicting mappings that maximizes the overall alignment score. Therefore, finding the best alignment is NP-hard. \blacksquare

Now that we proved extracting the best alignment is NP-hard, the next part is to describe how we tackle this problem. In the first step, we use the scores of mappings represented by H^i and the definition of conflict (Definition 2) to create a vertex weighted undirected graph $G_c = (V_c, E_c, w)$, which we name as *the conflict graph* as follows. We define a one-to-one correspondence from the mappings (R_i, \bar{R}_j) to the vertices in V_c . We also set the weight of each vertex $a = (R_i, \bar{R}_j) \in V_c$ (i.e., $w(a)$) to the similarity between R_i and \bar{R}_j as computed in H^i . We draw an undirected edge between two

	Subnetwork 1 (\mathcal{P})	Subnetwork ($\bar{\mathcal{P}}$)	H^i
a :	$R_1 = \{r_1, r_2\}$	$\bar{R}_1 = \{\bar{r}_1\}$	0.7
b :	$R_2 = \{r_1\}$	$\bar{R}_2 = \{\bar{r}_2\}$	0.6
c :	$R_3 = \{r_3\}$	$\bar{R}_3 = \{\bar{r}_1\}$	0.4
d :	$R_4 = \{r_4\}$	$\bar{R}_4 = \{\bar{r}_3, \bar{r}_4, \bar{r}_5\}$	0.9
e :	$R_5 = \{r_4, r_5\}$	$\bar{R}_5 = \{\bar{r}_5\}$	0.8

(a)



(b)

Figure 3: (a) Each row corresponds to a possible mapping between subnetworks from two hypothetical metabolic pathways. The first column is the unique label for each mapping. Second and third columns are the reactions in the two subnetworks that can be mapped. The last column is the similarity between two subnetworks. (b) The conflict graph G_c for the mappings in (a).

vertices $a = (R_i, \bar{R}_j)$ and $b = (R_u, \bar{R}_v)$ if $(R_i \cap R_u) \cup (\bar{R}_j \cap \bar{R}_v) \neq \emptyset$ (i.e., a and b conflict). For instance, in Figure 3 there is an edge between a and b representing that they conflict since reaction r_1 is common to both a and b .

In the second step, we explain the greedy vertex-selection strategy we adopt from Sakai *et al.* [38] in order to extract the MWIS of G_c as our alignment. Let $N(v)$ denote the set of vertices that are connected to $v \in V_c$. At each iteration of this algorithm, we pick a vertex v that maximizes $f(v) = \sum_{u_i \in N(v)} \frac{w(v)}{w(u_i)}$. This strategy implies that a vertex is more likely to be picked if the mapping it represents has large similarity score and conflicts with small number of other mappings with small similarity scores. After picking a vertex v , we put v into the result set and remove v and all the vertices connected to it (i.e., $v \cup N(v)$). We also remove all the edges incident to at least one of the removed vertices. When there are no more vertices to remove from G_c , the result set contains a maximal weight independent set. For our alignment problem, this vertex set corresponds to a list of non-conflicting subnetwork mappings. As an example, in Figure 3, d is the first vertex to be picked. Then, we remove d and $e \in N(d)$ from the graph and put d in the result set. Next, we pick the vertex b as $f(b) = \frac{0.6}{0.7} > f(a) = \frac{0.7}{0.6+0.4} > f(c) = \frac{0.4}{0.7}$. We remove b and $a \in N(b)$ and put b in the result set. Finally, only c is left and taking it into our result set, we have our alignment as the mappings $b = (r_1, \bar{r}_2)$, $c = (r_3, \bar{r}_1)$ and $d = (r_4, \{\bar{r}_3, \bar{r}_4, \bar{r}_5\})$.

3 Experimental Results

In this section, we experimentally evaluate the performance of SubMAP.

Dataset: We use the metabolic pathways of 20 organisms taken from the KEGG database. Our dataset contains 1,842 pathways in total. The average number of reactions per pathway is 21 and the largest pathway has 72 reactions. We also combined 12 different pathways under the metabolism of cofactors and vitamins for 10 different organisms. The average number of reactions of the largest connected components of these 10 pathways is 98 and the biggest one has 130 reactions.

3.1 Alternative Subnetworks

Different organisms can perform the same function through different subnetworks. We name such altered parts that have similar functions as *alternative subnetworks*. An accurate alignment should reveal alternative subnetworks in different pathways. In our first experiment we evaluate whether SubMAP can find them in real metabolic pathways. We align the pathway pairs which are known to contain functionally similar parts with different reaction sets and topologies. Table 2 presents a subset of reaction subnetwork mappings that are found by our algorithm. Figure 4 visualizes the topologies of these mappings by using an enzyme based representation.

The first row of Table 2 corresponds to alternative subnetworks in Figure 4(a) (also in Figure 1). The reaction R07613 represents the top path in Figure 4(a) that plants and Chlamydia use to produce LL-2,6- Diaminopimelate from 2,3,4,5- Tetrahydrodipicolinate. This path is discovered and reported as a shortcut on the L-Lysine synthesis path for plants and Chlamydia which is not present in humans or *E.coli* [23, 32]. Also, Watanabe *et al.* [23] suggested that since humans lack the catalyzer of the reaction R07613, namely LL-DAP aminotransferase (EC:2.6.1.83), this is an attractive target for the development of new drugs (antibiotics and herbicides). When we aligned the Lysine biosynthesis pathways of *H.sapiens* and *A.thaliana* (a plant), our algorithm mapped the reaction R07613 of *A.thaliana* to the three reactions that *H.sapiens* has to use to transform 2,3,4,5- Tetrahydrodipicolinate to LL-2,6- Diaminopimelate (R02734, R04365, R04475). In other words, SubMAP successfully identified the alternative subnetworks of different size (1 for *A.thaliana* and 3 for *H.sapiens*) that perform the same function.

Table 2: Alternative subnetworks that produce same or similar output compounds from the same or similar input compounds in different organisms. ¹ Main input compound utilized by the given set of reactions. ² Main output compound produced by the given set of reactions. ³ Reactions mappings that corresponds to alternative paths. Reactions are represented by their KEGG identifiers.

Pathway	Organisms	Input Comp. ¹	Output Comp. ²	Reaction Mappings ³
Lysine biosynthesis	<i>A.thaliana</i>	2,3,4,5-Tetrahydrodipico.	LL-2,6-Di-aminopimelate	R07613 \Leftrightarrow R02734 + R04365 + R04475
	<i>H.Sapiens</i>			
Lysine biosynthesis	<i>A.thaliana</i>	L-Saccharo. meso-2,6-Di.	L-Lysine	R00451 + R00715 + R00716 \Leftrightarrow R00451
	<i>H.sapiens</i>			
Pyruvate metabolism	<i>E.coli</i>	Pyruvate	Oxaloacetate	R00199 + R00345 \Leftrightarrow R00344
	<i>H.sapiens</i>			
Pyruvate metabolism	<i>E.coli</i>	Oxaloacetate	Phosphoenolpyruvate	R00341 \Leftrightarrow R00431 + R00726
	<i>H.sapiens</i>			
Pyruvate metabolism	<i>T.acidophilum</i>	Pyruvate	Acetyl-CoA	R01196 \Leftrightarrow R00472 + R00216 + R01257
	<i>A.tumefaciens</i>			
Glycine,serine, threonine met.	<i>H.sapiens</i>	Glycine	Serine	R00945 \Leftrightarrow R00751 + R00945 + R06171
	<i>R.norvegicus</i>		L-Threonine	
Fructose and mannose met.	<i>E.coli</i>	L-Fucose	L-Fucose 1-p	R03163 + R03241 \Leftrightarrow R03161
	<i>H.sapiens</i>		L-Fuculose 1-p	
Citrate cycle	<i>S.aureus N315</i>	Isocitrate	2-Oxoglutarate	R00268 + R01899 \Leftrightarrow R00709
	<i>S.aureus COL</i>			
Citrate cycle	<i>H.sapiens</i>	Succinate	Succinyl-CoA	R00432 + R00727 \Leftrightarrow R00405
	<i>A.tumefaciens</i>			
Citrate cycle	<i>H.sapiens</i>	Isocitrate	2-Oxoglutarate	R00709 \Leftrightarrow R00362
	<i>A.tumefaciens</i>	Citrate	Oxaloacetate	

Another interesting example is the second row that is extracted from the same alignment described above. In this case, the three reactions that can independently produce L-Lysine for *A.thaliana* are aligned to the only reaction that produces L-Lysine for *H.sapiens* (Figure 4(b)). R00451 is common to both organisms and it utilizes meso-2,6- Diaminopimelate to produce L-Lysine. The reactions R00715 and R00716 take place and produce L-Lysine in *A.thaliana* in the presence of L-Saccharopine [39].

For the alignment of Pyruvate metabolisms of *E.coli* and *H.sapiens*, the third and fourth rows show two mappings that are found by SubMAP. The first one maps the two step process in *E.coli* that first converts Pyruvate to Orthophosphate (R00199) and then Orthophosphate to Oxaloacetate (R00345) to the single reaction that directly produces Oxaloacetate from Pyruvate (R00344) in *H.sapiens* (Figure 4(c)). The second one shows another mapping in which a single reaction of *E.coli* is replaced by two reactions of *H.sapiens* (Figure 4(d)). The first two rows for Citrate cycle also report similar mappings for other organism pairs (Figure 4(e)).

Note that all the above examples (also the rows from 6 to 9, Figure 4(f)-(i)) are one-to-many reaction mappings and hence a merit of the new al-

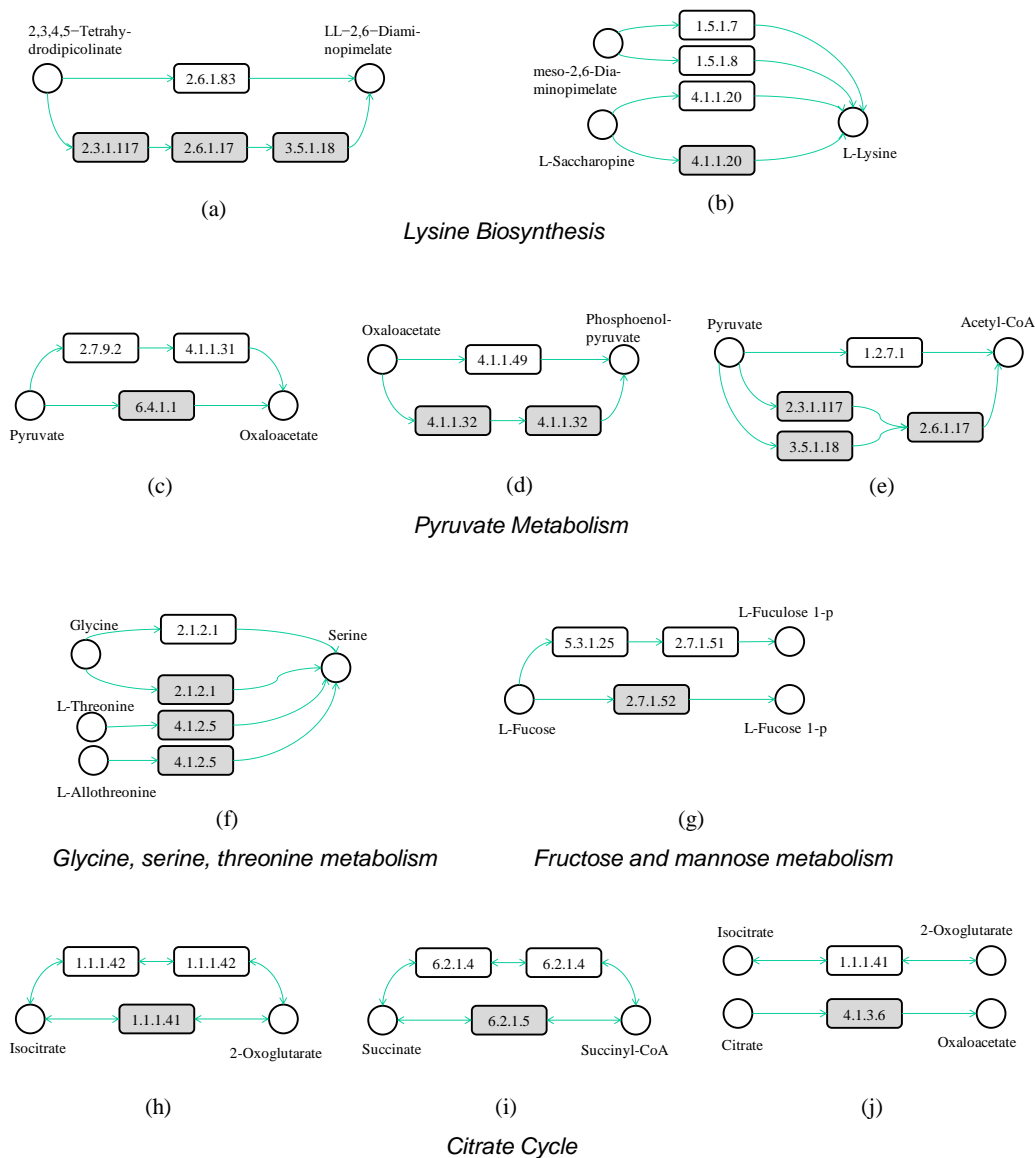


Figure 4: Visual representations of subnetwork mappings reported in Table 2. Figures (a) through (j) correspond to rows 1 through 10 of Table 2. Enzymes are represented by their Enzyme Commission (EC) numbers [30].

gorithm we propose here. Our algorithm SubMAP also reports one-to-one mappings. The last row of Table 2 is an example in which one reaction of an organism is replaced by exactly one reaction of another organism. Aligning Citrate cycles of *H.sapiens* and *A.tumefaciens* reveals that even though both the input and output compounds of two reactions R00709 and R00362 are different SubMAP maps these reactions (Figure 4(j)). Also, if we look at the EC numbers of the enzymes catalyzing these reactions (1.1.1.41 and 4.1.3.6) their similarity is zero (see Information content enzyme similarity [18]). If we were to consider only the homological similarities, these two reactions could not have been mapped to each other. However, both these reactions are the neighbors of two other reactions R01325 and R01900 that are present in both organisms. The mappings of R01325 to R01325 and R01900 to R01900 support the mapping of their neighbors R00709 to R00362. Therefore, by incorporating the topological similarity our algorithm is able to find meaningful mappings with similar topologies and distinct homologies. An algorithm not considering pathway topologies would fail to identify such mappings.

These results suggest: (i) By allowing one-to-many mappings, our method identifies functionally similar subnetworks even if they have different number of reactions. (ii) Incorporation of topological similarity makes it possible to find mappings that can be missed by only considering homological similarity.

3.2 Number of Connected Subnetworks

Given the parameter k , our algorithm enumerates all connected reaction subnetworks of size at most k for each query pathway. One question that we need to answer is: How many such subnetworks exist? Figure 5 plots this average for different pathway sizes in our dataset. When $k = 1$, the figure shows the number of reactions. For $k > 1$ the results demonstrate that the number of subnetworks increase exponentially with k . However, the increase is significantly lower than the theoretical worst case $\sum_{i=1}^k \binom{n}{i}$ (i.e., n choose i). For instance, the number of subnetworks we obtained for $n = 72$ and $k = 5$ is around 750 times less than the theoretical worst case. For $n = 130$ and $k = 5$ the number of subnetworks is 0.027% of the worst case.

The figure also suggests that the number of subnetworks increase linearly with the size of the pathway. This is mainly because the average number of edges (i.e., neighbors) of a node (i.e., subnetwork) remains roughly same as the size of the network increases. As a result, we conclude that for $k \leq 4$, we can enumerate and store all the subnetworks of the pathways in KEGG

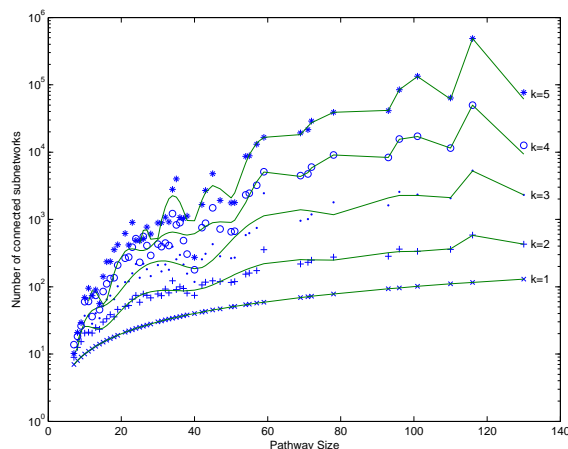


Figure 5: The number of subnetworks with at most k reactions for pathways of different sizes.

dataset. In practice it is unlikely for a single reaction to replace a subnetwork with more than 3 or 4 reactions. Therefore, we expect that using $k \leq 4$ would be sufficient to find most of the biologically relevant alternative subnetworks.

3.3 One-to-many Mappings within and across Major Clades

In Section 3.1, we demonstrated that our algorithm can find alternative subnetworks on a number of examples. An obvious question that follows is: How frequent are such alternative subnetworks and what are their characteristics? In other words, is there really a need to allow one-to-many mappings in alignment. In this experiment we aim to answer these questions.

We conduct an experiment as follows. We first pick 9 different organisms 3 from each major phylogenetic clade. These organisms are:

T.acidophilum, *Halobacterium sp.*, *M.thermoautotrophicum* from Archaea;
H.sapiens, *R.norvegicus*, *M.musculus* from Eukaryota;
E.coli, *P.aeruginosa*, *A.tumefaciens* from Bacteria.

We then extract 10 common pathways for these 9 organisms from KEGG. For each of these common pathways, we choose all possible pairs of the 9 organisms ($\binom{9}{2} = 36$) and align that specific pathway for all organism pairs. In these alignments we exclude the self alignments and the alignment with

Table 3: Percentages of 1-to-1, 1-to-2, 1-to-3 and 1-to-4 mappings in between and across three major clades. (**A**: Archaea, **E**: Eukaryota, **B**: Bacteria).

	1-to-1	1-to-2	1-to-3	1-to-4
E-E	89.6	8.8	1.1	0.5
B-B	80.1	16.0	3.1	0.8
A-A	78.3	15.7	4.7	1.3
B-E	69.1	23.1	6.3	1.5
A-B	60.5	28.3	8.5	2.7
A-E	55.8	31.0	10.4	2.8

parameter $k = 1$ since those will definitely incur a bias favoring the number of one-to-one alignments. We computed all possible alignments ($10 \times 36 = 360$) for $k = 2, 3$ and 4 ($360 \times 3 = 1,080$ alignments in total). Finally, we calculated the number of four possible types of subnetwork mappings which are 1-to-1, 1-to-2, 1-to-3 and 1-to-4. We hypothesize that the metabolisms of the organisms within a clade will tend to perform the same function through the same (or similar) sized sets of reactions while those across different clades will perform from alternative subnetworks of varying sizes.

Table 3 summarizes the results of this experiment. The percentages of each mapping type between two clades is shown as a row in this table. The first three rows corresponds to alignments within a clade and the last three represents alignments across two different clades. An important outcome of these results is that there are considerably large number of one-to-many mappings between organisms of different clades. In the extreme case (last row), nearly half of the mappings are one-to-many. The results also support our hypothesis that one-to-one mappings is more frequent for alignments within the clades compared to across clades due to high similarity between the organisms of the same clade. For instance, for both the first and last row one side of the query set is the Eukaryota. However, going from first row to last, we see around 40% decrease in the number of one-to-one mappings and 250%, 850% and 450% increase in the number of 1-to-2, 1-to-3 and 1-to-4 mappings respectively. Considering Archaea are single-cell microorganisms (e.g., Halobacteria) and Eukaryota are complex organisms with cell membranes (e.g., animals and plants), these jumps in the number of one-to-many mappings suggest that the individual reactions in Archaea are replaced by a number of reactions in Eukaryota. These results have two

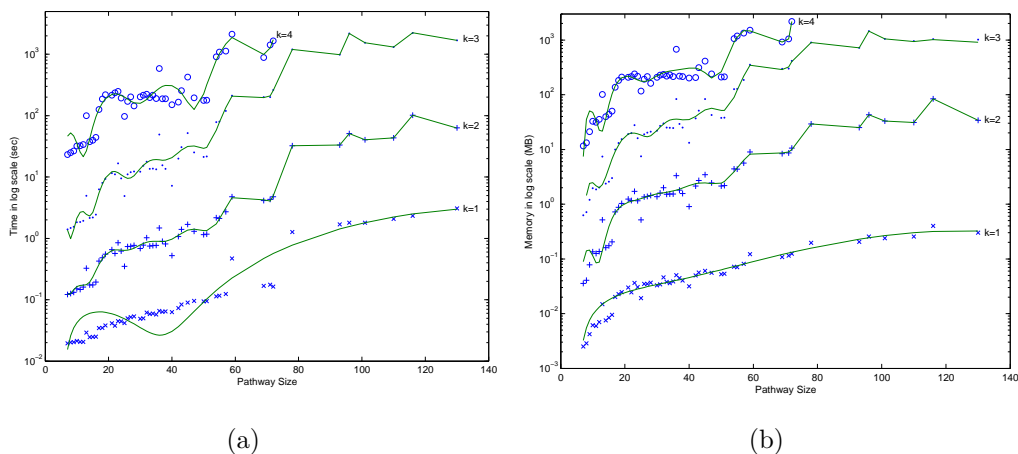


Figure 6: Average running time and memory use of SubMAP when a query pathway is aligned with a database containing 50 pathways. Pathway size is measured by the number of reactions.

major implications. (i) *One-to-many mappings are frequent in nature. To obtain biologically meaningful alignments we need to allow such mappings.* (ii) *The characteristics of the alternative subnetworks can help in inferring the phylogenetic relationship among different organisms.*

3.4 Running Time and Memory Utilization

SubMAP allows one-to-many mappings to find biologically relevant alignments. This however comes at the expense of increased computational cost. Theoretically, this increase can be exponential in k . The worst case happens when the pathway is highly connected. Metabolic pathways however are sparse and their connectivity follows power law distribution [40]. In order to understand the capabilities and limitations of our method we examine its performance on real datasets in terms of its running time and memory usage.

We evaluate the performance of our method for querying a database of pathways as follows. We create a query set by selecting 50 pathways of varying sizes from KEGG and adding it the combined pathways for cofactors and vitamins metabolism of 10 different organisms as described at the beginning of this section. We then select another 50 pathways of different sizes to use as our database set for this experiment. We pick the latter 50 pathways such that the average reactions per pathway is 21.4, which is very close to that of

the entire database. We then align each of the 60 query pathways with all the database pathways one by one for different values of k . We measure the average running time and the average memory usage for each query pathway and k value combination. Note that we do not present any performance comparison with an existing method as the existing methods do not allow one-to-many mappings. However, our results for $k = 1$, shows the performance of our algorithm when we restrict it to one-to-one mappings and can be compared with traditional alignment methods.

Figure 6(a) shows the average running time of SubMAP for query pathways with increasing number of reactions. When $k = 1$ (i.e., only one-to-one mappings as in existing methods), it runs in less than a few seconds even for the largest query pathway in our query set. As k increases, the running time increases significantly. This is because the number of subnetworks and the average number of forward and backward neighbors of subnetworks increase with k . We observe that our method can perform alignments in practical time even when $k = 4$ for pathways with around 70 reactions. Hence, it is scalable for all the individual pathways in KEGG. However, for the 10 combined pathways of cofactors and vitamins metabolism, the running time and memory use becomes a bottleneck when we consider $k > 3$.

We also measure the actual memory usage of our algorithm for real pathways of varying sizes and k values in Figure 6(b). For $k = 1$ or 2, the memory usage is negligible (100 MB or less) for all pathways. Although the memory usage increases quickly with k , it remains feasible for query pathways with around 70 reactions for $k = 4$. For $k = 3$, the biggest pathway of 130 reactions requires 1GB of memory per query on the average. *These results show that, SubMAP can run on a standard computer for aligning real-sized metabolic pathways for $k \leq 4$.*

4 Conclusion

In this paper, we considered the problem of aligning two metabolic pathways. The distinguishing feature of our work from the literature is that we allow mapping one molecule of one pathway to a set of molecules of the other. To address this problem, given two metabolic pathways \mathcal{P} and $\bar{\mathcal{P}}$ and an upper bound k on the size of the connected subnetworks, we developed the SubMAP algorithm that can find the consistent mapping of the subnetworks of \mathcal{P} and $\bar{\mathcal{P}}$ with the maximum similarity. We transformed the alignment problem to an eigenvalue problem. The solution to this eigenvalue problem

produced a good mixture of homological and topological similarities of the subnetworks. Using these similarity values, we constructed a vertex weighted graph that connects conflicting mappings with an edge. Then, our alignment problem is transformed into finding the maximum weight independent subset of this graph. We employed a heuristic method that is used to solve MWIS problem. The result of this method provided us an alignment that has no conflicting pair of mappings (i.e., consistent). Our experiments on real datasets suggested that our method can identify biologically relevant mappings of alternative subnetworks that are missed by traditional alignment methods. Furthermore, even though SubMAP does not restrict the topologies of query pathways, it is still scalable for real size metabolic pathways when the reaction subsets of size at most four are considered.

References

- [1] Edwards JS and Palsson BO. Robustness analysis of the Escherichia coli metabolic network. *Biotechnology Progress*, 16:927–939, 2000.
- [2] Ay F, Xu F, and Kahveci T. Scalable Steady State Analysis of Boolean Biological Regulatory Networks. *PLoS ONE*, 4(12):e7992, 2009.
- [3] Schuster S, Pfeiffer T, Koch I, Moldenhauer F, and Dandekar T. Exploring the Pathway Structure of Metabolism: Decomposition into Subnetworks and Application to Mycoplasma pneumoniae. *Bioinformatics*, 18:351–361, 2002.
- [4] Ay F, Dinh TN, and Thai MT Kahveci T. Finding dynamic modules of biological regulatory networks. In *BIBE*, 2010.
- [5] Wernicke S and Rasche F. FANMOD: a tool for fast network motif detection. *Bioinformatics*, 22(9):1152–3, 2006.
- [6] Grochow JA and Kellis M. Network motif discovery using subgraph enumeration and symmetry-breaking. In *RECOMB*, pages 92–106, 2007.
- [7] Michal G. On representation of metabolic pathways. *Biosystems*, 47(1-2):1–7, 1998.
- [8] Babur O, Dogrusoz U, Demir E, and Sander C. ChiBE: interactive visualization and manipulation of BioPAX pathway models. *Bioinformatics*, 26(3):429–31, 2010.

- [9] Koyuturk M, Grama A, and Szpankowski W. An efficient algorithm for detecting frequent subgraphs in biological networks. In *ISMB*, pages 200–207, 2004.
- [10] Qian X and Yoon B. Effective Identification of Conserved Pathways in Biological Networks Using Hidden Markov Models. *PLoS ONE*, 4(12):e8070, 2009.
- [11] Berg J and Lassig M. Local graph alignment and motif search in biological networks. *PNAS*, 101(41):14689–94, 2004.
- [12] Berg J and Lassig M. Cross-species analysis of biological networks by Bayesian alignment. *PNAS*, 103(29):10967–72, 2006.
- [13] Dandekar T, Schuster S, Snel B, Huynen M, and Bork P. Pathway alignment: application to the comparative analysis of glycolytic enzymes. *Biochemistry Journal*, 343:115–124, 1999.
- [14] Pinter RY, Rokhlenko O, Yeger-Lotem E, and Ziv-Ukelson M. Alignment of metabolic pathways. *Bioinformatics*, 21(16):3401–8, 2005.
- [15] Singh R, Xu J, and Berger B. Pairwise Global Alignment of Protein Interaction Networks by Matching Neighborhood Topology. In *RECOMB*, pages 16–31, 2007.
- [16] Koyuturk M, Grama A, and Szpankowski W. Pairwise Local Alignment of Protein Interaction Networks Guided by Models of Evolution. In *RECOMB*, pages 48–65, 2005.
- [17] Ay F, Kahveci T, and de Crecy-Lagard V. Consistent alignment of metabolic pathways without abstraction. In *Computational Systems Bioinformatics Conference (CSB)*, volume 7, pages 237–48, 2008.
- [18] Ay F, Kahveci T, and de Crecy-Lagard V. A fast and accurate algorithm for comparative analysis of metabolic pathways. *Journal of Bioinformatics and Computational Biology (JBCB)*, 7(3):389–428, 2009.
- [19] Tohsato Y and Nishimura Y. Metabolic Pathway Alignment Based on Similarity of Chemical Structures. *Information and Media Technologies*, 3:191–200, 2008.

- [20] Tohsato Y, Matsuda H, and Hashimoto A. A Multiple Alignment Algorithm for Metabolic Pathway Analysis Using Enzyme Hierarchy. In *ISMB*, pages 376–383, 2000.
- [21] Cheng Q, Harrison R, and Zelikovsky A. MetNetAligner: a web service tool for metabolic network alignments. *Bioinformatics*, 25(15):1989–90, 2009.
- [22] Sridhar P, Kahveci T, and Ranka S. An iterative algorithm for metabolic network-based drug target identification. In *Pacific Symposium on Bio-computing (PSB)*, volume 12, pages 88–99, 2007.
- [23] Watanabe N, Cherney MM, van Belkum MJ, Marcus SL, Flegel MD, Clay MD, Deyholos MK, Vederas JC, and James MN. Crystal structure of LL-diaminopimelate aminotransferase from *Arabidopsis thaliana*: a recently discovered enzyme in the biosynthesis of L-lysine by plants and *Chlamydia*. *Journal of Molecular Biology*, 371(3):685–702, 2007.
- [24] Francke C, Siezen RJ, and Teusink B. Reconstructing the metabolic network of a bacterium from its genome. *Trends in Microbiology*, 13(11):550–558, 2005.
- [25] Clemente JC, Satou K, and Valiente G. Reconstruction of Phylogenetic Relationships from Metabolic Pathways Based on the Enzyme Hierarchy and the Gene Ontology. *Genome Informatics*, 16(2):45–55, 2005.
- [26] Heymans M and Singh A. Deriving phylogenetic trees from the similarity analysis of metabolic pathways. *Bioinformatics*, 19:138–146, 2003.
- [27] Ogata H, Fujibuchi W, Goto S, and Kanehisa M. A heuristic graph comparison algorithm and its application to detect functionally related enzyme clusters. *Nucleic Acids Research*, 28:4021–8, 2000.
- [28] Green ML and Karp PD. A Bayesian method for identifying missing enzymes in predicted metabolic pathway databases. *BMC Bioinformatics*, 5:76, 2004.
- [29] Damaschke P. Graph-Theoretic Concepts in Computer Science. *Lecture Notes in Computer Science*, 484:72–78, 1991.
- [30] Webb EC. *Enzyme nomenclature 1992*. Academic Press, 1992.

- [31] Deutscher D, Meilijson I, Schuster S, and Ruppin E. Can single knock-outs accurately single out gene functions? *BMC Systems Biology*, 2:50, 2008.
- [32] McCoy AJ, Adams NE, Hudson AO, Gilvarg C, Leustek T, and Maurelli AT. L,L-diaminopimelate aminotransferase, a trans-kingdom enzyme shared by Chlamydia and plants for synthesis of diaminopimelate/lysine. *PNAS*, 103(47):17909–14, 2006.
- [33] Ogata H, Goto S, Sato K, Fujibuchi W, Bono H, and Kanehisa M. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids Research*, 27(1):29–34, 1999.
- [34] Hattori M, Okuno Y, Goto S, and Kanehisa M. Development of a chemical structure comparison method for integrated analysis of chemical and genomic information in the metabolic pathways. *Journal of the American Chemical Society (JACS)*, 125(39):11853–65, 2003.
- [35] Lovasz L. Stable set and polynomials. *Discrete Mathematics*, 124:137–53, 1994.
- [36] Austrin P, Khot S, and Safra M. Inapproximability of Vertex Cover and Independent Set in Bounded Degree Graphs. In *IEEE Conference on Computational Complexity*, pages 74–80, 2009.
- [37] Berman P and Karpinski M. On some tighter inapproximability results. *Lecture notes in Computer Science*, 1999.
- [38] Sakai S, Togasaki M, and Yamazaki K. A note on greedy algorithms for the maximum weighted independent set problem. *Discrete Applied Mathematics*, 126:313–22, 2003.
- [39] Saunders PP and Broquist HP. Saccharopine, an intermediate of amino adipic acid pathway of lysine biosynthesis. *Journal of Biological Chemistry*, 241:3435–40, 1966.
- [40] Jeong H, Tombor B, Albert R, Oltvai ZN, and Barabasi AL. The large-scale organization of metabolic networks. *Nature*, 407(6804):651–4, 2000.