

VENOM: Approximate Subgraph Matching with Enhanced Neighbourhood Structural Information

SHUBHANGI AGARWAL[†]
sagarwal@cse.iitk.ac.in

SOURAV DUTTA^{††}
sourav.dutta2@huawei.com

ARNAB BHATTACHARYA[†]
arnabb@cse.iitk.ac.in

CODS COMAD 2024
IIIT Bangalore, India

[†]Department of Computer Science and Engineering,
Indian Institute of Technology Kanpur, **India**

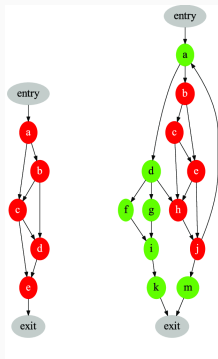
^{††}Huawei Research Centre, Dublin, **Ireland**

Approximate Subgraph Matching (ASM)

Objective: Search in a graph dataset for subgraph *similar* to the pattern of interest.

Approximate Subgraph Matching (ASM)

Objective: Search in a graph dataset for subgraph *similar* to the pattern of interest.



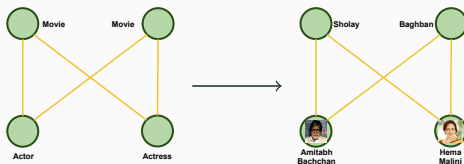
Malware Detection

with Control Flow Graph representations

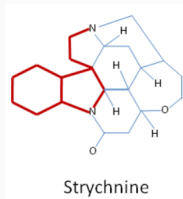
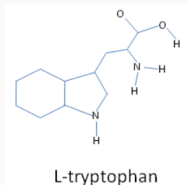
of a Program

Source:

<https://ieeexplore.ieee.org/document/6838703>



Question-Answering in Knowledge Graphs



Compound Analysis

Source: <https://www.cs.tau.ac.il/~roded/sigma.pdf>

Goal of Study

Factors affecting ASM performance:

- Number of nodes
- Density of graph
- Label distribution
- Query size
- \vdots
- Size of neighborhood matched
- Degree distribution

Goal of Study

Factors affecting ASM performance:

- Number of nodes
- Density of graph
- Label distribution
- Query size
- \vdots
- **Size of neighborhood matched**
- Degree distribution

Goal of Study

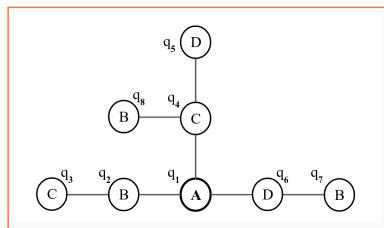
Factors affecting ASM performance:

- Number of nodes
- Density of graph
- Label distribution
- Query size

⋮

- **Size of neighborhood matched**

- Degree distribution



EXAMPLE GRAPH

→ **Depth of neighborhood (h)**

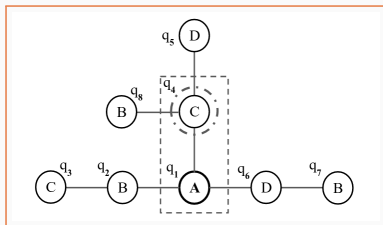
→ **Breadth of neighborhood (k)**

VENoM

Parametrizing size of neighborhood

$h \implies$ depth

$k \implies$ breadth

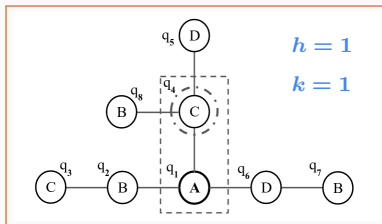


VENoM

Parametrizing size of neighborhood

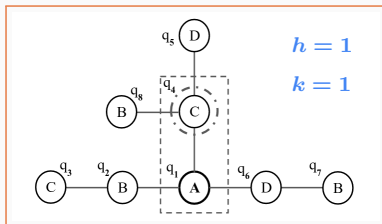
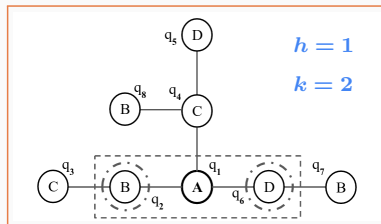
$h \Rightarrow$ depth

$k \Rightarrow$ breadth



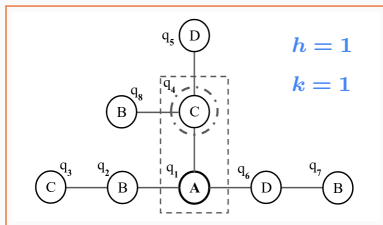
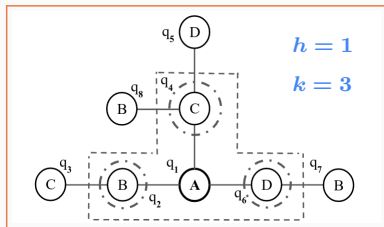
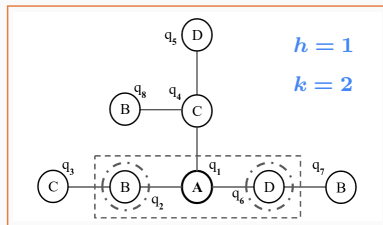
VENOM

Parametrizing size of neighborhood


 $h \Rightarrow \text{depth}$
 $k \Rightarrow \text{breadth}$


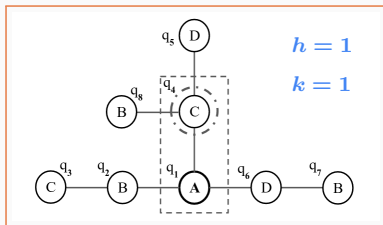
VENOM

Parametrizing size of neighborhood

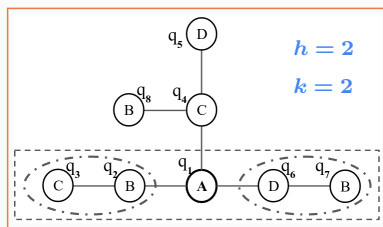
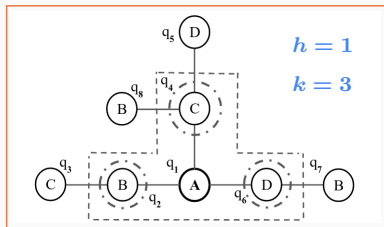
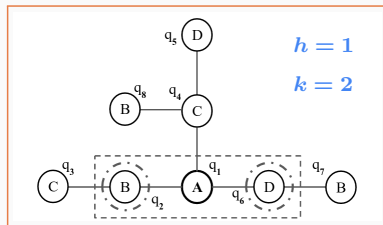

$$h \implies \text{depth}$$
$$k \implies \text{breadth}$$


VENOM

Parametrizing size of neighborhood



$h \implies \text{depth}$

 $k \implies \text{breadth}$ 

Units and Groups

VENoM- (k, h)

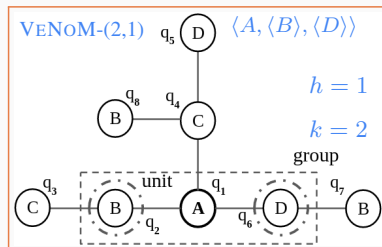
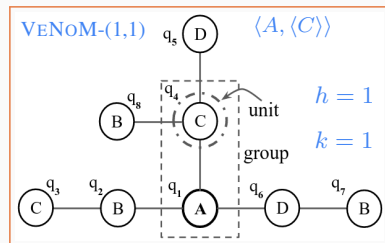
Different instances of VENoM based on breadth (k) and depth (h) of the neighborhood.

Unit

An ordered collection of neighbor labels of a vertex that forms a path of length h .

Group

A set of k units of a vertex along with its label where the units correspond to different neighbors.



Units and Groups

VENoM- (k, h)

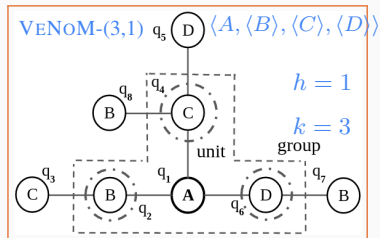
Different instances of VENoM based on breadth (k) and depth (h) of the neighborhood.

Unit

An ordered collection of neighbor labels of a vertex that forms a path of length h .

Group

A set of k units of a vertex along with its label where the units correspond to different neighbors.



Units and Groups

VENoM- (k, h)

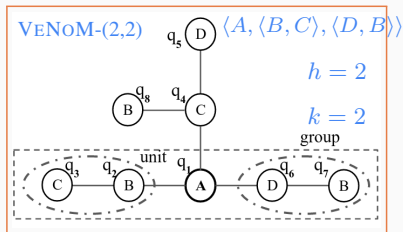
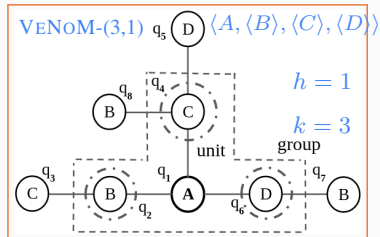
Different instances of VENoM based on breadth (k) and depth (h) of the neighborhood.

Unit

An ordered collection of neighbor labels of a vertex that forms a path of length h .

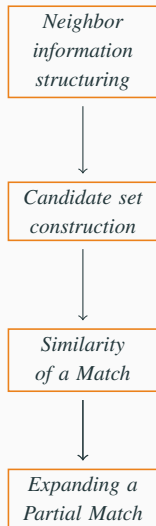
Group

A set of k units of a vertex along with its label where the units correspond to different neighbors.



VENoM

VENoM STAGES



VENoM

VENoM STAGES

*Neighbor
information
structuring*



*Candidate set
construction*



*Similarity
of a Match*

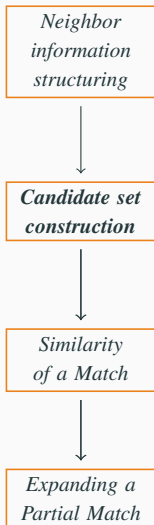


*Expanding a
Partial Match*

1. Number of unique labels in target graph \mathcal{G}
2. Degree of each vertex $v \in \mathcal{G}$
3. Enumerate *units* and *groups* of v

VENoM

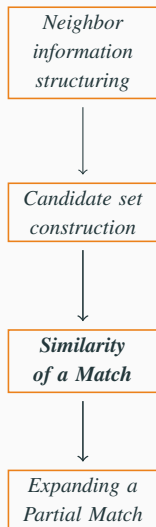
VENoM STAGES



- Candidate vertex: target vertices with label same as query node label
- Label similarity: e.g., Jaccard similarity

VENOM

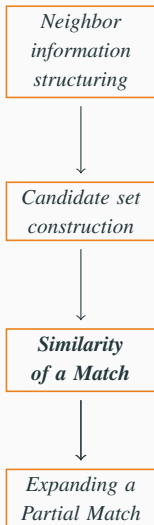
VENOM STAGES



- Goal: find most statistically significant match
- VENOM-(2,2): 5 group match categories
 - Ex: $\langle \underline{A}, \langle B, C \rangle, \langle D, B \rangle \rangle$
 - s_0 : $\langle \underline{A}, \langle \times, \times \rangle, \langle \times, \times \rangle \rangle$
 - s_1 : $\langle \underline{A}, \langle B, \times \rangle, \langle \times, \times \rangle \rangle$
 - s_2 : $\langle \underline{A}, \langle B, \times \rangle, \langle \times, B \rangle \rangle, \langle \underline{A}, \langle B, C \rangle, \langle \times, \times \rangle \rangle$
 - s_3 : $\langle \underline{A}, \langle B, \times \rangle, \langle D, B \rangle \rangle$
 - s_4 : $\langle \underline{A}, \langle B, C \rangle, \langle D, B \rangle \rangle$
- $s_i \succ s_j$, if $i > j$
- Enumerate groups and compare
 - Ordered comparison based on label overlap
 - Categorize group matches (s_i)
 - Large overlap preferred

VENoM

VENoM STAGES



- Similarity metric
 - Deviation of observed behavior from expected
 - Pearson χ^2 statistic

$$\chi^2 = \sum_i \frac{[O(s_i) - E(s_i)]^2}{E(s_i)}$$

- Random vertex pair \rightarrow Low similarity
 - Higher deviation \implies Exceptional similarity
-
- Expected behavior: based on
 - Number of unique labels
 - Degree of the node
 - Number of groups of query node

VENoM

VENoM STAGES

*Neighbor
information
structuring*



*Candidate set
construction*



*Similarity
of a Match*



*Expanding a
Partial Match*

1. Greedy expansion
2. Seed: candidate pair with highest χ^2
3. Prefer unmatched neighbors with highest χ^2
4. Repeat
 - Until terminal condition met

Setup

Datasets

REAL-WORLD DATASETS

Dataset	#Vertices	#Edges	#Labels
Human	4.6K	86.2K	44
HPRD	9.4K	37K	307
Flickr	80.5K	5.9M	195
PPI	12K	10.74M	2.4K

SYNTHETIC BARABÁSI-ALBERT GRAPHS

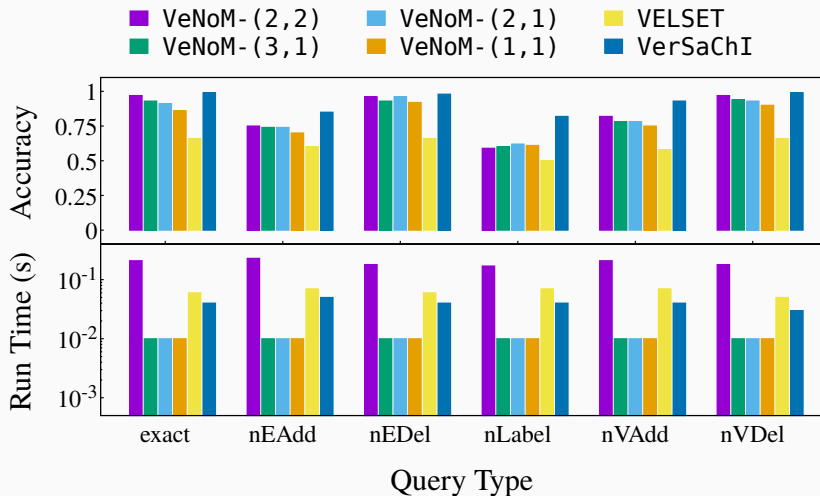
(200 DIFFERENT GRAPHS)

(#Vertices,	Density,	#Labels)
• 2K	• 4	• 2
• 10K	• 10	• 5
• 50K	• 40	• 10
• 250K	• 100	• 25
• 1000K	• 250	• 50
		• 150
		• 500
		• 5K

Query set

- Random connected subgraph extraction
- Perturbations for noisy set

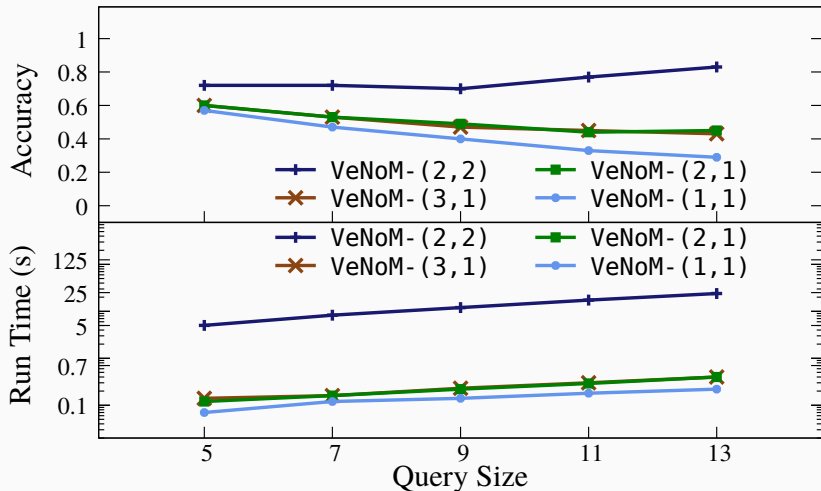
Real Graphs - HPRD



VELSET - Dutta *et al*, WWW 2017

VERSAChI- Agarwal *et al*, CIKM 2021

BA Graphs - Query Size

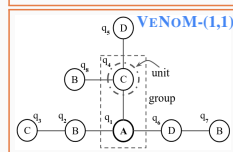
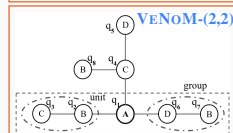
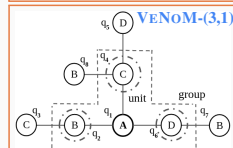
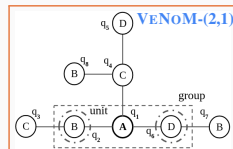


Results

- Increase in graph size and density
 - Accuracy decreases with increase in runtime
- Increase in label set
 - Accuracy increases and runtime decreases
- VENoM-(2,2) overall more robust than counterparts
 - Higher accuracy due to structural look-ahead ability
 - Accuracy drop slower
- Limited improvement in accuracy
 - VENoM-(1,1), VENoM-(2,1), VENoM-(3,1)
- Runtime decreases with larger group size
 - For graphs with lower average degree

Summary

1. **Key contribution:** *Units and Groups*
2. Effects of neighborhood size
 - Increase in depth \implies trade-off between time and accuracy
3. Runtime decreases with larger group size
 - Increase in breadth with depth may increase runtime efficiency
4. Diversified experiments over various parameters on both real and synthetic datasets
 - Size of graph, vertex degree, number of labels, query degree etc.



Questions?

Thank you!

Paper



Code



Homepage

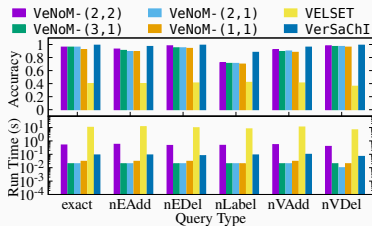


LinkedIn

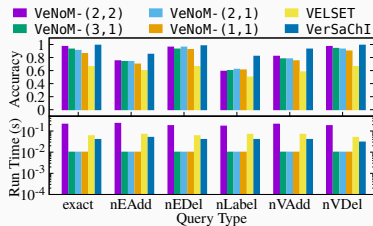


RESULTS: REAL-WORLD GRAPHS

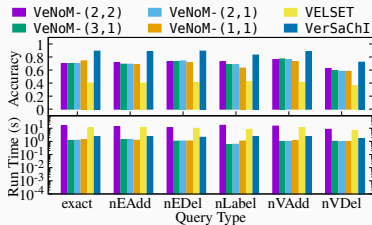
Human



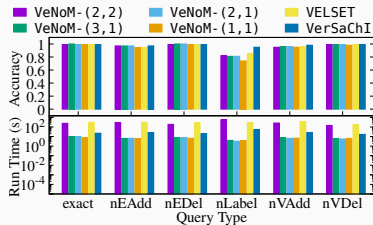
HPRD



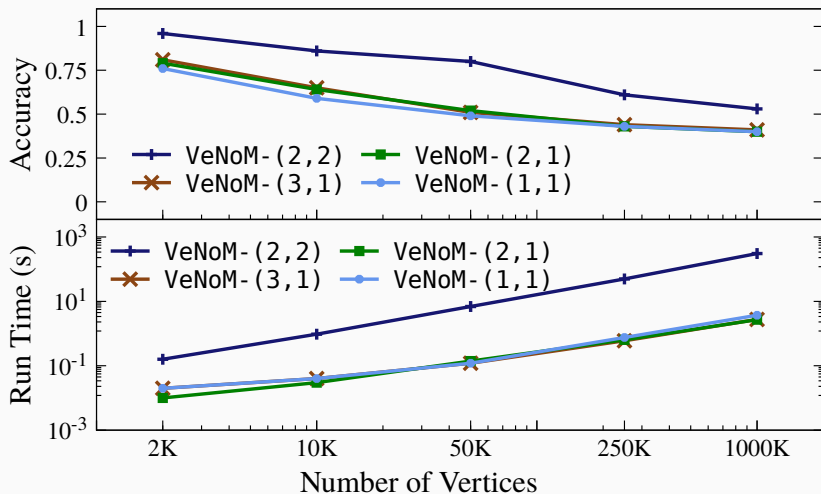
Flickr



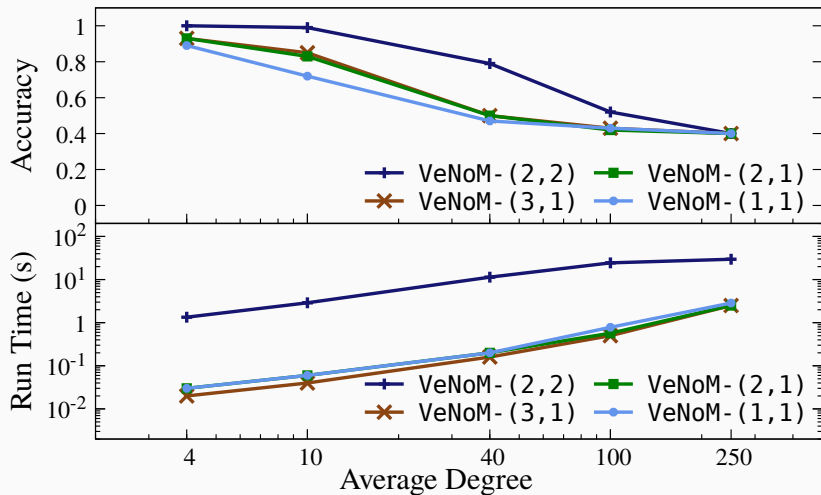
PPI



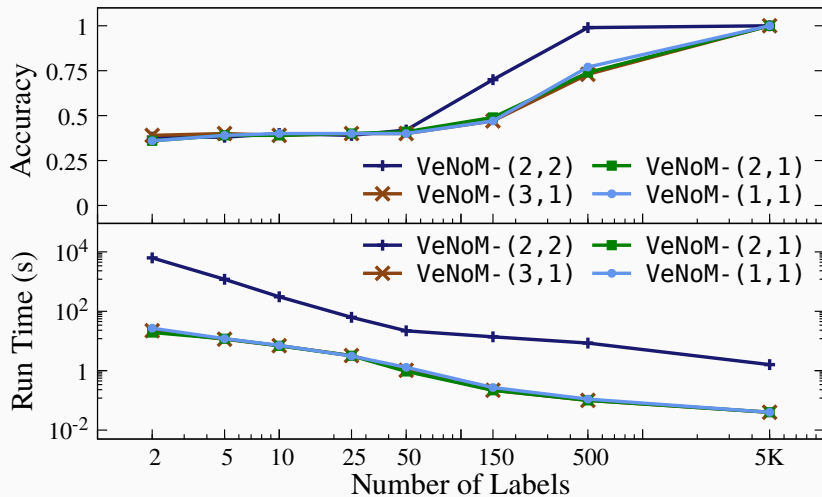
BA GRAPHS - VERTEX SCALING



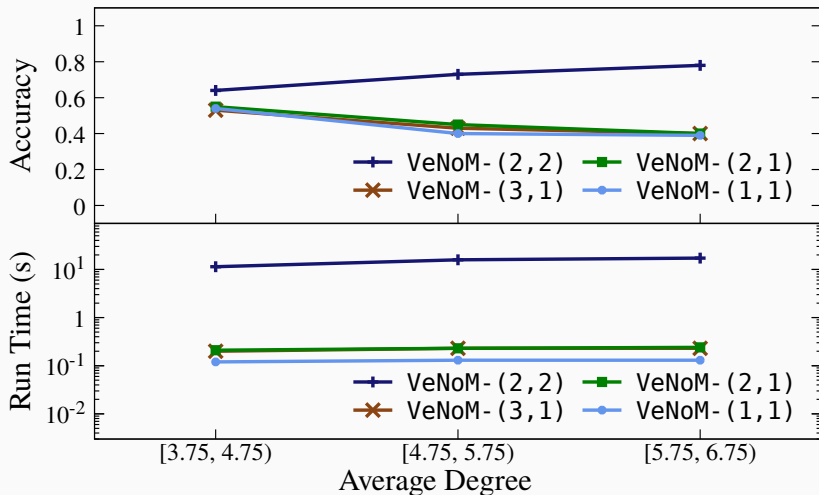
BA GRAPHS - DEGREE SCALING



BA GRAPHS - LABEL SCALING



BA GRAPHS - QUERY DEGREE



BA GRAPHS - NOISY QUERY

