

Hierarchical High Level Information Fusion Using Graph Structures, Subgraph Matching and State Space Search

Moises Sudit

Rakesh Nagi

Kedar Sambhoos

February 7, 2007

In Nearly Every Application:

We Are Drowning in Data

- More Sources
- More Source Types: NTM, Tactical, Commercial
- Wider Bandwidths
- More Connectivity: GCCS, DI/COE



DATA:
Observations & Measurements

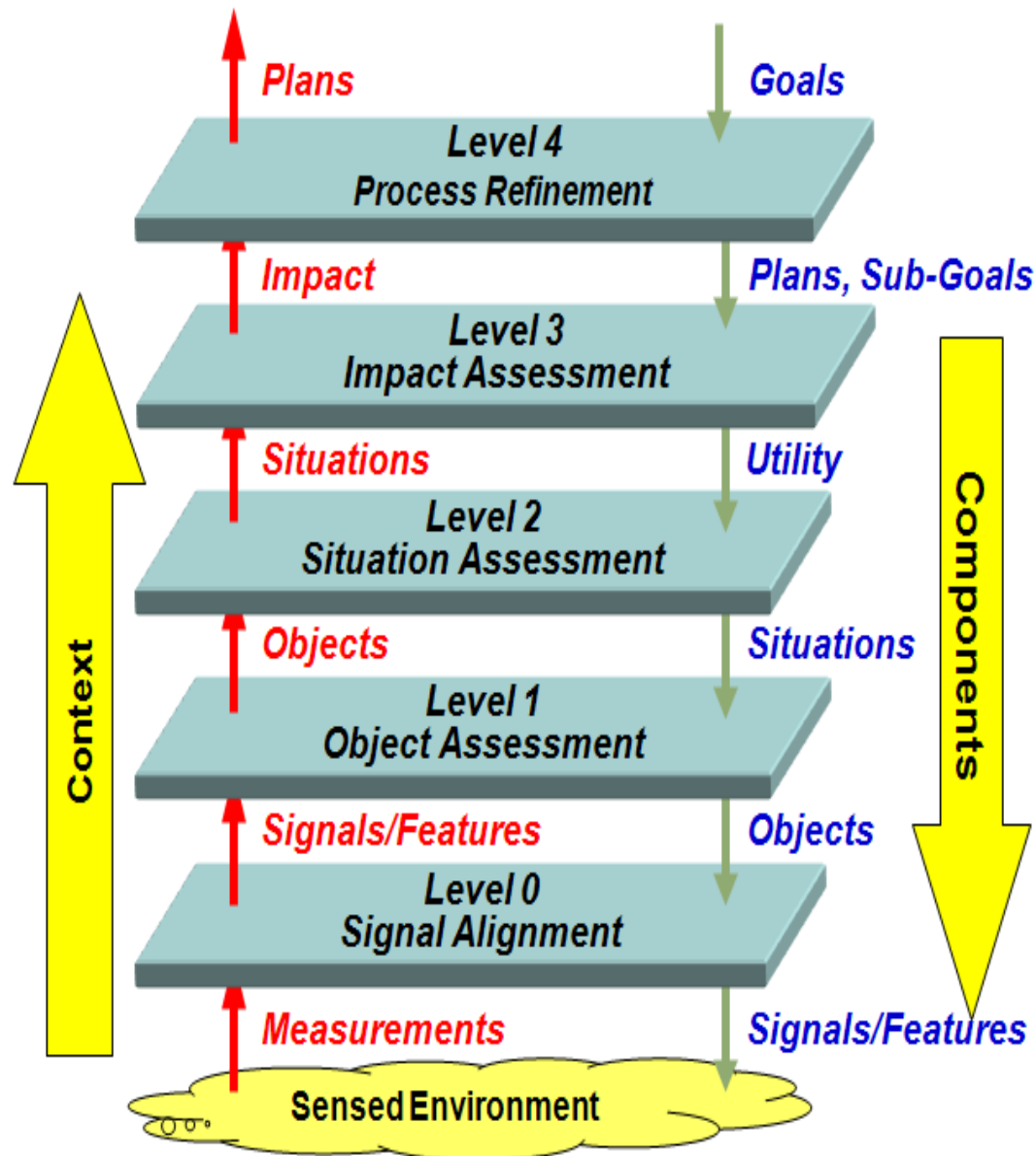


But Starving for Information

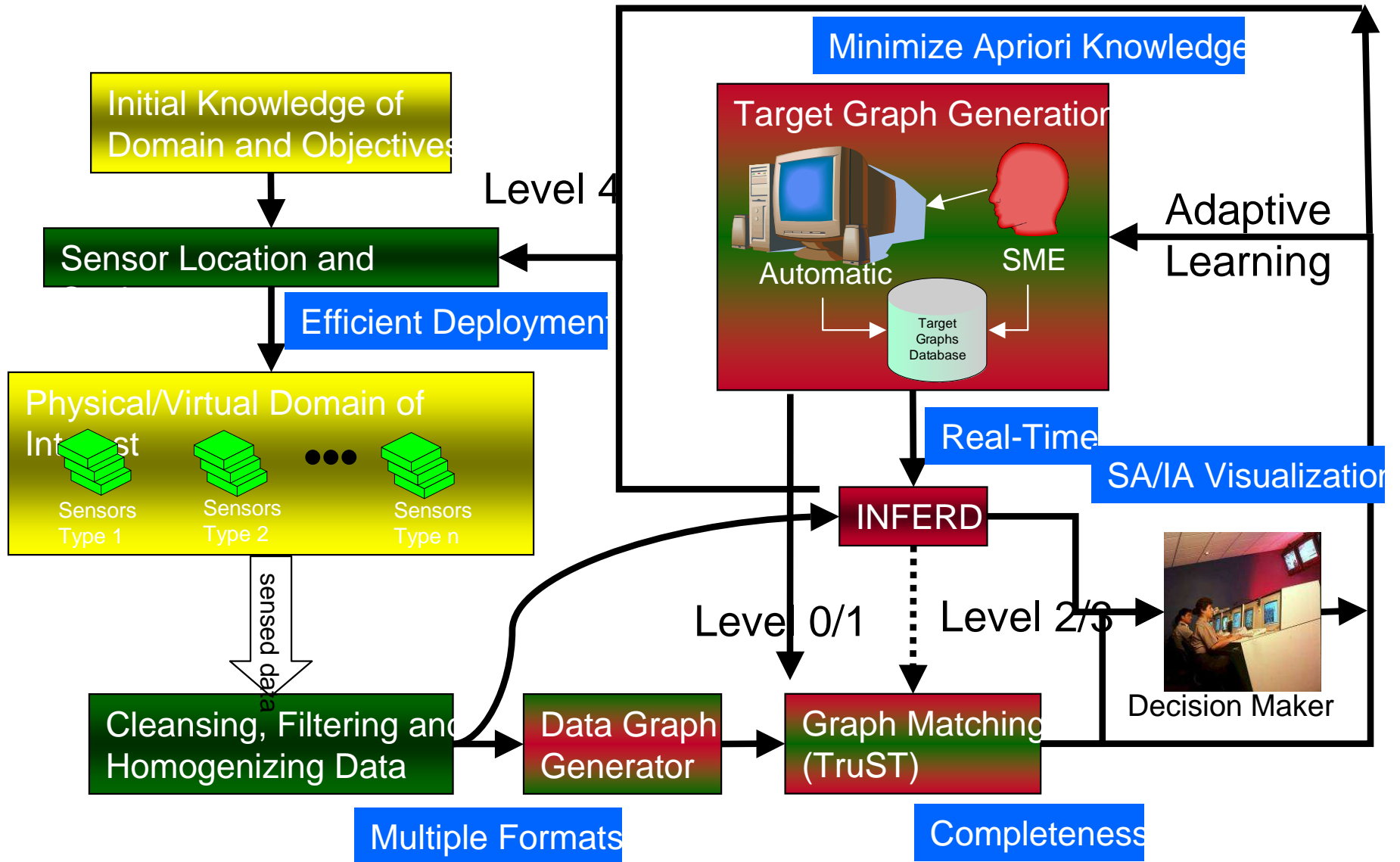
- How does data relate to the situation (context)?
- What is its significance (impact on my plans)?
- What is its quality? Can I trust it?
- How can it be confirmed, refuted or refined?



INFORMATION:
Organized Data

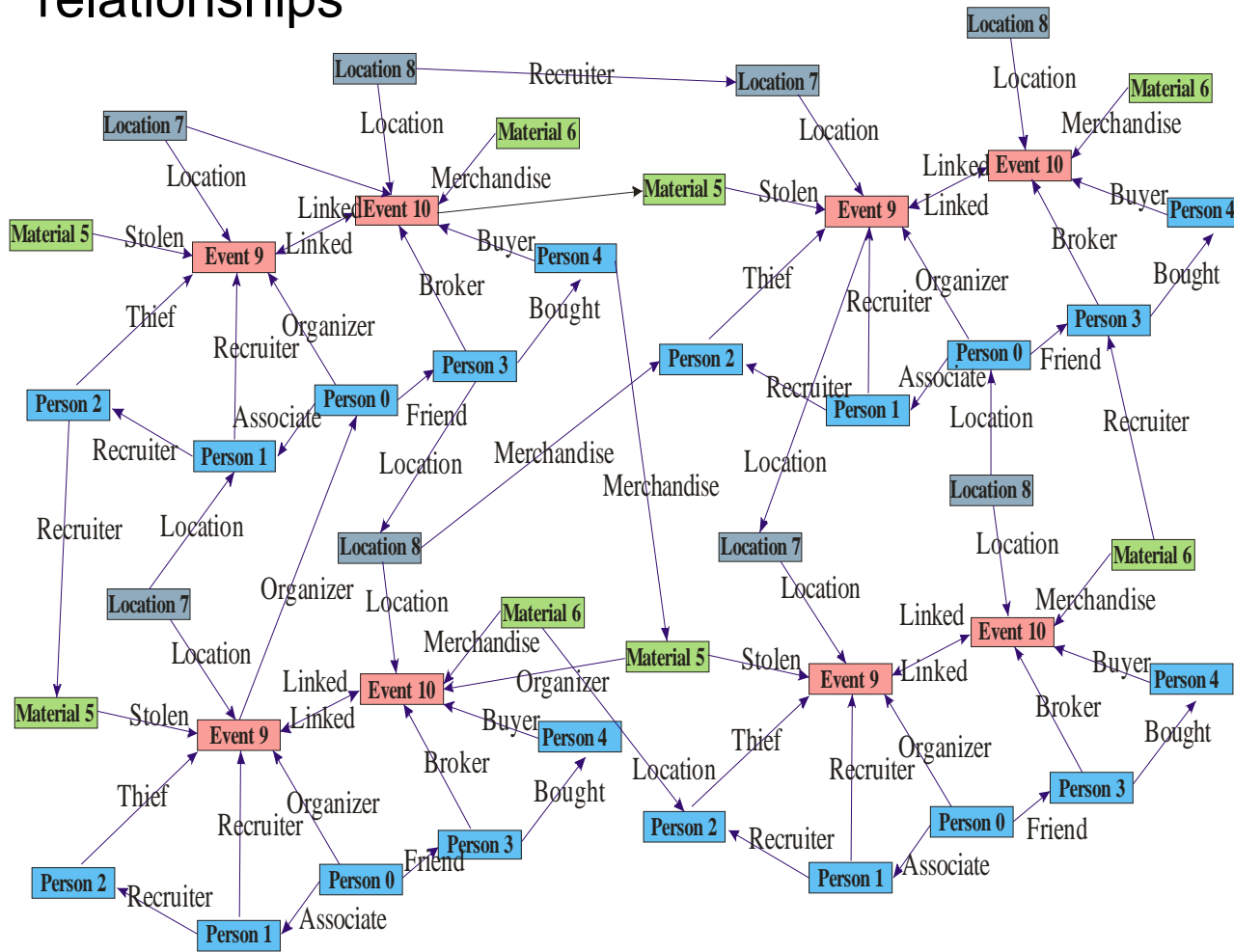


Technology Approach and Justification



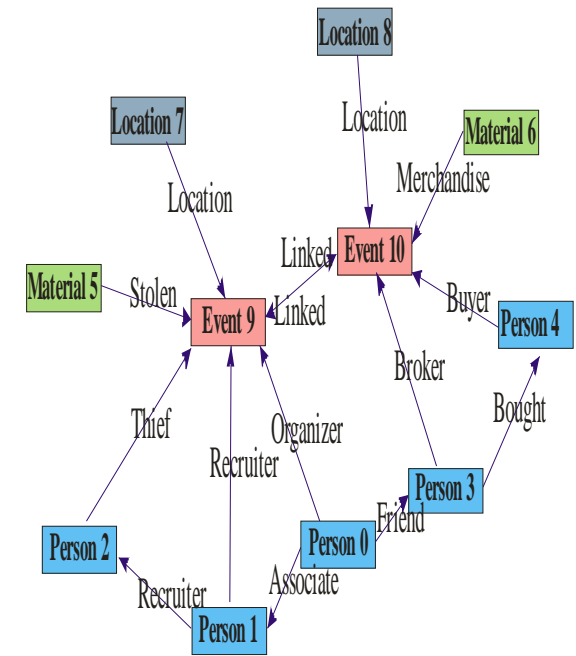
Data Graph and Template Graph

- Decision-maker often encounter complex uncertain situations and they need to develop *relationship* between situational elements.
- Attributed Relational Graphs* (ARGs) represent situational elements and relationships



Data Graph

UNCLASSIFIED



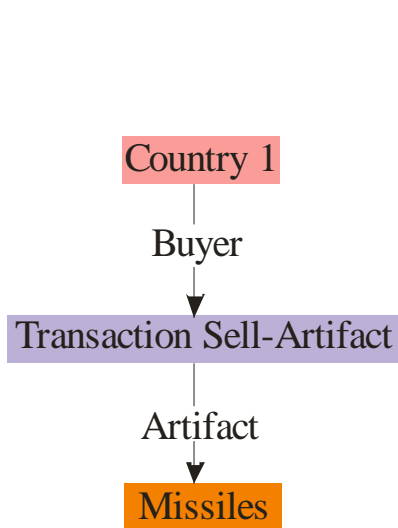
Template Graph

Graph Matching Structures

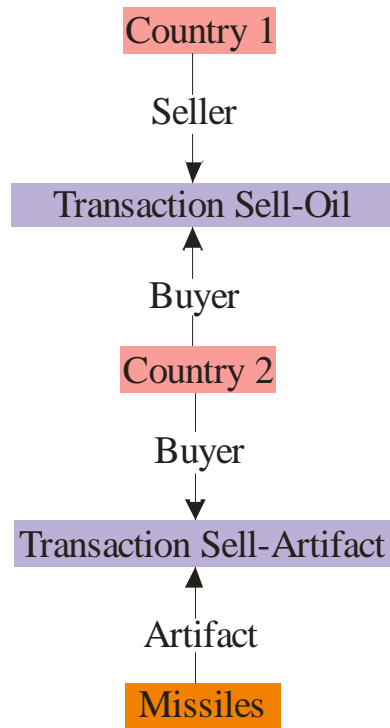
		Exact	Inexact
Syntactic (labelled-graphs)		Graph Isomorphism	Graph Homeomorphism
Semantic (attributed-graphs)	Structured	- attributed-vertex graph - Greedy Search	- attributed-vertex graph - Greedy Search
	Semistructured	- attributed-vertex graph - Graph Isomorphism - Graph Bundling	- attributed-vertex graph - Graph Homeomorphism - Graph Bundling
	Unstructured	Heuristics	Heuristics

- **Structured** data is rigidly organized & well defined: **Predictable**
- **Semistructured** data is organized enough to be predictable
 - Data is organized in semantic entities
 - Similar entities are grouped together but,
 - Entities in the same group may not have the same attributes
 - The order of the attributes is not necessarily important
 - The presence of some attributes may not always be required
 - The size of same attributes of entities in a same group may not be the same
 - The type of the same attributes of entities in a same group may not be of the same type
- **Unstructured** data is disordered and unruly: **Unpredictable**
- Difference between Graph Matching and **Subgraph Matching**

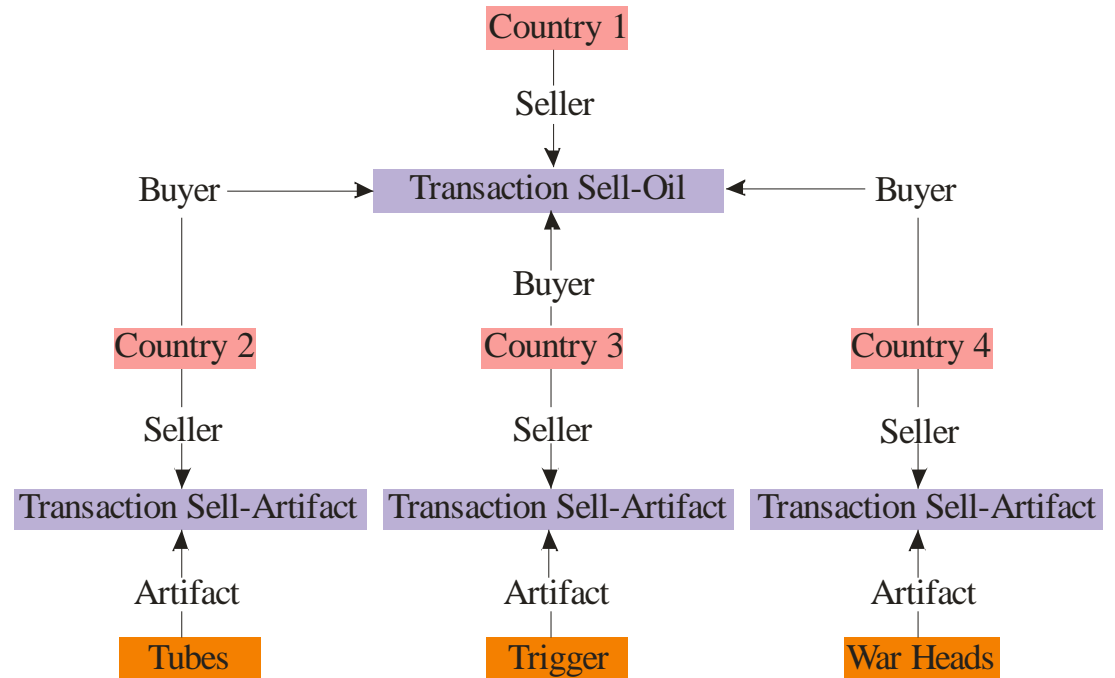
Isomorphism, Homeomorphism and Bundle Matching



Isomorphism



Homeomorphism



Bundle matching

Motivation for our Approach

- Enhancements in Level 2 and 3 fusion capabilities through a new class of models and algorithms in graph matching.
- Create a Framework to Span Temporal Decision-Making (from “real-time” to forensics”)
- Most of literature is in subgraph isomorphism (Cordella *et al.*, 2004)
- Most of the matching techniques are restricted to Simple Graphs
- *Truncated branch and bound* performs several times better than the local search (Zhang, 2000) motivates for Unstructured Graphs
- Attributed Relational Graphs are Visually Appealing, but don't give us much value-added compared with Classical Graph Structures in Solving Graph Matching Problems.

Theorem: Any graph matching problem over an attributed-graph (attributed-vertices and attributed-edges) can be polynomially transformed in to an equivalent attributed-vertex graph (no attributes on the edges)

1. L. P. Cordella, P. Foggia, C. Sansone, and M. Vento. A (sub)graph isomorphism algorithm for matching large graphs. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 26(10):1367, 2004.
2. W. Zhang. Depth-first branch-and-bound versus local search: A case study. In *Proc. 17th National Conf. on Artificial Intelligence*, pages 930–935, 2000.

Definitions and Notation

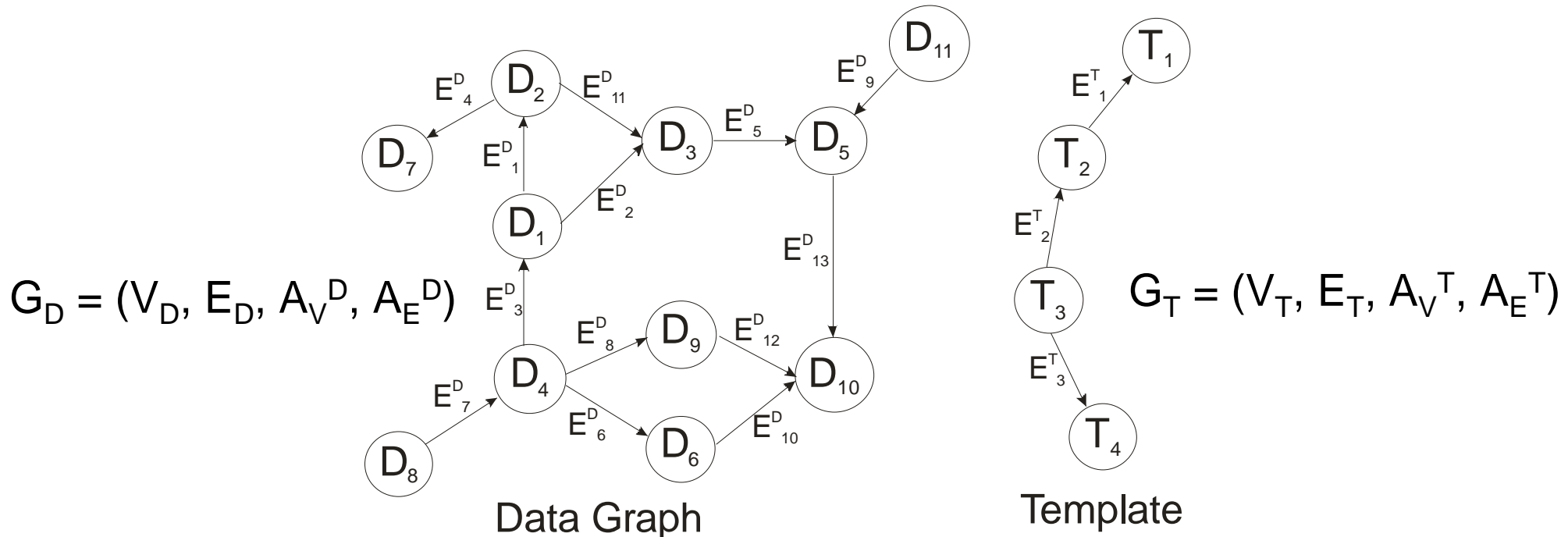
- Attributed Graph Structure

$$G = (V, E, A_V, A_E)$$

where V - the set of nodes;

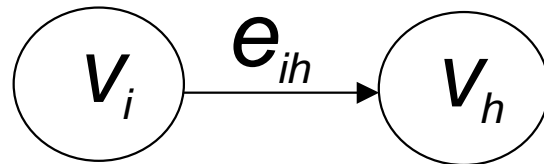
E - the set of arcs;

A_V - the set of node attributes; A_E - the set of arc attributes.

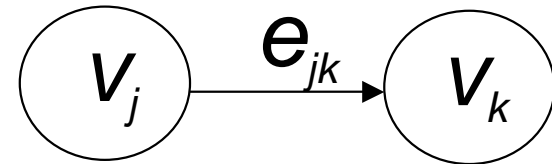


Most Graph Matching Techniques based on Triplets

- Triplets:



Template



Data Graph

- One-to-one Scores of Nodes and Arcs:

$$S(v_i, v_j) = 0.9 \quad S(v_h, v_k) = 0.7 \quad S(e_{ih}, e_{jk}) = 0.8$$

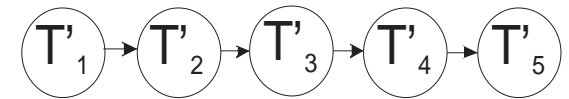
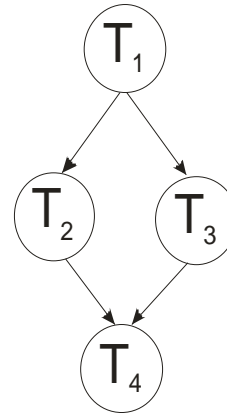
- Similarity Functions for Triplets:

$$\begin{aligned} f_{ihjk} &= f(S(v_i, v_j), S(v_h, v_k), S(e_{ih}, e_{jk})) \\ &= 0.8 \end{aligned}$$

Triplet Deficiency

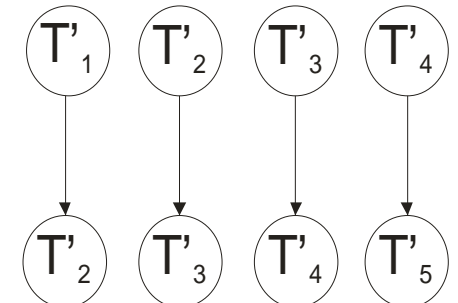
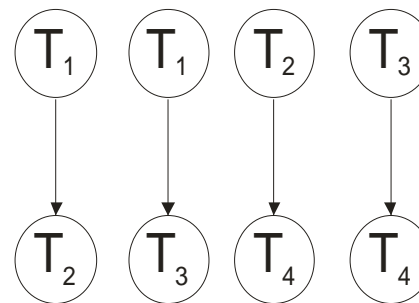
Problem:

- Two totally different topological templates could receive exactly the *same* similarity values.



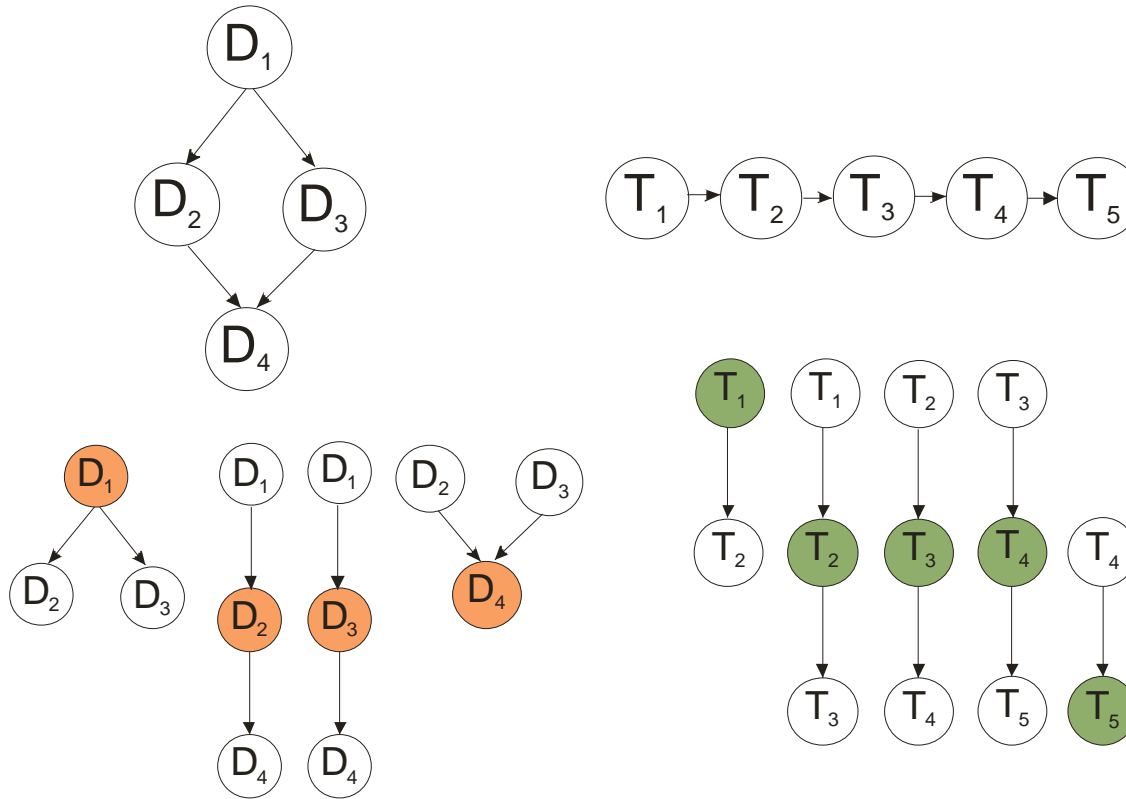
Goal:

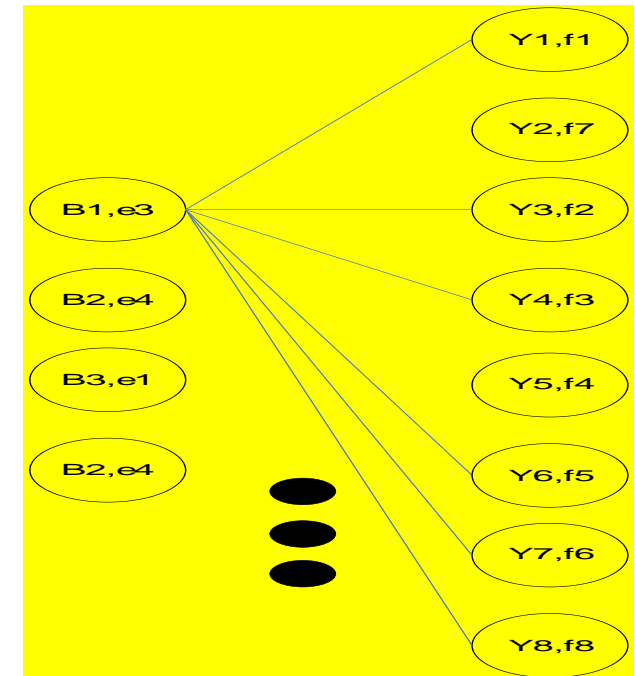
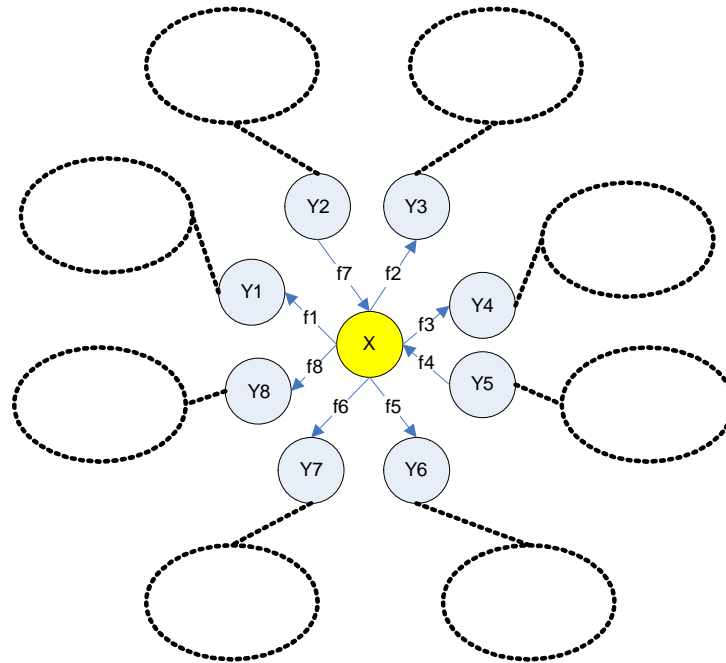
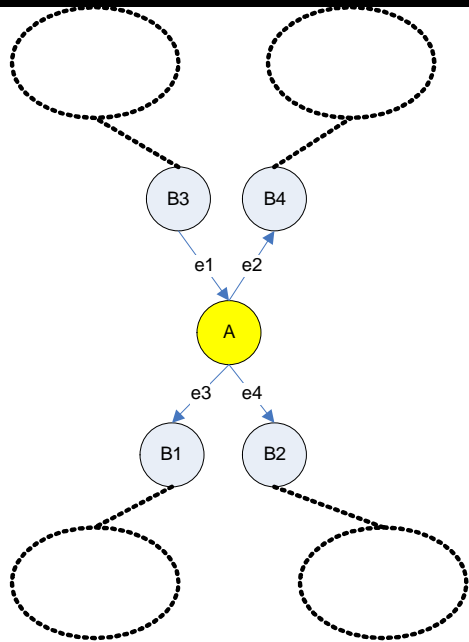
- Develop an algorithm that uses the current system's rule based scoring engine while increasing its *topological accuracy* during matching.



1-Hop Neighborhood Breakdown

- 1-Hop Neighborhood
 - A root node and all other nodes of edge distance 1 away.
 - Benefit:





• Algorithm

- **Step 1:** Compute a node score, denoted as C_{ij} , for each node in the template graph to each node in the data graph.
- **Step 2:** Compute the scores, denoted as W_{ij} , for the 1-Hop neighbors of each root node pair.

- The score is given by $\alpha C_{ij} + (1-\alpha) W_{ij}$

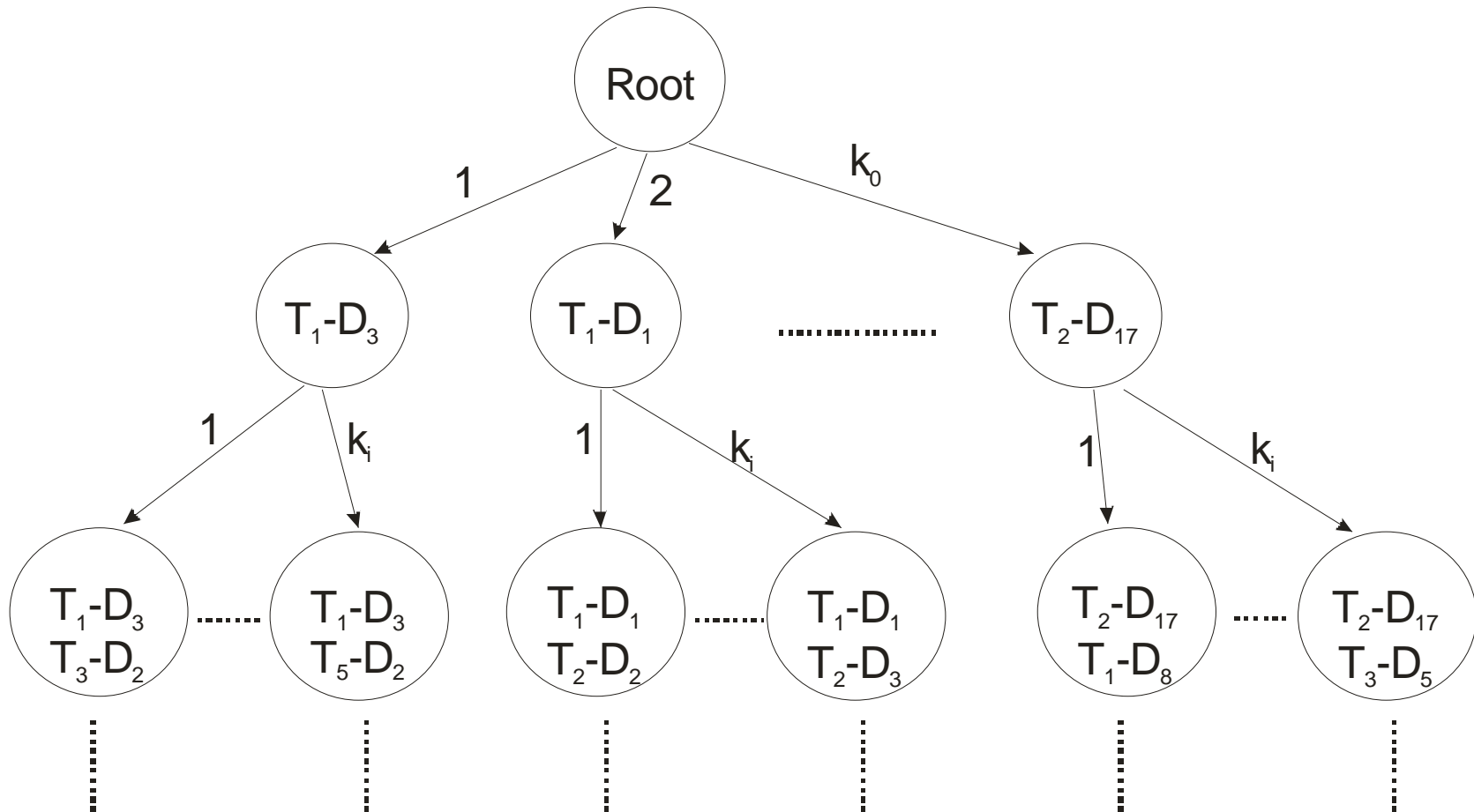
UNCLASSIFIED

“ α ” is the Score vs.
Topology Parameter

TruST (Truncated Search Tree) Algorithm

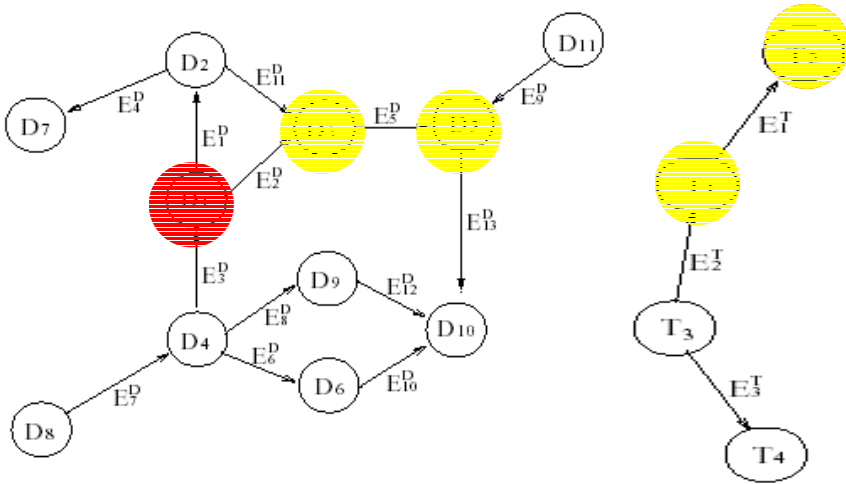
- Advantages:
 - Easier to implement, faster to execute and requires less computing resources.
 - User can decide tradeoff between time and quality.
 - Converges to optimality.
- Disadvantages:
 - Not guaranteed to yield an optimal solution.
 - Parametric approach
 - Greedy in nature.
- Parameters of controlling state space
 - k_0 : The number of child problems of the root.
 - k_i : The number of child problems of the problem at level $i-1$.
 - β : The total number of problems.
 - δ : The total number of levels.

TruST Algorithm (cont'd)



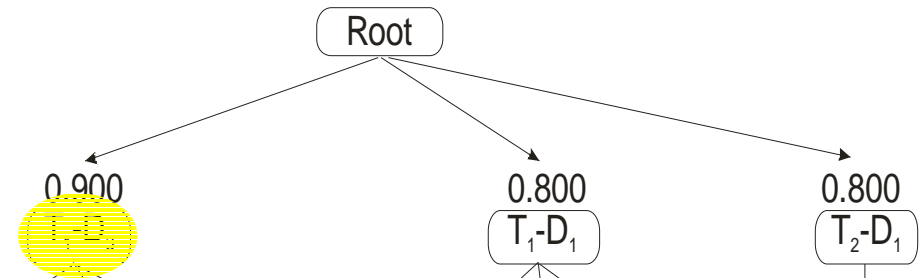
UNCLASSIFIED

Example of Best-Bound with Topology



Data Graph

Template



T ₁		T ₂		T ₃		T ₄	
D ₃	0.9	D ₁	0.8	D ₆	0.77	D ₁₀	0.6
D ₁	0.8	D ₅	0.75	D ₁	0.75	D ₇	0.45
D ₉	0.75	D ₃	0.71	D ₂	0.63	D ₂	0.4
D ₆	0.73	D ₇	0.68	D ₄	0.52	D ₆	0.35
D ₅	0.7	D ₉	0.55	D ₈	0.44	D ₄	0.28
D ₄	0.65	D ₆	0.53	D ₉	0.27	D ₈	0.18
D ₂	0.4	D ₄	0.23	D ₃	0.23	D ₁	0.15
D ₇	0.17	D ₂	0.21	D ₇	0.17	D ₃	0.14
D ₁₀	0.15	D ₁₀	0.13	D ₁₁	0.11	D ₅	0.12
D ₁₁	0.11	D ₈	0.09	D ₅	0.09	D ₉	0.1
D ₈	0.09	D ₁₁	0.05	D ₁₀	0.05	D ₁₁	0.07

1-Hop Neighbor Optimal Assignments

Design of experiment

- But we still don't know what are the *optimal parameters* to be set.
 - Responses: *Runtime* and *Maximum Heuristic Score*.
 - 2-level factorial design: 9 {512 (2^9) runs}
 - Screening experiment: 1/4th factorial design with 128 runs.
 - depth $\delta = 6$

Factors	Levels		Range	
Number of nodes in data graph	25	50	-	-
Number of nodes in template graph	6	10	-	-
Data graph density	10%	15%	0%	100%
Template graph density	24%	50%	0%	100%
t	0.1	1	0	1
α	0.1	1	0	1
k_0	50	100	0	mn
k_i	50	100	0	mn
β	50	300	0	$\prod_i k_i$

m = The number of data graph nodes.

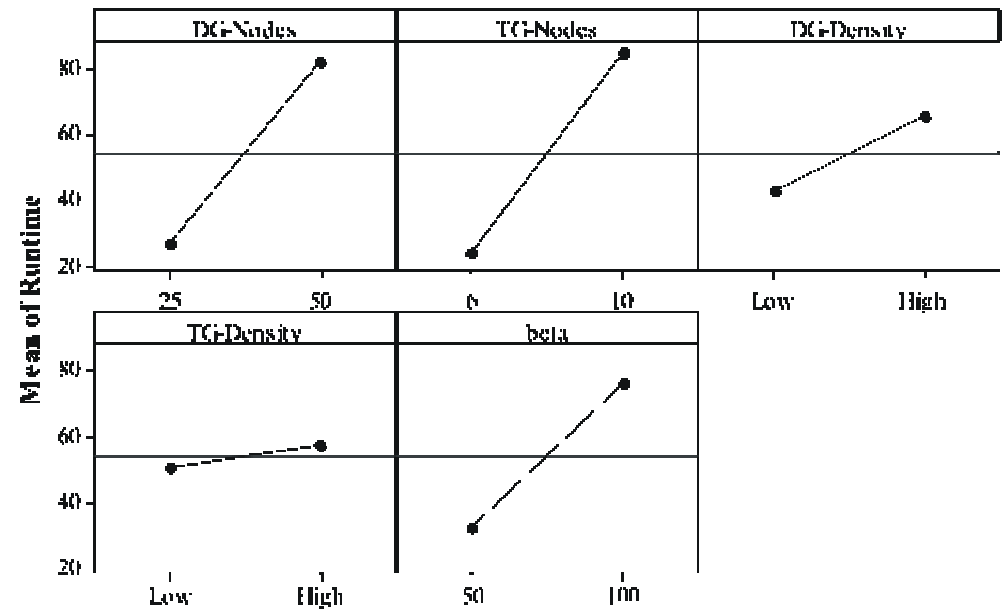
n = The number of data graph

Design of experiment

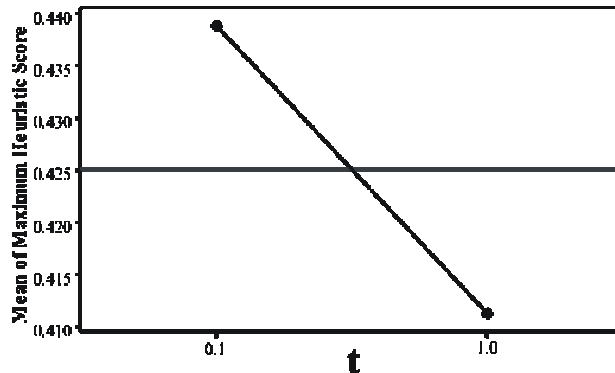
INTUITIVE CONCLUSIONS

- Template graph density has *less effect* on runtime as compared to data graph density.
- β dominates all other parameters for runtime.
- The constant α are much more significant in obtaining quality of the heuristic.

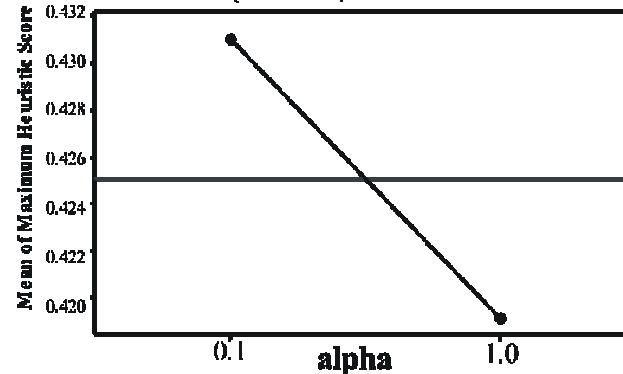
Main Effects Plot (data means) for Runtime



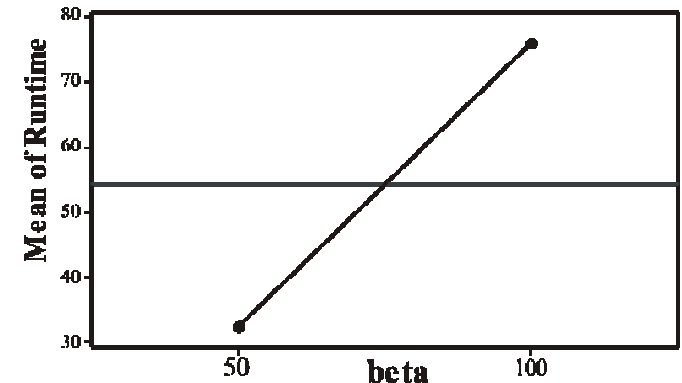
Main Effects Plot (data means) for Maximum Heuristic Score



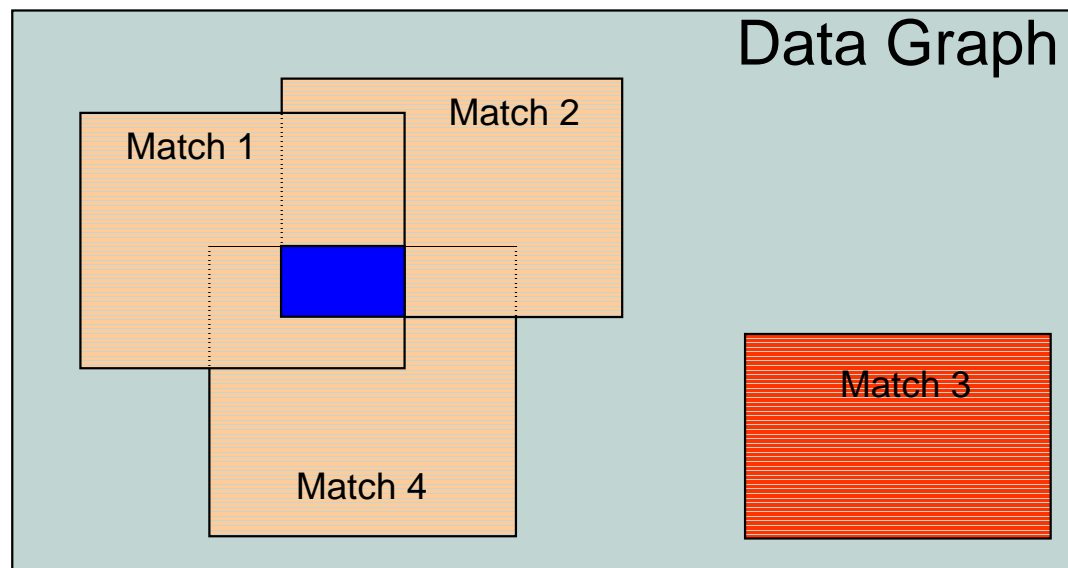
Main Effects Plot (data means) for Maximum Heuristic Score



Main Effects Plot (data means) for Runtime



- TruST results in *large number* of matches.
- What do the results tell us?



- To analyze the results we group the results using K means clustering

K-Means Clustering: Hypercube vs. Fuzzy Hamming Distance

- Hypercube Distance
 - Result becomes a point M in the N dimensional hypercube
 - The metric to minimize distance will determine the shape of the optimum clusters.

$$(\Delta(A, B))$$

— *Hypercube distance measure*

$$\Delta(A, B) = 1 - \frac{\mu(A \cap B)}{\mu(A \cup B)}$$

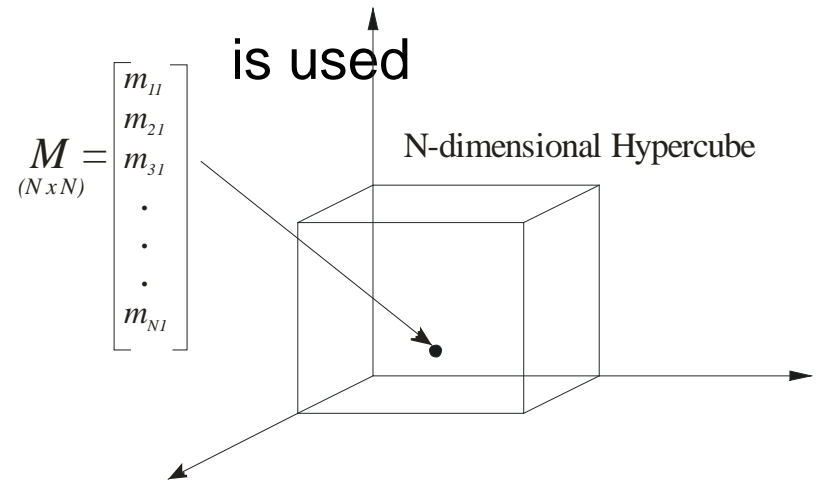
where,

$$0 \leq \Delta(A, B) \leq 1;$$

$$\mu(A) = \sum_{i=1}^n m_A(y_i);$$

$A \cap B$ = Element - wise minima of A and B

$A \cup B$ = Element - wise maxima of A and B



- Fuzzy Hamming Distance

• *Fuzzy Cardinality* $CardA$ of a fuzzy set $A = \{x_i, \mu_i\}$

$$CardA = \sum_{i=0}^n \frac{i}{\mu_{CardA}(i)}$$

where $\mu_{CardA}(i) = \min(\mu_i, (1 - \mu_{i+1}))$. μ_i denotes the i^{th} largest value of μ .

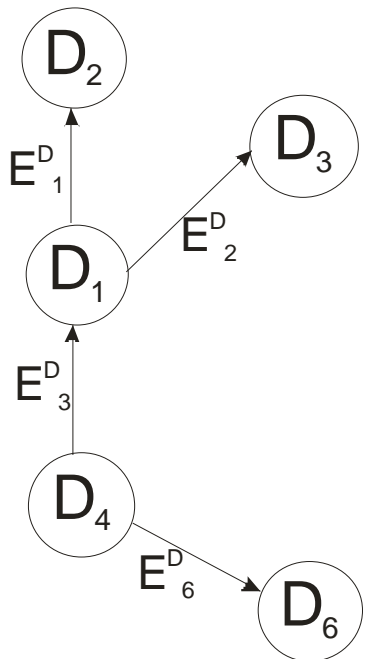
Compared using Silhouette validation technique and Mann-Whitney test and **Hypercube Distance** performed better.

- Aggregate the matches to find the correlated nodes & edges in the formed clusters.
- The matches are combined using **union** and intersection operation

D_1	D_1
E^D_2	E^D_1
D_3	D_2
E^D_3	E^D_3
D_4	D_4
E^D_6	E^D_6
D_6	D_6

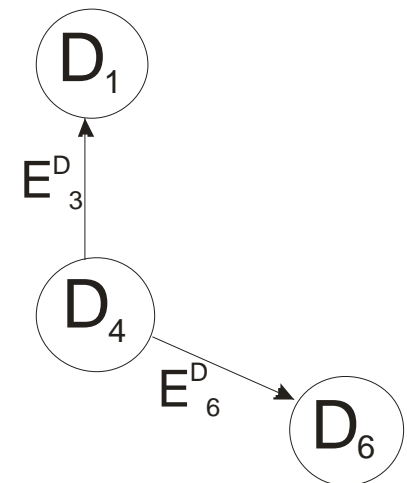
Diversification (Union)

Intensification (Intersection)

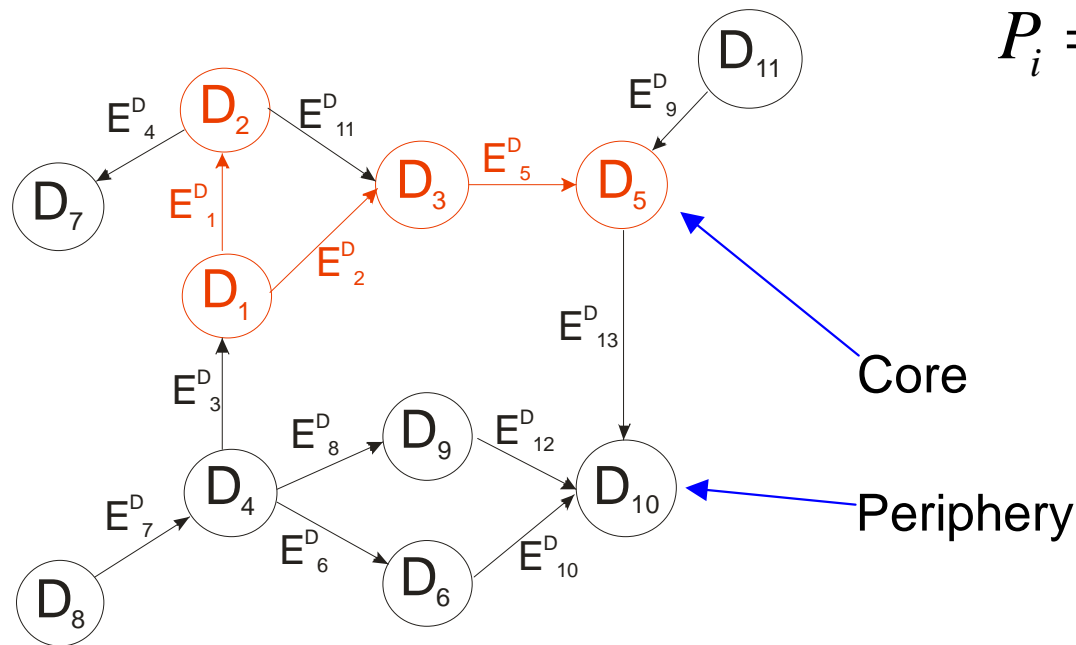


Nodes	Edges
D_1	E^D_1
D_2	E^D_2
D_3	E^D_3
D_4	E^D_6
D_5	
D_6	

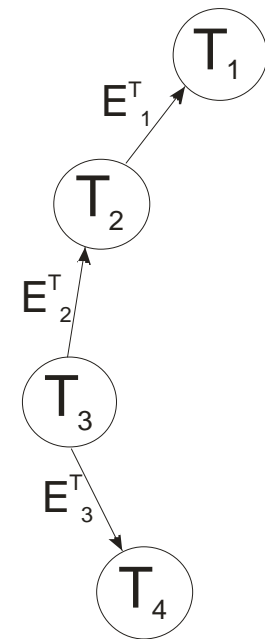
Nodes	Edges
D_1	E^D_3
D_4	E^D_6
D_6	

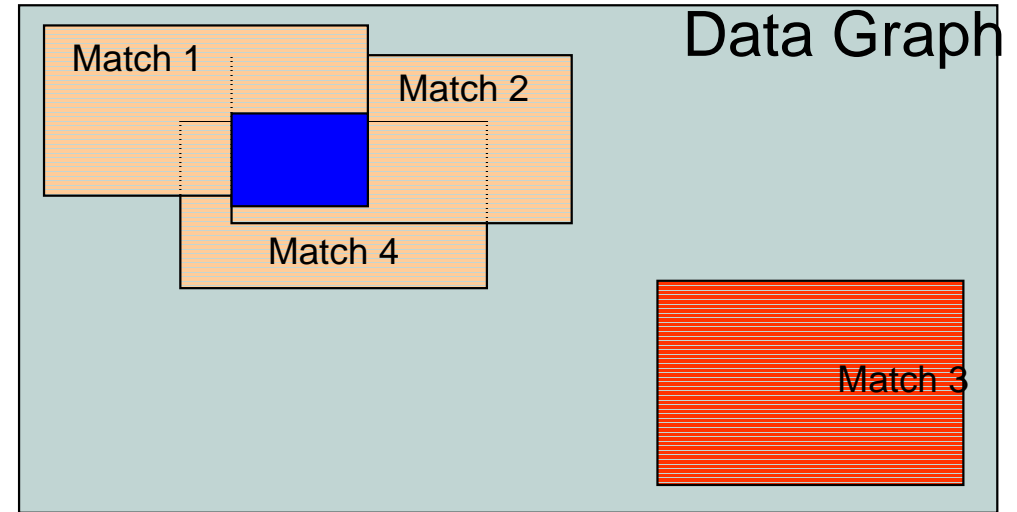
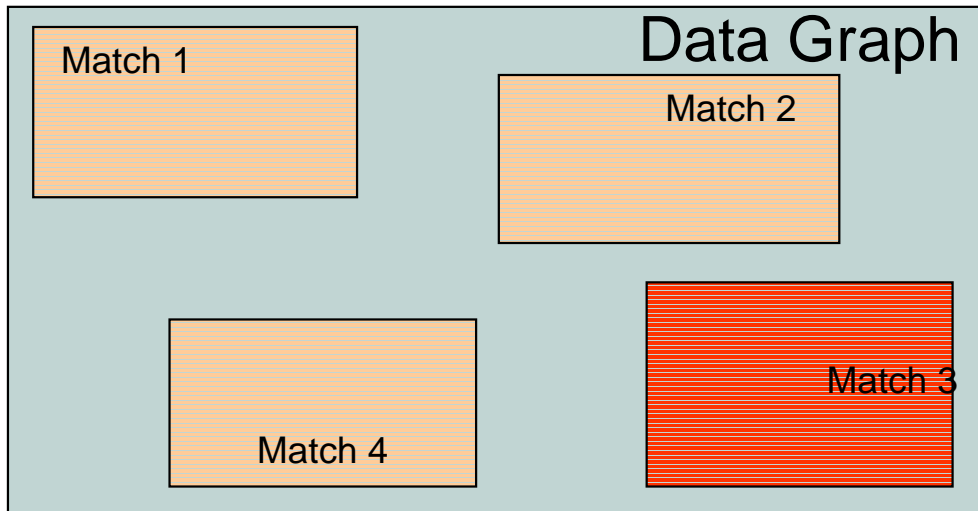


- P_i is the periphery node score
- N_T is the number of nodes in the match
- d_{ij} is the shortest path distance (*Floyd-Marshall*) between periphery node i and core node j



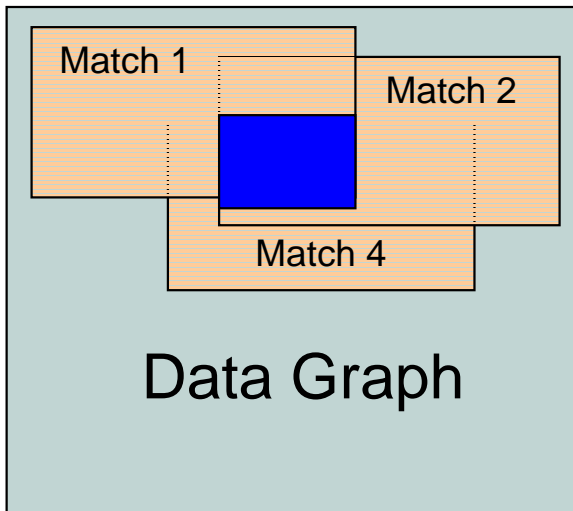
$$P_i = \frac{\sum_j d_{ij}}{N_T}$$



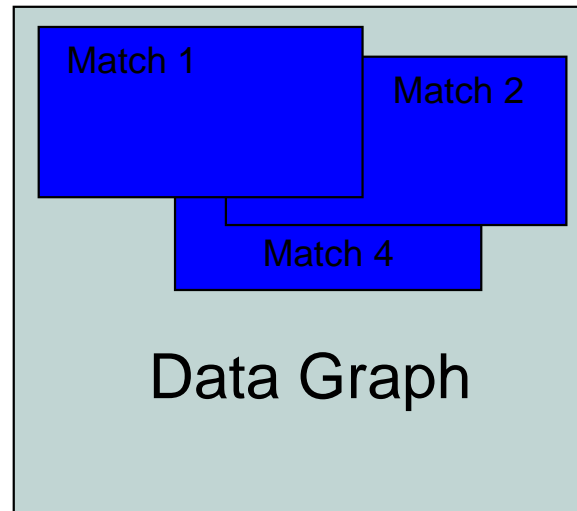


Ranking of Patterns of Interest

Clustering Patterns of Interest

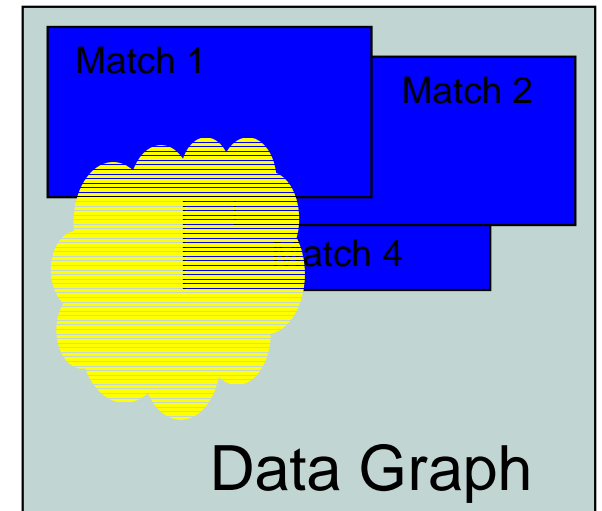


Intensification



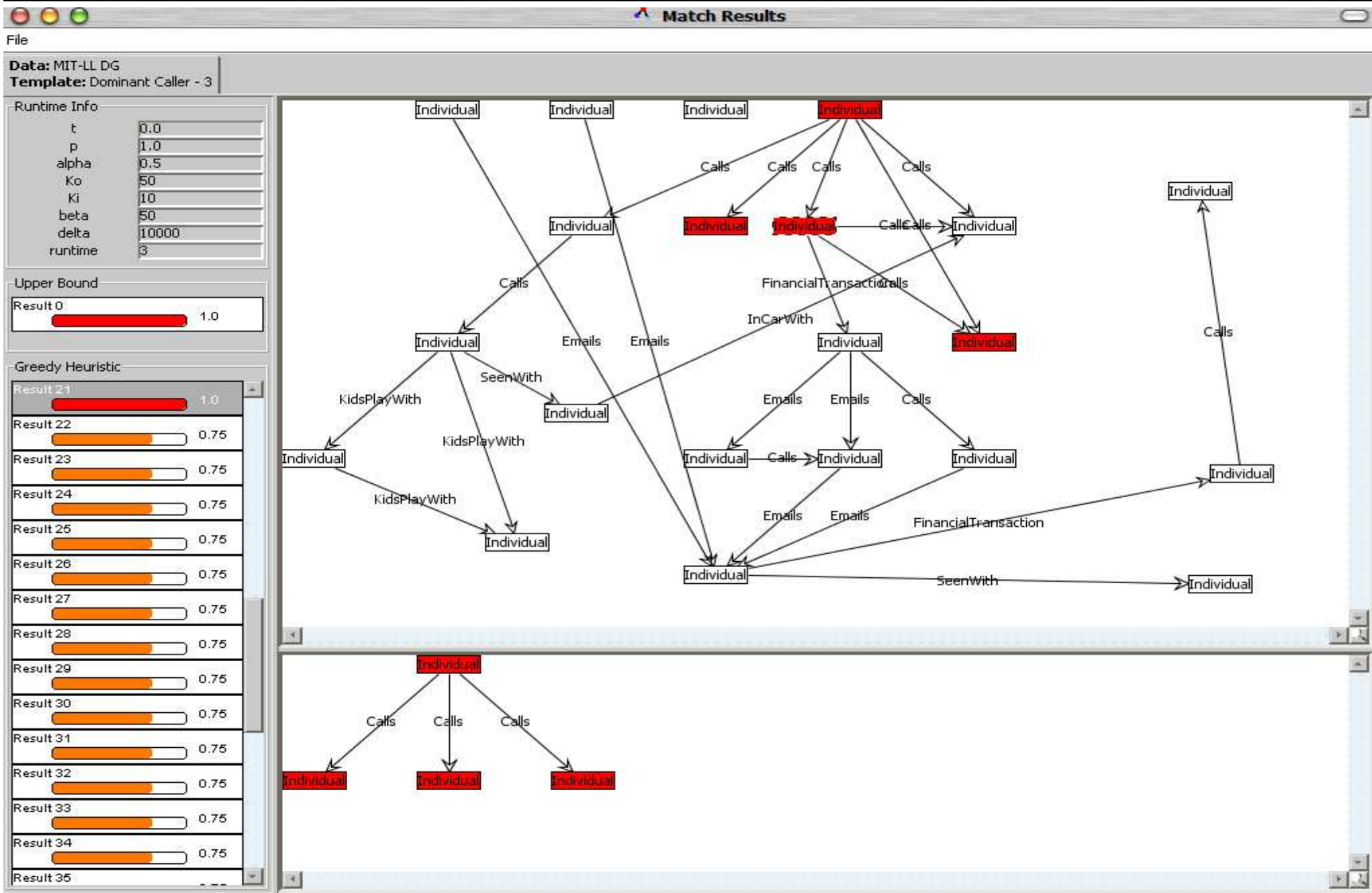
Diversification

UNCLASSIFIED

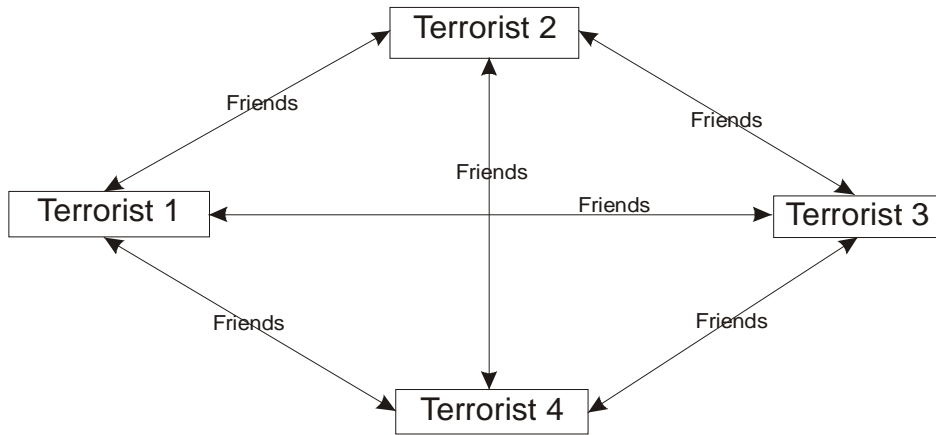


Pattern Discovery

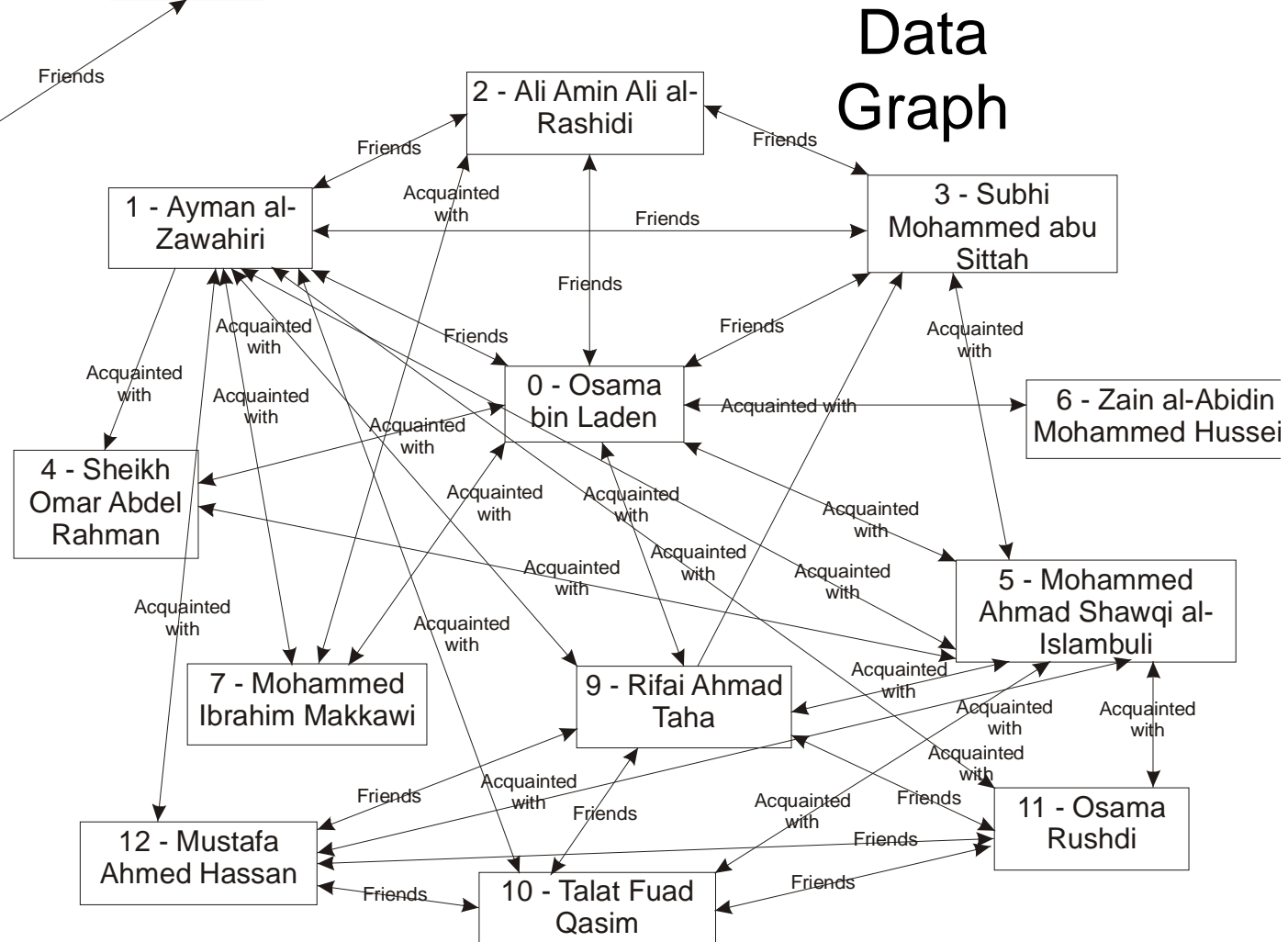
TruST Implementation



Example with Marc Sageman Data



Template Graph

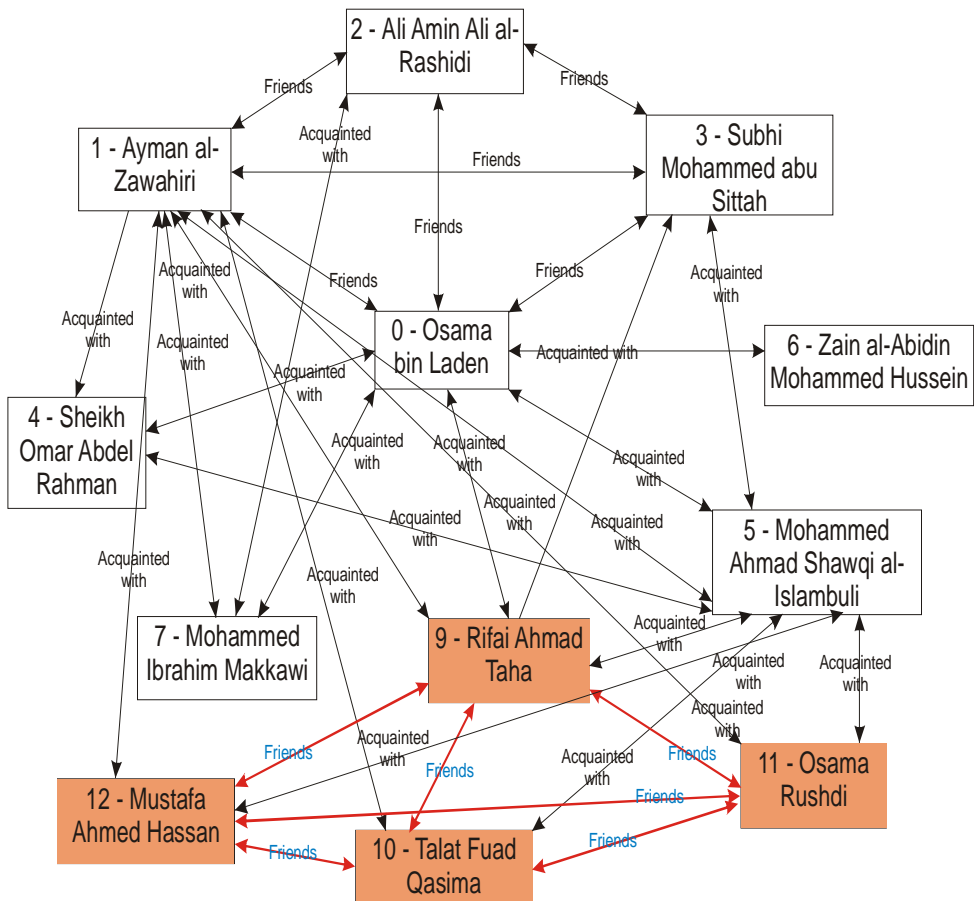


Data Graph

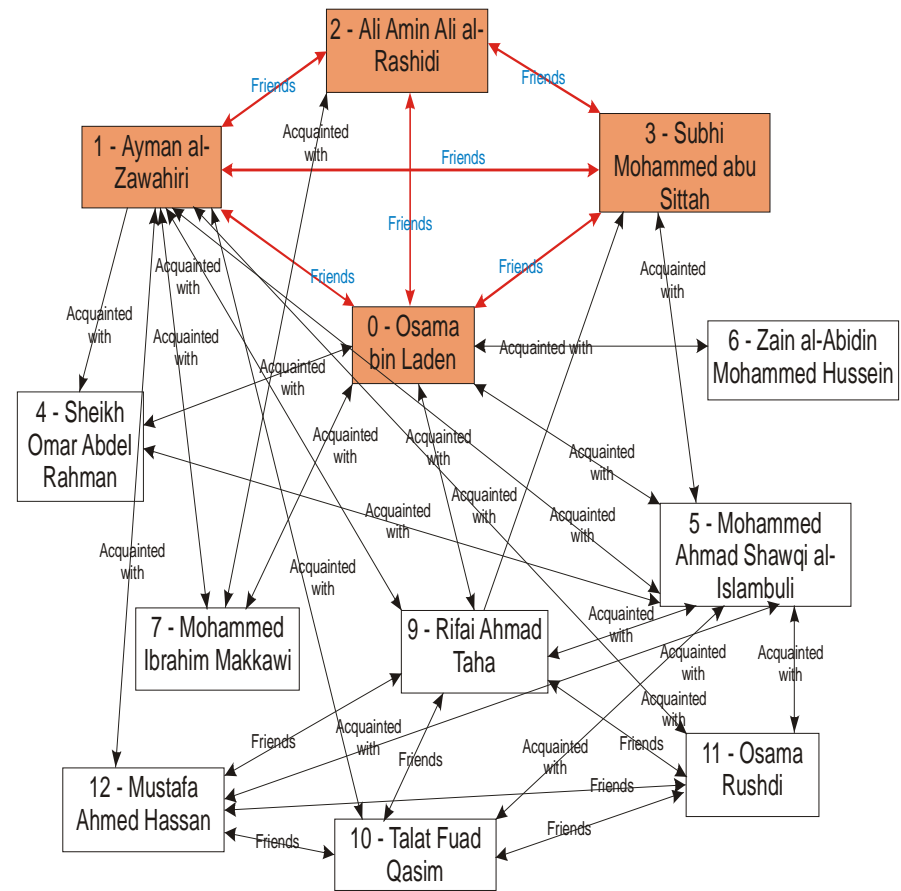
UNCLASSIFIED

TruST matches

Match 1

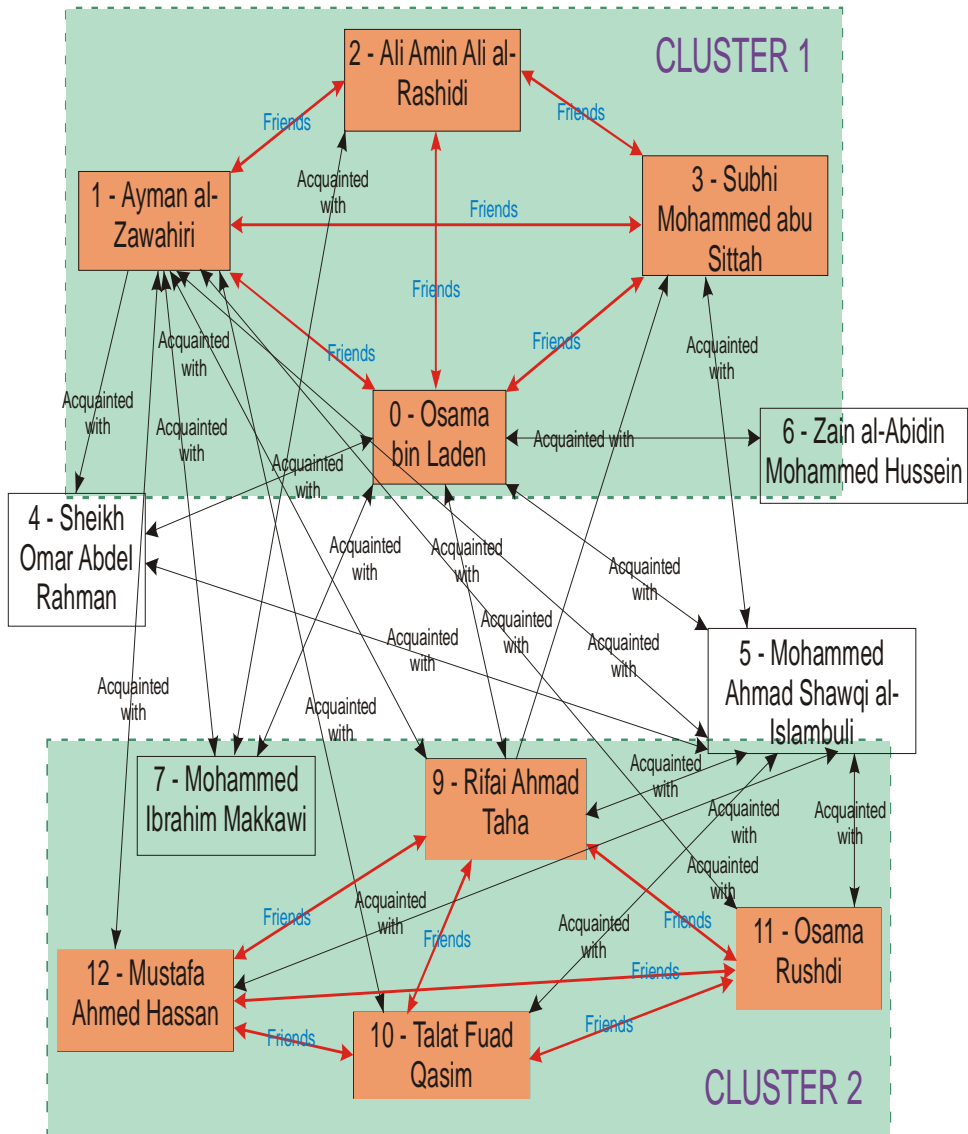


Match 2

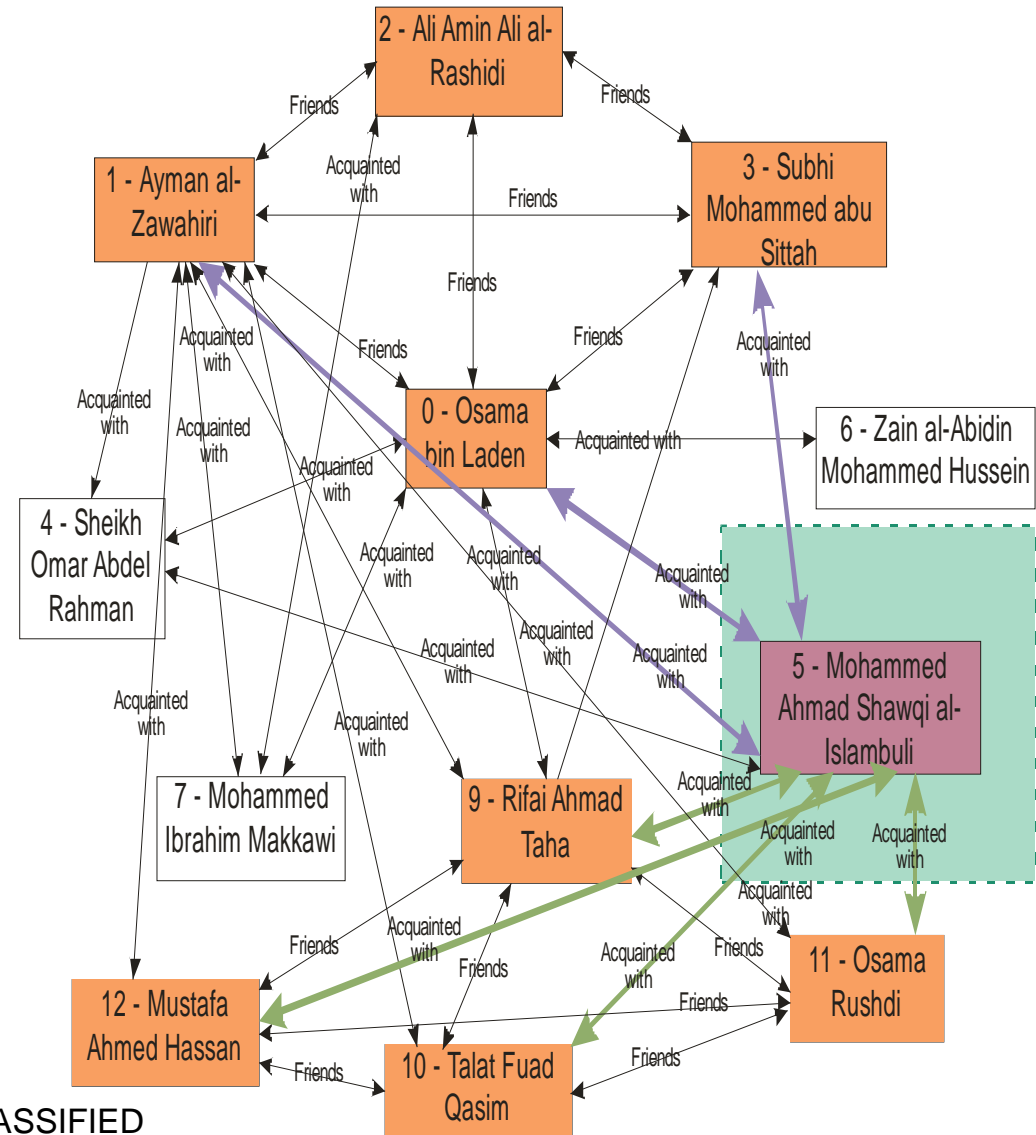


Clustering and Neighborhood Information

Hypercube Distance Clustering



Discovery of Neighborhood Information



- Millions of possible Templates (Hypothesis) of interest
 - Information Fusion Engine for Real Time Decision-making (**INFERD**)
- Ad-hoc choice of **TruST** parameters
 - We have done design of experiments to bound the best choice
 - We are working on trying to Characterize problem structure
- Graph Matching requires initial similarity values
 - We are working on automated process using the semantic features of the nodes and edges (**Conceptual Spaces**)
- Graph Matching could take to long for some domains
 - We span temporal decision-making process with **INFERD** and **TruST**
- Performance on multiple domains
 - Implemented in Asymmetric Warfare, Chem/Bio Warfare, Maritime Domain, Cyber Security, Sensor Management COA, etc.