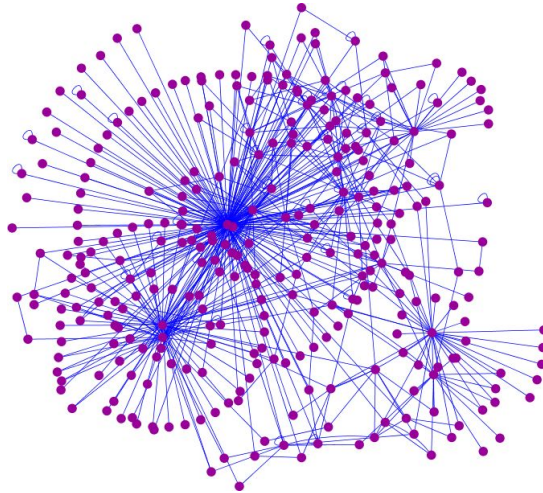# Top-k Similar Graph Matching in Biological Networks

Surabhi S Nath
Utsav Rohilla

# Introduction

- Biological networks are one of most complex systems
- Relationships between between molecules, entities, structural properties
- Crucial to understand life
- Graph based data structures efficient in capturing dynamics of networks
- A variety of static, dynamic graph based methods have been applied

# Problem Statement

- Graph matching important to retrieve information embedded in graphs
- Efficient alternative to exact graph matching - approximate graph matching
- Tram algorithm for large graph networks
- Given a graph database and multiple query graphs, find similar substructures from the data graph that match the query graph
- Applications include querying in protein networks - searching, comparing

# TraM Algorithm

- Computes similarity between two labelled graphs
- Algorithm computes the topmost k similar graphs based on different similarity criteria. This algorithm's strength lies in the fact that multiple domain dependent similarity scores can be incorporated into an overall similarity calculation of nodes
- Conceptual likeness of two sets of nodes is captured using a domain-dependent similarity function called the $\sigma$-similarity function
- Random walk captures structural information of the graph and similarity function is called the $\beta$-similarity function

# TraM Algorithm

- Before we get into the main algorithm, we need to get familiar with some definitions and helping algorithms

---

**Algorithm 1:** Random Walk

---

**Input**: Graph $G = (V_g, E_g)$ and Restart Probability $\beta$.

**Output**: Random Walk Score $P_s(V_g)$.

1 Let $r_s(V_g)$ be the restart vector with all entries having value $\frac{1}{|V_g|}$;

2 Let $A$ be the column normalized adjacency matrix defined by $E$;

3 Initialize $P_s(V_g) = r_s(V_g)$;

4 **while** $P_s(V_g)$ *has not converged* **do**

5     $P_s(V_g) = (1 - \beta) * A * P_s(V_g) + \beta * r_s(V_g))$;

---

# TraM Algorithm

**Definition 2 ($\beta$-signature).** *Let $G = (V_g, E_g)$ be a graph and $\mathbb{B}$ be a set of $\beta$ values. Random walk score of a node is a function of the form $\omega : V_g \times \mathbb{B} \to [0,1]$. An $n$-dimensional vector $(\omega(v, \beta_1), \omega(v, \beta_2), \ldots, \omega(v, \beta_n))$ is called a $\beta$-signature of node $v \in V_g$, denoted $\vec{\beta}(v)$ and the set $\beta(G) = \bigcup_{v \in V} \vec{\beta}(v)$ is called the $\beta$-signature of $G$.*

**Definition 4 ($\beta$-similarity).** *Let $G_1 = (V_{g_1}, E_{g_1})$ and $G_2 = (V_{g_2}, E_{g_2})$ be two graphs and $\beta(G_1)$ and $\beta(G_2)$ be their $\beta$-signatures. For any $v_1 \in V_{g_1}$ and $v_2 \in V_{g_2}$, their structural or $\beta$-similarity, denoted $\hat{\beta}(v_1, v_2)$, is defined by*

$$\hat{\beta}(v_1, v_2) = 1 - \sqrt{\sum_{i=1}^{k} (a_i - b_i)^2},$$

*where $\vec{\beta}(v_1) = (a_1, a_2, \ldots, a_k) \in \beta(G_1)$ and $\vec{\beta}(v_2) = (b_1, b_2, \ldots, b_k) \in \beta(G_2)$.*

# TraM Algorithm

**Definition 5 (Graph similarity).** Let $G_1 = (V_{g_1}, E_{g_1})$ and $G_2 = (V_{g_2}, E_{g_2})$ be two graphs and $\phi$ be a mapping function. Then, the similarity $\gamma$ between two graphs $G_1$ and $G_2$ under a mapping function $\phi$ is

$$\gamma(G_1, G_2) = \sum_{\forall v_1, v_2 (v_1 \in V_{g_1}, \phi(v_2) \in V_{g_2})} \sigma(v_1, \phi(v_2)) \times \hat{\beta}(v_1, \phi(v_2)). \tag{2}$$

# TraM Algorithm

**Definition 6 ($\delta$-neighborhood of nodes).** *Let $Q$ be a query graph, and $r$ be its radius.*[1] *Let $D = (V_d, E_d)$ be any graph, and $v \xrightarrow{j} u$ represent $j$-hop reachability from node $v$ to $u$. Then, $\delta$-neighborhood of $v$ is the set $\{u | v, u \in V_d \wedge v \xrightarrow{j} u \wedge j \leq r\}$ (including zero hop reachability, i.e., $v$ itself), denoted $\delta(r, v)$.*

**Definition 7 (Induced subgraphs).** *Let $G = (V_g, E_g)$ be any graph, and $N \subseteq V_g$ be an arbitrary set of nodes. Then, $G' = (N, E')$ is called an induced subgraph, denoted $\chi(N) = G' = (N, E')$, such that whenever $v, v' \in N$ and $e = (v, v') \in E_g$, $e$ is also in $E'$, and nothing else is in $E'$, i.e., $G'$ is a subgraph of $G$, denoted $G' \sqsubseteq G$.*

# TraM Algorithm

**Algorithm 2:** GraphMatch

**Input**: Data graph $D = (V_d, E_d)$ and query graph $Q = (V_q, E_q)$. Thresholds $k$, $\mu_v$, $\mu_s$ and $\lambda$.

**Output**: Top-$k$ matches of $Q$.

1  Initialize priority queue $PQ$ as empty;
2  Calculate $\beta$-signature $\beta(Q)$ for $Q$;
3  Compute radius $r$ of $Q$;
4  **for** $\forall v_d (v_d \in V_d)$ **do**
5       Compute $\delta(r, v_d)$;
6       **if** $|Filter(\delta(r, v_d), Q, \mu_v, \mu_s)| > |V_q|$ **then**
7           Top-$k$ Match($Q$, $\beta(Q)$, $D$, $\beta(\chi(\delta(r, v_d)))$, $\lambda$, $PQ$);

8  **return** *All top k graphs* $g \in PQ$ ;

# TraM Algorithm

**Algorithm 4:** Filter

**Input**: Set of nodes $\delta(r, v_d)$ and Graph $Q = (V_q, E_q)$.
Thresholds $\mu_v$ and $\mu_s$.

**Output**: Pruned $\delta(r, v_d)$.

1   **if** $|\delta(r, v_d)| > |V_q|$ **then**
2      **for** *every* $v_n \in \delta(r, v_d)$ **do**
3         Compute candidate subgraph as $\chi(\delta(r, v_d))$;
4         Compute $\beta$-signature $\beta(\chi(\delta(r, v_d)))$ for $\chi(\delta(r, v_d))$;
5         **if** $\forall v_q(v_q \in V_q(\max\{\sigma(v_n, v_q)\} < \mu_v$ *or*
6         $\max\{\hat{\beta}(v_n, v_q)\} < \mu_s))$ **then**
7           remove $v_N$ from $\delta(r, v_d)$;

8 **return** $\delta(r, v_d)$;

# TraM Algorithm

**Algorithm 3:** Top-$k$ Match

**Input**: Query graph $Q = (V_q, E_q)$, $\beta(Q)$, candidate graph $C = (V_c, E_c)$, $\hat{\beta}(C)$, threshold $\lambda$. Queue $PQ$

**Output**: Updated priority queue $PQ$.

1 **for** $i = 1$ *to Number of Nodes in time-stamp 1* **do**
2      initialize $S^{(i,1)} = \emptyset$ and $sim(i, S^{(i,1)})^{(1)} = 0$;

3 **for** $t = 2$ *to* $2 \times |V_q|$ **do**
4      **for** $r = 1$ *to Number of Nodes in time-stamp* $t$ **do**
5          **if** $t$ *is Even* **then**
6              **for** $p = 1$ *to Number of Nodes in time-stamp* $t - 1$ **do**
7                  **if** $(r \notin S^{(p,t-1)})$ *and* $((\exists i (i \in S^{(p,t-1)}$ *and* $r \xrightarrow{1} i$ *and* $i \in V_c)))$ *and* $(sim(p, S^{(p,t-1)})^{(t-1)} + \sigma(p,r) \times \hat{\beta}(p,r) \geq MAX)$ **then**
8                      $MAX = sim(p, S^{(p,t-1)})^{(t-1)} + \sigma(p,r) \times \hat{\beta}(p,r)$;
9                      $k = p$;
10                     $newMAX = true$;

11              **if** $newMAX = true$ **then**
12                  $S^{(r,t)} = S^{(k,t-1)} \cup \{r\}$;
13                  $sim(r, S^{(r,t)})^{(t)} = MAX$;

14          **if** $t$ *is Odd or* $newMAX = false$ **then**
15              $S^{(r,t)} = S^{(r,t-1)}$;
16              $sim(r, S^{(r,t)})^{(t)} = sim(r, S^{(r,t-1)})^{(t-1)}$;

17          **if** $t = 2 \times |V_q|$ *and last* $r$ **then**
18              $\gamma(Q,C) = sim(r, S^{(r,t)})^{(t)}$;

19 **if** $\gamma(Q,C) \geq \lambda$ **then**
20      Add $\langle C, \gamma(Q,C) \rangle$ to queue $PQ$;

21 **return** $PQ$