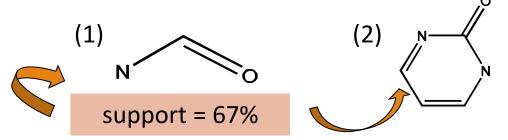


Frequent (Sub)Graph Patterns

- Given a labeled graph dataset D = {G₁, G₂, ..., G_n), the supporting graph set of a subgraph g is D_g = {G_i | $g \subseteq G_i$, G_i \in D}
 - \square support(g) = $|D_g|/|D|$
- \triangle A (sub)graph g is **frequent** if support(g) \ge min_sup
- Ex.: Chemical structures
- Alternative:
 - Mining frequent subgraph patterns from a single large graph or network

 $min_sup = 2$

Frequent Graph Patterns

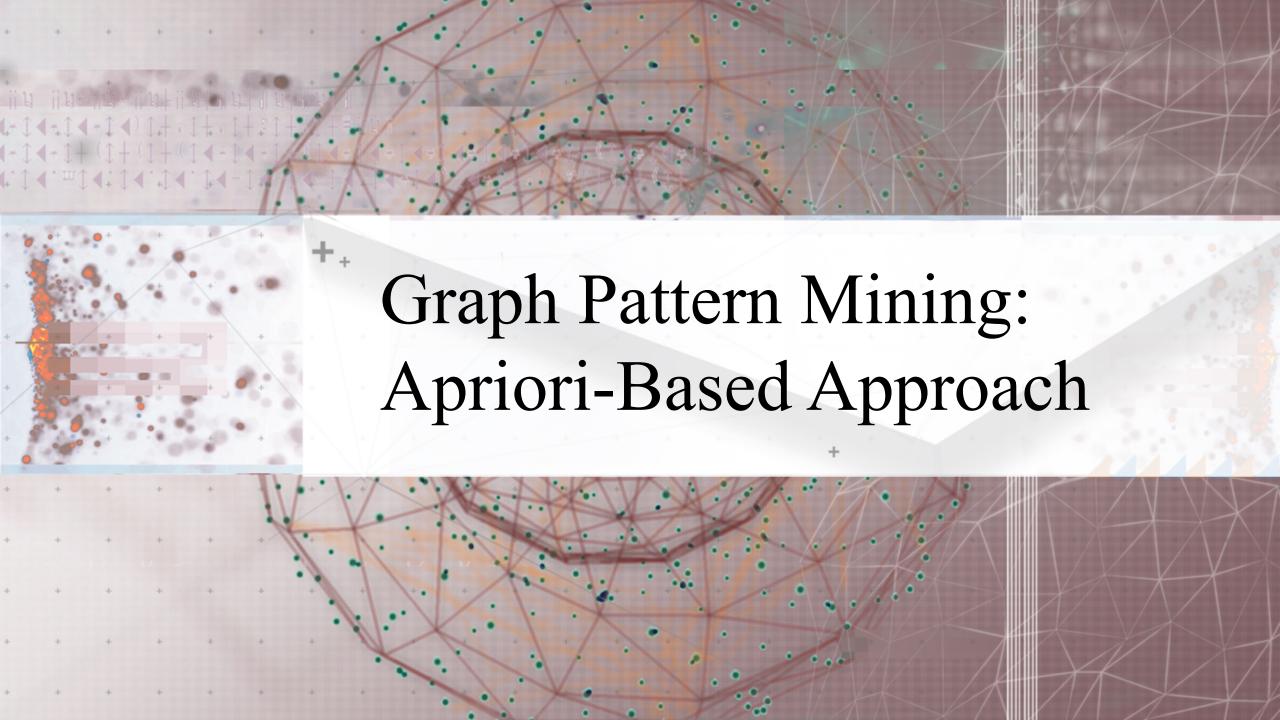


Applications of Graph Pattern Mining

- Bioinformatics
 - Gene networks, protein interactions, metabolic pathways
- Chem-informatics: Mining chemical compound structures
- Social networks, web communities, tweets, ...
- Cell phone networks, computer networks, ...
- Web graphs, XML structures, Semantic Web, information networks
- Software engineering: Program execution flow analysis
- Building blocks for graph classification, clustering, compression, comparison, and correlation analysis
- Graph indexing and graph similarity search

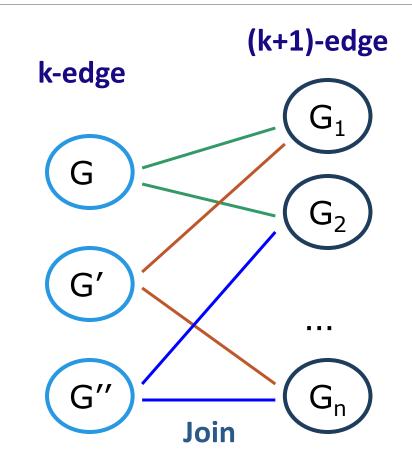
Graph Pattern Mining Algorithms: Different Methodologies

- Generation of candidate subgraphs
 - Apriori vs. pattern growth (e.g., FSG vs. gSpan)
- Search order
 - Breadth vs. depth
- Elimination of duplicate subgraphs
 - Passive vs. active (e.g., gSpan [Yan & Han, 2002])
- Support calculation
 - Store embeddings (e.g., GASTON [Nijssen & Kok, 2004], FFSM [Huan, Wang, & Prins, 2003], MoFa [Borgelt & Berthold, ICDM'02])
- Order of pattern discovery
 - □ Path → tree → graph (e.g., GASTON [Nijssen & Kok, 2004])



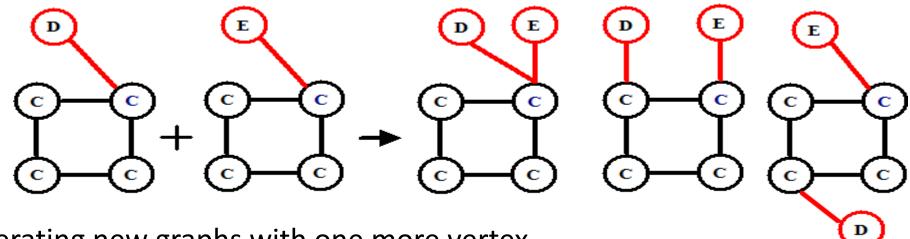
Apriori-Based Approach

- □ The Apriori property (anti-monotonicity): A size-k subgraph is frequent if and only if all of its subgraphs are frequent
- □ A candidate size-(k+1) edge/vertex subgraph is generated if its corresponding two k-edge/vertex subgraphs are frequent
- Iterative mining process:
 - □ Candidate-generation → candidate pruning → support counting → candidate elimination



Candidate Generation: Vertex Growing vs. Edge Growing

- ☐ Methodology: Breadth-search, Apriori joining two size-k graphs
 - \square Many possibilities at generating size-(k+1) candidate graphs

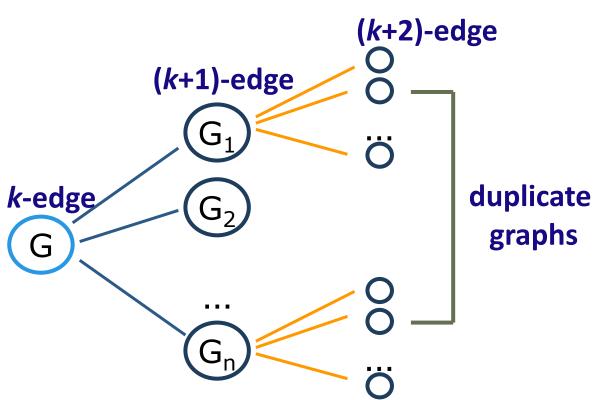


- ☐ Generating new graphs with one more vertex
 - AGM (Inokuchi, Washio, & Motoda, PKDD'00)
- Generating new graphs with one more edge
 - FSG (Kuramochi & Karypis, ICDM'01)
- ☐ Performance shows *via edge growing* is more efficient



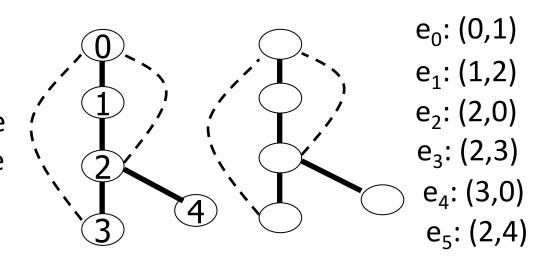
Pattern-Growth Approach

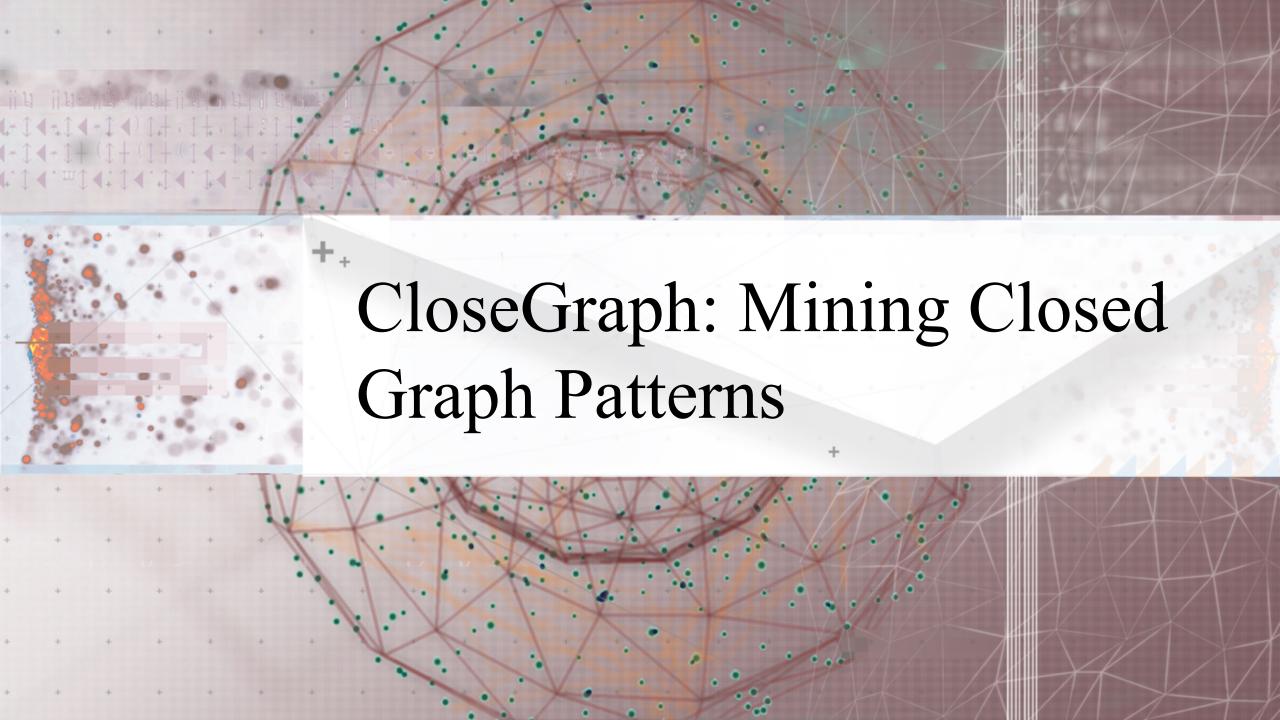
- □ Depth-first growth of subgraphs from k-edge to (k+1)-edge, then (k+2)-edge subgraphs
- Major challenge
 - Generating many duplicate subgraphs
- Major idea to solve the problem
 - Define an order to generate subgraphs
 - DFS spanning tree: Flatten a graph into a sequence using depth-first search
 - gSpan (Yan & Han, ICDM'02)



gSPAN: Graph Pattern Growth in Order

- Right-most path extension in subgraph pattern growth
 - Right-most path: The path from root to the right-most leaf (choose the vertex with the smallest index at each step)
 - Reduce generation of duplicate subgraphs
- Completeness: The enumeration of graphs using right-most path extension is <u>complete</u>
- DFS code: Flatten a graph into a sequence using depth-first search





Why Mine Closed Graph Patterns?

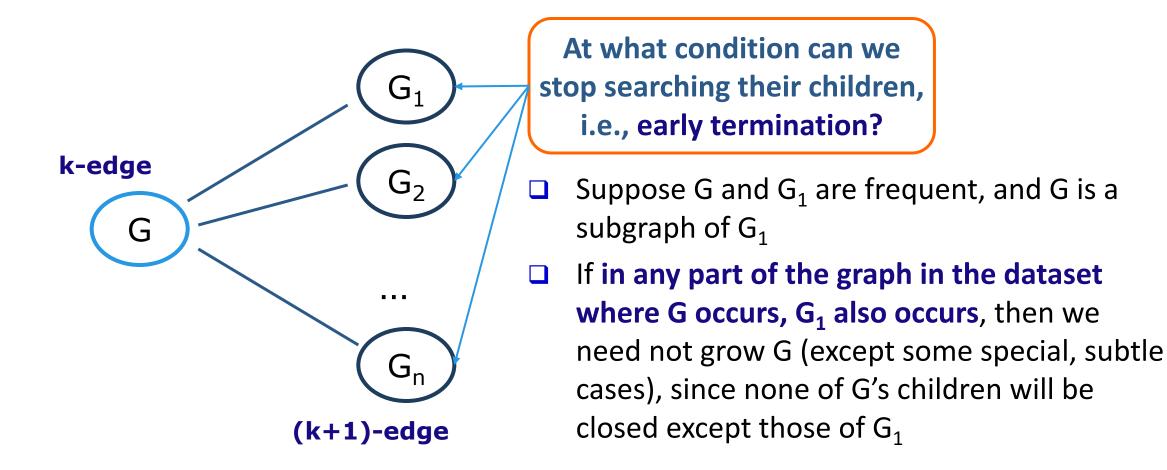
- □ Challenge: An **n**-edge frequent graph may have 2ⁿ subgraphs
- Motivation: Explore closed frequent subgraphs to handle graph pattern explosion problem
- □ A frequent graph G is *closed* if there exists no supergraph of G that carries the same support as G

If this subgraph is *closed* in the graph dataset, it implies that none of its frequent super-graphs carries the same support

- Lossless compression: Does not contain non-closed graphs, but still ensures that the mining result is complete
- Algorithm CloseGraph: Mines closed graph patterns directly

CloseGraph: Directly Mining Closed Graph Patterns

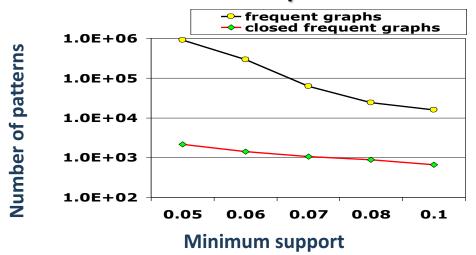
CloseGraph: Mining closed graph patterns by extending gSpan (Yan & Han, KDD'03)



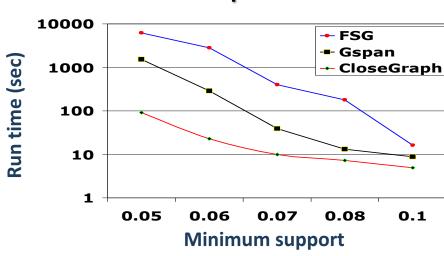
Experiment and Performance Comparison

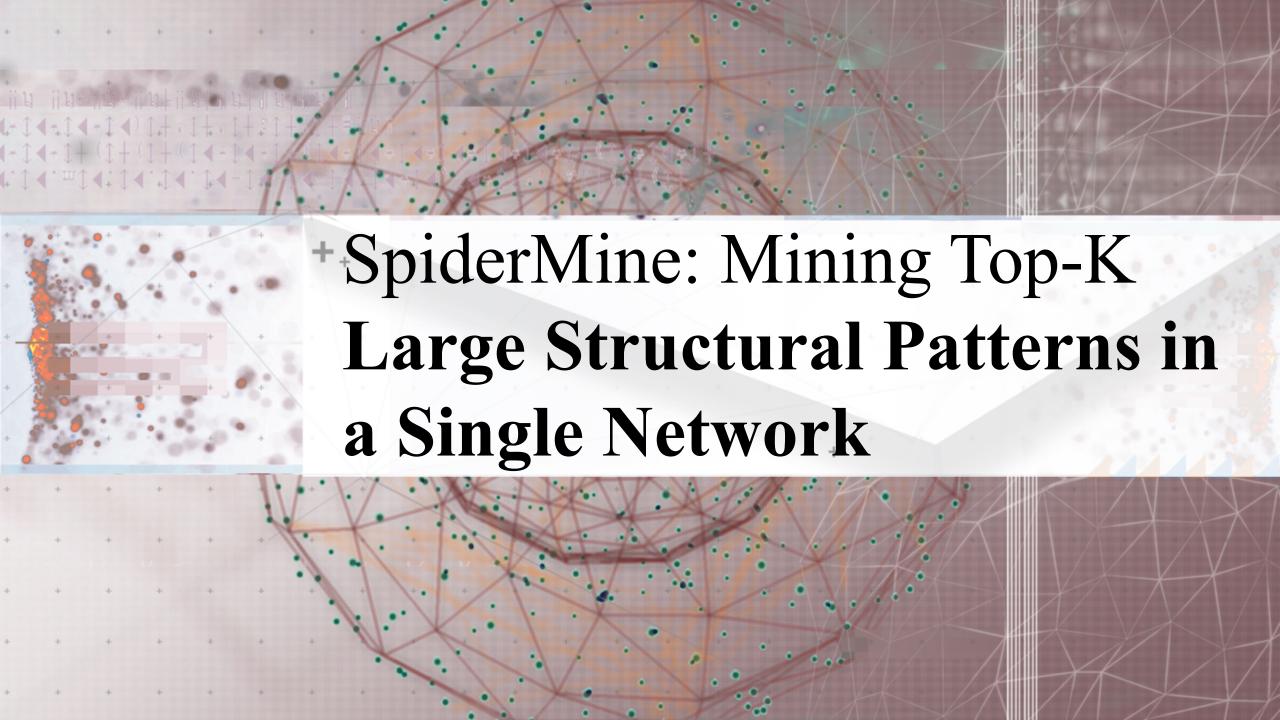
- The AIDS antiviral screen compound dataset from NCI/NIH
- ☐ The dataset contains 43,905 chemical compounds
- Discovered patterns: The smaller minimum support, the bigger and more interesting subgraph patterns discovered

of Patterns: Frequent vs. Closed



Runtime: Frequent vs. Closed





SpiderMine: Mining Top-K Large Structural Patterns in a Massive Network

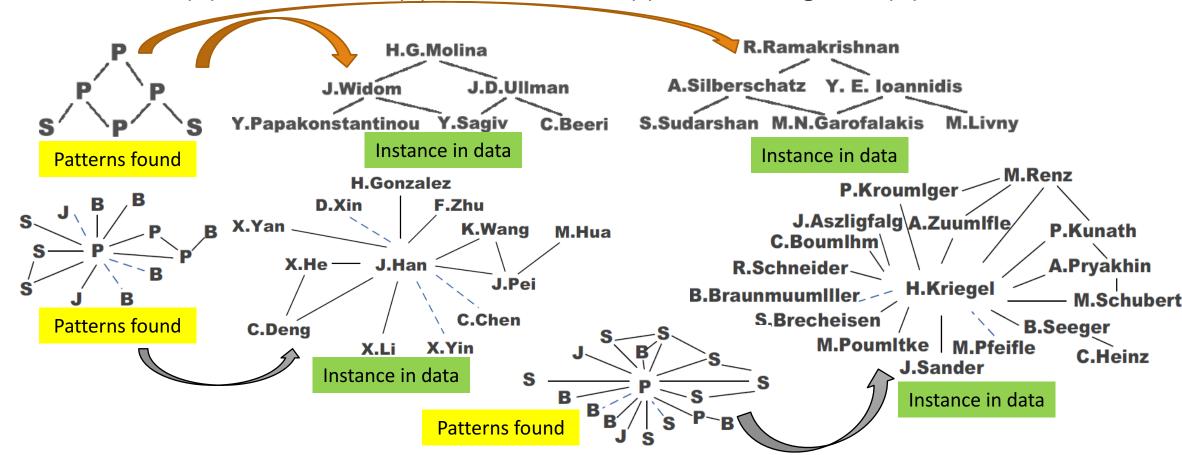
- □ Large patterns are informative to characterize a large network (e.g., social network, web, or bio-network)
- Similar to pattern fusion, mining large patterns should not aim for completeness but for representativeness of the target results
- SpiderMine (Zhu et al., VLDB'11): Mine top-K largest frequent substructure patterns whose diameter is bounded by D_{max} with a probability at least $1-\epsilon$
- General idea: Large patterns are composed of a number of small components ("spiders"), which will eventually connect together after some rounds of pattern growth
- □ **r-Spider:** An r-spider is a frequent graph pattern P such that there exists a vertex u of P, and all other vertices of P are within distance r from u

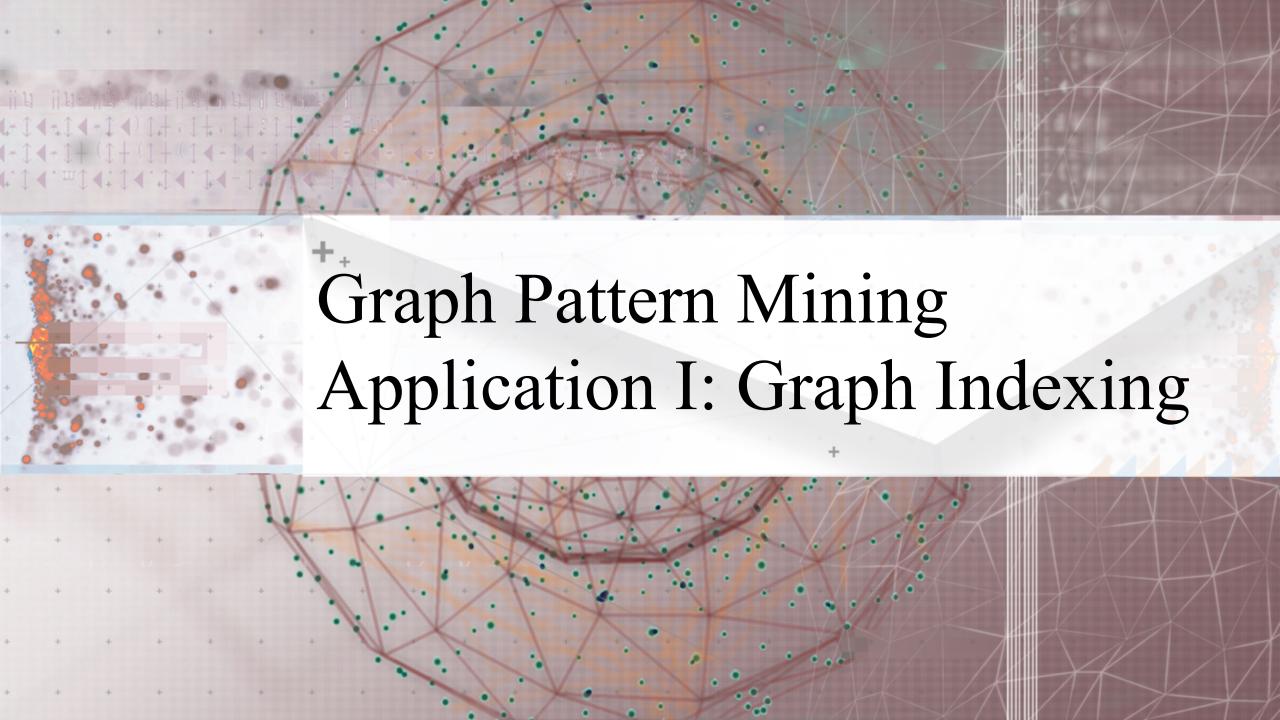
Why Is SpiderMine Good for Mining Large Patterns?

- The SpiderMine algorithm
 - Mine the set S of all the r-spiders
 - Randomly draw M r-spiders
 - Grow these M r-spiders for $t = D_{max}/2$ iterations, and merge two patterns whenever possible
 - Discard unmerged patterns
 - Continue to grow the remaining ones to maximum size
 - Return the top-K largest ones in the result
- Why is SpiderMine likely to retain large patterns and prune small ones?
 - Small patterns are much less likely to be hit in the random draw
 - Even if a small pattern is hit, it is even less likely to be hit multiple times
 - ☐ The larger the pattern, the greater the chance it is hit and saved

Mining Collaboration Patterns in DBLP Networks

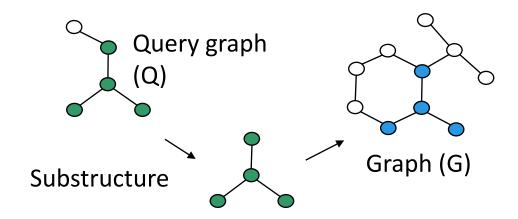
- □ Data description: 600 conferences, 9 major CS areas, 15,071 authors in DB/DM
- Author labeled by # of papers published in DB/DM
 - Prolific (P): \geq 50, Senior (S): 20~49, Junior (J): 10~19, Beginner(B): 5~9

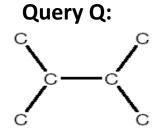




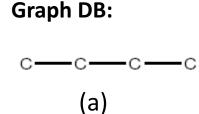
Application of Pattern Mining I: Graph Indexing

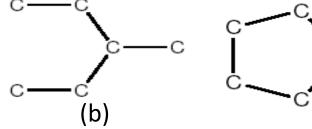
- ☐ Graph query: Find all the graphs in a graph DB containing a given query graph
- Index should be a powerful tool
- Path-index may not work well
- Solution: Index directly on substructures (i.e., graphs)

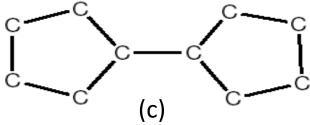




Only graph (c) contains Q



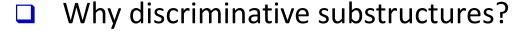




Path-indices: C, C-C, C-C-C, C-C-C cannot prune (a) & (b)

gIndex: Indexing Frequent and Discriminative Substructures

- Why index frequent substructures?
 - Too many substructures to index
 - Size-increasing support threshold
 - Large structures will likely be indexed well by their substructures

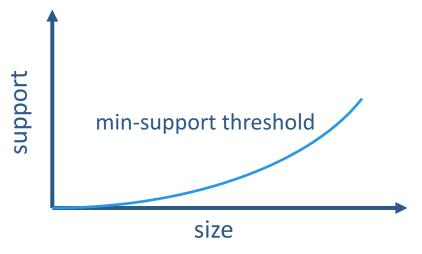


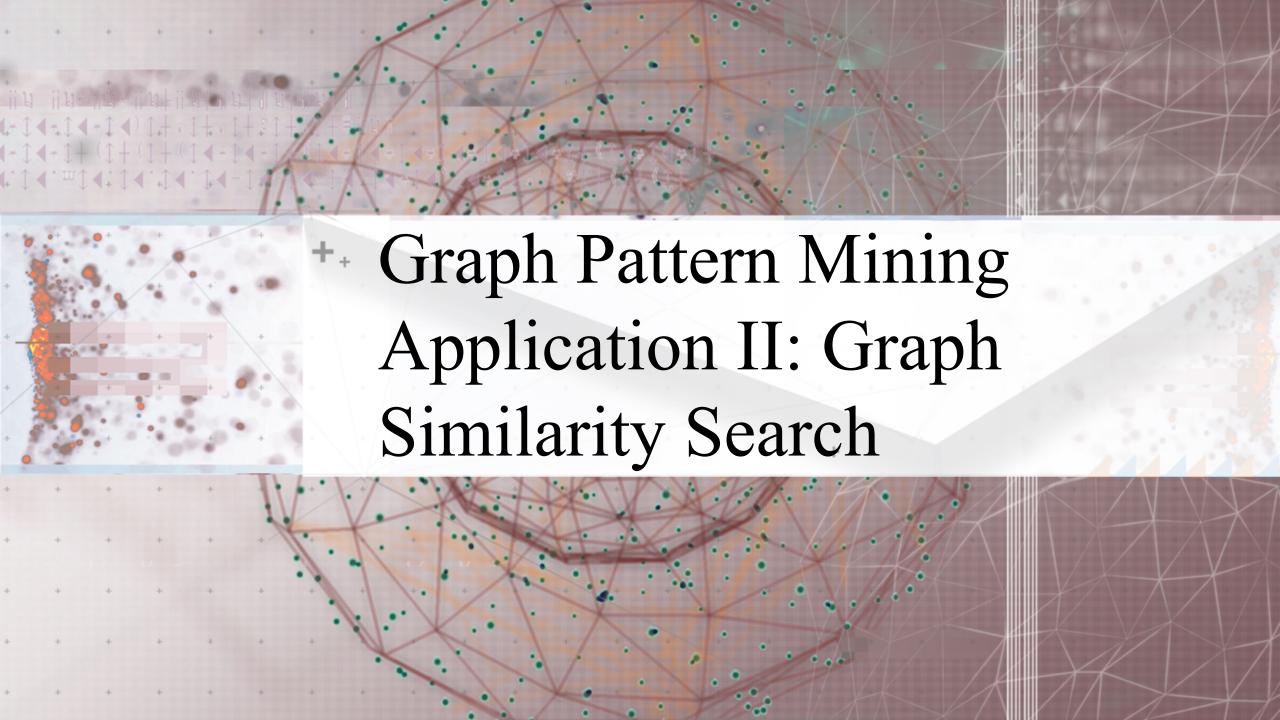
- Reduce the index size by an order of magnitude
- Selection: Given a set of selected structures f_1 , f_2 , ... f_n , and a new structure x, the extra indexing power is measured by

$$\Pr(x|f_1, f_2, \dots f_n), f_i \subset x$$

when $Pr(x|f_1, f_2, ..., f_n)$ is small enough, x is a discriminative structure and should be included in the index

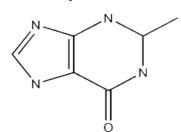
Experiments show that glndex is small, effective, and stable

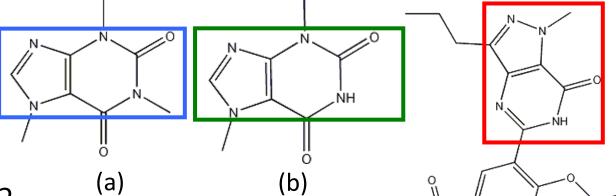




Application II: Support Substructure Similarity Search

- ☐ Find graphs in a graph DB containing substructures similar to a given query graph
- Ex. Data: A chemical compound DB
 - A query graph q:





(c)

- How to do similarity search efficiently?
 - No indexing? Sequential scan + computing subgraph similarity – too costly!
 - Build graph indices to support approximate search?
 - Need an explosive number of subgraphs to cover all the similar subgraphs!
- An elegant solution (Yan, Yu, & Han, SIGMOD'05):
 - Keep the graph index structure, but select features in the query space

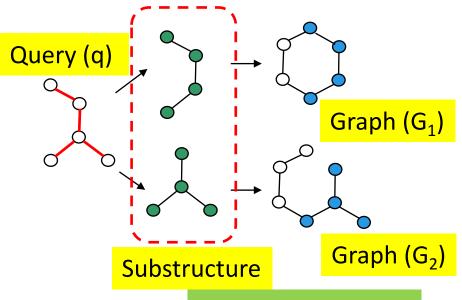
Feature-Based Similarity Search

- Decompose a query graph into a set of features
- Feature-based similarity measure
 - Each graph is represented as a feature vector $X = \{x_1, x_2, ..., x_n\}$
 - Similarity is defined by the distance of their corresponding vectors
- If graph G contains the major part of a query graph q,
 G should share a number of common features with q
 - Given a relaxation ratio, one can calculate the maximal number of features that can be missed!

Assume: Query graph has 5 features

Relaxation threshold: Can miss at most 2 features

Then: G₁, G₂, G₃ are pruned



Graphs in database

	G_1	G_2	G_3	G_4	G_5
f_1	0	1	0	1	1
f ₂	0	1	0	0	1
f_3	1	0	1	1	1
f ₄	1	0	0	0	1
f_5	0	0	1	1	0

