

kMatrix: A Space Efficient Streaming Graph Summarization Technique

Oshan Mudannayake

University of Colombo School of Computing

Colombo, Sri Lanka

oshan.ivantha@gmail.com

Nalin Ranasinghe

University of Colombo School of Computing

Colombo, Sri Lanka

dnr@ucsc.cmb.ac.lk

Abstract—The amount of collected information on data repositories has vastly increased with the advent of the internet. It has become increasingly complex to deal with these massive data streams due to their sheer volume and the throughput of incoming data. Many of these data streams are mapped into graphs, which helps discover some of their properties. However, due to the difficulty in processing massive streaming graphs, they are summarized such that their properties can be approximately evaluated using the summaries. gSketch, TCM, and gMatrix are some of the major streaming graph summarization techniques. Our primary contribution is devising kMatrix, which is much more memory efficient than existing streaming graph summarization techniques. We achieved this by partitioning the allocated memory using a sample of the original graph stream. Through the experiments, we show that kMatrix can achieve a significantly less error for the queries using the same space as that of TCM and gMatrix.

Index Terms—graph querying, streaming graphs, summarization

I. INTRODUCTION

Massive-scale datasets are becoming increasingly common today. The growth of the number of users who are actively using digital devices connected to the internet has vastly affected this phenomenon. Also, there lies an interest in researchers to solve the problems which involve large datasets. Most of these datasets could be mapped into graphs to extract useful information, giving rise to the need for processing massive scale graphs. There are many practical scenarios where massive scale graphs are applied such as social networks, network traffic data, and road networks. Large scale dynamic natural graphs are used by many companies today. Google uses the PageRank algorithm [1], [2] to map the links between the web pages. Facebook has a massive graph with trillions of edges [3], depicting the interactions of each user on the platform.

It is much easier to work with graphs when they are static and small. However, most of the natural graphs that are being encountered in the real world are dynamic. It becomes increasingly complex to handle the graph as the velocity with which its edges get updated increases. Determining the properties of streaming graphs is a relatively strenuous task than static graphs as they are continually evolving. Thus the traditional graph algorithms cannot be run on streaming graphs due to their dynamic nature. High throughput of update queries requires any other types of queries to be run efficiently and as fast as possible in an unblocking manner. Therefore if there

is a need for processing the graph while streaming, separate streaming graph algorithms have to be devised [4].

Since real-world streaming graphs could grow very large in size, they are often stored as partitions in different machines over a network rather than in a single location. It is difficult to evaluate the properties of a graph with high volume and throughput even after the partitioning process, as the whole graph would have to be processed despite the partitioning. Graph summarization is a technique used in dealing with these massive graphs taking the limitations mentioned above into account so that important information regarding the underlying dataset can be inferred easily. In graph summarization, we reduce the complexity of a graph while retaining only some of its properties. These summaries often incur an error when queried due to the loss of information. When the same algorithm is executed on a graph summary and its original graph, the two results are expected to be approximately equal. Here, the error depends on the compression ratio and various other factors. This tradeoff in accuracy is usually worth it for real-world graphs such as social networks when considering the computational cost incurred in obtaining exact answers. Most of the time, the cost of obtaining an exact solution is so high that it is impossible to do so even if the need arises.

Being applied in a wide range of industrial and research applications, realtime property evaluation of streaming and dynamic natural graphs is a critical requirement in many scenarios. Graph summarization plays a significant role in this as it reduces the computational resources required to evaluate the properties in a rather massive scale streaming graph. It would be beneficial for many sectors if the process of summarizing streaming graphs were made efficient.

In this work, we propose an improved streaming graph summarization technique; kMatrix. It can outperform the existing state of the art summarization sketches by efficiently using the available memory to answer the queries more accurately. We also show that kMatrix is generally faster than the other sketches in handling the graph streams. Despite the number of methods that have been devised for streaming graph summarization, they still lack the accuracy to be used in most real-world scenarios [5]. Our motivation in improving the existing sketching techniques lies in increasing the efficiency while maintaining the same resource constraints of the application domains, such as real-time property evaluation of the social

networks where streaming graph summarization is critical.

The remainder of this paper is organized as follows. We explore the related work for this research in Section II. In Section III and IV, we will focus on the methodology and the implementation respectively. We will summarize all the results obtained during the experiments in Section V. Section VI will address the remaining work to be done before deploying the kMatrix in a real world application. We will conclude the paper in Section VII, highlighting the importance of this work to the graph summarization domain.

II. RELATED WORK

Summarizing a graph can have many benefits [6] apart from the speedup of graph algorithms and queries, such as, reduction of data volume and storage [7], visualization [8], [9], noise elimination [10], privacy preservation [11]. Graph summarization has a wide range of industrial and research applications as well. Some of them are clustering [12], classification [13], community detection [14], outlier detection [15], [16], pattern set mining [17] and finding sources of infection in large graphs [18]. Throughout this work, our aim lies in query optimization through graph summarization.

Streaming graph summarization is much more complex than summarizing a static graph due to the constant data flow. Since the underlying graph is updated continuously, the summarization process also has to be done in realtime. Almost any static graph summarization technique can be used with a streaming graph snapshot in a specific timestamp. However, mining information using aggregate time snapshots of data could prove to be a less than ideal solution when considering massive data streams. Thus sophisticated sparsification techniques have to be derived in order to summarize streaming graphs.

A. CountMin

CountMin [19] is a 2-dimensional data structure that is used for frequency approximation queries. It has a width of $w = \lceil e/\epsilon \rceil$ and a depth of $d = \lceil \ln(1/\delta) \rceil$. Here the e is the base of the natural logarithm while ϵ and δ are user-specified constants. The underlying idea is to hash the aggregated frequencies of the edges using multiple hash functions into predefined blocks, as indicated in Fig. 1. Any incoming edge e_t at timestamp t will get hashed into each row using its hash function h_d . A CountMin sketch will have a fixed memory allocation of $w \cdot d$ throughout its lifespan. Irrespective of the volume of the data stored in the sketch, the initial memory allocation will not change. Thus the accuracy of the queries will decrease as more and more data is inserted into the sketch. Despite the weaknesses, CountMin can be considered as a good generalized summarization sketch as many other current techniques are geared towards specific graph computation scenarios. However, the CountMin approach is not restricted to streaming graphs but other applications as well [19].

B. gSketch

gSketch [20] is an extension of CountMin data structure. But unlike the CountMin sketch, this is specifically geared

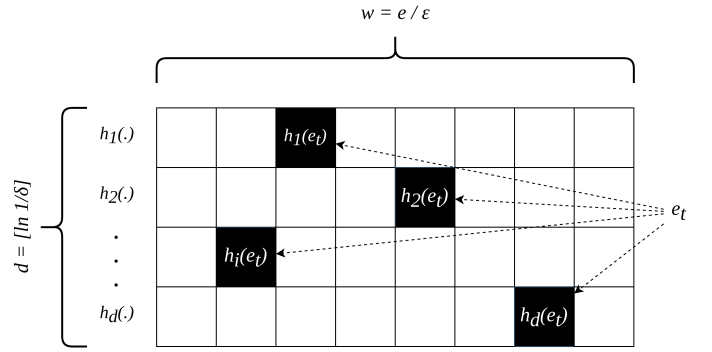


Fig. 1. CountMin sketch [20]

towards summarizing graph streams. gSketch is based on one of the below two assumptions.

- A sample of the graph stream is available.
- Samples of both the graph stream and the query workload are available.

In CountMin, one global sketch is created for the entire stream. By doing so, it fails to take advantage of any structural properties present in the graph stream. gSketch tries to avoid this by considering the underlying structure of the graph stream using a sample. It then proceeds to partition its allocated space, as indicated in Fig. 2. The sum of the widths of the partitions w_1 and w_2 is equal to the original width of the sketch, w . The goal of this partitioning step aims to maintain a sufficient frequency uniformity within each localized sketch in a way such that the combined error of the quarry estimations over the entire graph is kept at a minimum.

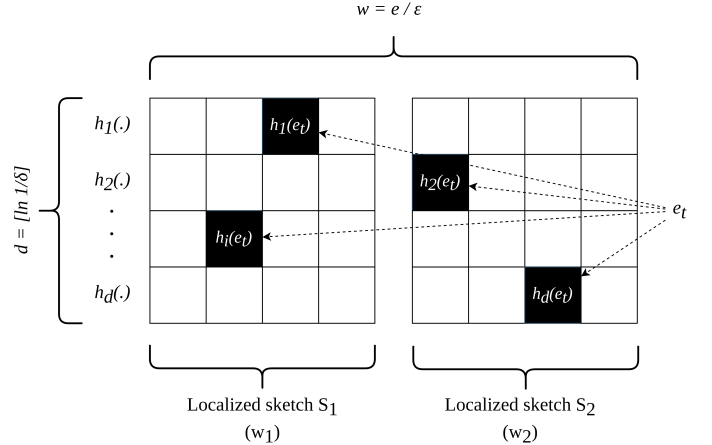


Fig. 2. gSketch sketch

C. TCM

A disadvantage posed by all the approximate frequency count sketches like CountMin or gSketch is that they do not store the locality of the nodes. Therefore CountMin and gSketch cannot be used for conditional node queries or queries involving node connectivity. If these queries were to be run,

the locality of the nodes has to be retained in the graph synopses. TCM [21] aims to solve this issue by storing the connectivity of the nodes in its data structure. TCM can summarize both node and edge information in constant time. Thus, it can answer a wide range of queries, unlike its predecessors. The structure of a TCM sketch is depicted in Fig. 3. In here w is the width of the sketch while d is the number of the hash functions. An edge (i, j) of a graph stream will get hashed into the bucket at the location $(h_r(i), h_r(j))$ in the r th layer. TCM sketch could be considered as one of the pioneering works in summarizing data streams, which is directly related to our work presented in this paper.

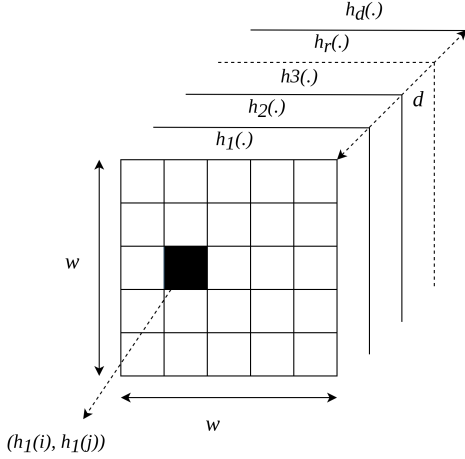


Fig. 3. TCM sketch [22]

D. gMatrix

The functionality of gMatrix [22] is very similar to the TCM sketch. However, gMatrix considers several aspects which the TCM sketch does not address.

- Reverse hashing queries through pairwise independent hash functions.
- Alternative options to extend sketch and space-saving synopses.

III. METHODOLOGY

Through this research, we propose kMatrix, which is an improvement over the traditional gMatrix algorithm. The idea behind the kMatrix is to partition the 3-dimensional frequency matrix using a sample of the original graph steam as proposed in gSketch [20]. This idea has already been discussed in the TCM [21] work to a certain degree. However, we explore this approach extensively with the gMatrix data structure, which can answer reachability queries, unlike the TCM sketch. kMatrix can also answer the reachability queries, which makes it suitable for more application scenarios than TCM. The significance of our approach is that kMatrix can answer all the queries that gMatrix is able to with much higher accuracy while occupying the same amount of space as its counterpart gMatrix sketch.

A. kMatrix

Let a stream be, $G = \langle e_1, e_2, \dots, e_m \rangle$. This can be mapped to a graph, $G = (V, E)$ where V is the set of nodes and E is a set of edges as $\{e_1, e_2, \dots, e_m\}$. We can summarize this graph using a 3-dimensional matrix sketch [22]. The straightforward choice would be to use a sketch similar to the one shown in Fig. 3. An edge, $(i, j) \in E$ will be hashed onto each layer of the sketch with has functions, h_r , $r \in \{1, \dots, d\}$. The coordinate of the cell where the edge value is preserved will be $(h_r(i), h_r(j))$. Since the kMatrix aims to use the gMatrix sketch's advantages over TCM, the hash functions should be pairwise independent of each other.

However, by constructing a global sketch for the entire graph stream, some critical information about the structural properties of the underlying graph is dismissed. It is possible to improve the performance of a sketch by retaining some of these properties. Sketch partitioning [20] is one of the techniques that allow us to improve the sketch using the properties of the graph stream. In sketch partitioning, the global sketch is partitioned using a sample of the original graph stream such that it is possible to maintain a sufficient frequency uniformity within each partition. In this work, we use the sketch partitioning process discussed in gSketch to increase the accuracy of the queries further.

Fig. 4 depicts the high-level view of kMatrix sketch after partitioning. Here, the sum of the memory occupied by all the localized sketches is equal to the memory allocated for the initial global sketch. The proceeding section will explain the partitioning algorithm in detail.

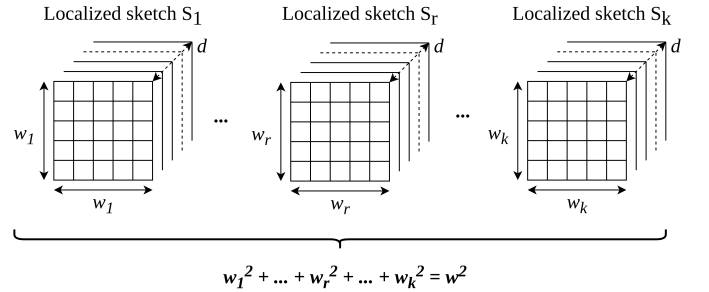


Fig. 4. kMatrix sketch

Partitioning Algorithm: Consider that the original sketch is partitioned into i sub-sketches. Let $F(S_i)$ be the sum of the edge frequencies in the i th sketch and w_i be its width. If (m, n) is an edge in the i th sketch, let $f(m, n)$ and $\bar{f}(m, n)$ be its frequency and expected frequency (1) respectively.

$$\bar{f}(m, n) = \frac{F(S_i) - f(m, n)}{w_i} \quad (1)$$

Then the expected relative error of the edge (m, n) is given by (2).

$$\bar{e}(m, n) = \frac{\bar{f}(m, n)}{f(m, n)} = \frac{F(S_i)}{f(m, n) \cdot w_i} - \frac{1}{w_i} \quad (2)$$

The overall relative error E_i of the sketch can be expressed as the sum of the expected relative errors of all the edges.

$$E_i = \sum_{(m,n) \in S_i} \bar{e}(m,n) \quad (3)$$

Let the average frequency of a vertex m be $\tilde{f}_v(m)$ and the estimated out-degree be $\tilde{d}(m)$. Then the average frequency of the vertex would be $\tilde{f}_v(m)/\tilde{d}(m)$. Therefore the total estimated frequencies of the partitioned sketch S_i can be expressed as (4).

$$\tilde{F}(S_i) = \sum_{m \in S_i; m \in V} \tilde{f}_v(m) \quad (4)$$

According to (2), (3) and (4), the overall relative error of the sketch can be simplified to (5).

$$E_i = \sum_{m \in S_i} \frac{\tilde{d}(m) \cdot \tilde{F}(S_i)}{w_i \cdot (\tilde{f}_v(m)/\tilde{d}(m))} - \sum_{m \in S_i} \frac{\tilde{d}(m)}{w_i} \quad (5)$$

$\tilde{d}(m)$ in the numerator accounts for the fact that $O(\tilde{d}(m))$ edges are coming out of the vertex m .

When a sketch of width w is partitioned into two sketches of widths w_1 and w_2 , the total error can be expressed as $E = E_1 + E_2$. Let $w_1 = w_2$. Then,

$$E = \sum_{m \in S_1} \frac{\tilde{d}(m) \cdot \tilde{F}(S_i)}{w_1 \cdot (\tilde{f}_v(m)/\tilde{d}(m))} + \sum_{m \in S_2} \frac{\tilde{d}(m) \cdot \tilde{F}(S_i)}{w_2 \cdot (\tilde{f}_v(m)/\tilde{d}(m))} - \sum_{m \in S_1 \cup S_2} \frac{\tilde{d}(m)}{w_1} \quad (6)$$

The (6) can be further simplified as,

$$E' = E \cdot w_1 + \sum_{m \in S_1 \cup S_2} \tilde{d}(m) \quad (7)$$

where the value of E' is,

$$E' = \sum_{m \in S_1} \frac{\tilde{d}(m) \cdot \tilde{F}(S_i)}{\tilde{f}_v(m)/\tilde{d}(m)} + \sum_{m \in S_2} \frac{\tilde{d}(m) \cdot \tilde{F}(S_i)}{\tilde{f}_v(m)/\tilde{d}(m)} \quad (8)$$

Thus it can be shown that the overall error in (7) can be minimized by choosing the smallest E' according to (8). Therefore the underlying idea behind the partitioning algorithm is to choose a data sample of the original stream and then repeatedly partition the available space between the vertices in the sample according to (8). After this partitioning phase, the streaming can begin and the edges that represented the vertices in the sample are put into their respective partitioned sketches.

A separate data structure has to be used in order to track the vertices belonging to different localized partition. However the extra cost of storing this information is negligible when compared with the advantages obtained with the sketch partitioning.

IV. IMPLEMENTATION

A. Experimental Setup

The implementation mainly consists of two components; the test suite and the sketching algorithms. The entire code-base has been written in Python 3.8. All the tests were run on a 12-core Ryzen 3900 machine with a base clock of 3.1GHz and 32 GB RAM. However, only one core was utilized in running the tests.

A sample of 30,000 edges has been extracted from relevant datasets for initializing kMatrix at the beginning of each experiment. This sample stream has been obtained using reservoir sampling.

B. Datasets

3 datasets were chosen to carry out the benchmarking process in this research. These were chosen to represent different application domains.

a) *unicorn-wget* [23]: unicorn-wget is a dataset created from capturing the packet information of the network activity of a simulated network. This dataset was created at Harvard University. The dataset consists of 5 parts. From them, Hour-Long Wget Benign Dataset (Base Graph) which consist of 17,778 nodes and 2,779,726 edges was chosen for the experiment. We filtered 10% of the edges using reservoir sampling for our experiments.

b) *email-EuAll* [24]: This data was extracted using email data from a large European research institution. The dataset consists of emails sent out in a period of 18 months. Each data item contains sender, receiver and the time of the origination of each email. The dataset consisted of 265,214 nodes and 420,045 edges [25].

c) *cit-HepPh* [26]: cit-HepPh citation graph is from the e-print arXiv regarding high energy physics phenomenology. It has 34,546 papers (nodes) and 421,578 citations (edges). We used the full dataset in our experiments.

C. Evaluation Metrics

1) *Average Relative Error (ARE)*: The relative error $er(Q)$ of a query Q is defined as (9) where $\tilde{f}'(Q)$ and $f(Q)$ is the estimated frequency and the true frequency of the query respectively.

$$er(Q) = \frac{\tilde{f}'(Q) - f(Q)}{f(Q)} = \frac{\tilde{f}'(Q)}{f(Q)} - 1 \quad (9)$$

Given a set of m queries, $\{Q_1, \dots, Q_m\}$, the average relative error is defined by taking the average of the relative error of all queries Q_i for $i \in [1, m]$.

$$e(Q) = \frac{\sum_{i=1}^k er(Q_i)}{m} \quad (10)$$

2) *Number of Effective Queries (NEQ)*: A query is said to be effective if the error, $\tilde{f}'(Q) - f(Q) \leq G_0$, where G_0 is a predefined value. The number of effective queries is defined as,

$$g(Q) = |\{q \mid (\tilde{f}'(q) - f(q)) \leq G_0, q \in Q\}| \quad (11)$$

V. RESULTS

This section will describe all the experiments conducted to measure the effectiveness of kMatrix against existing streaming graph sketching techniques.

We have considered CountMin, gSketch, TCM, gMatrix and kMatrix sketches in our experiments. These sketches can be categorized into two groups depending on the type of queries they are able to answer.

- 1) *Type I* - The sketches which support only the edge frequency queries, i.e. CountMin and gSketch.
- 2) *Type II* - The sketches which support many graph queries in general, i.e. TCM, gMatrix and kMatrix

Since *Type I* sketches cannot answer anything other than edge frequency queries, we have only included *Type II* sketches in our comparisons against kMatrix.

A. Build-time

Here we investigated the time to add the entire dataset to the sketch. The sketches were allocated a constant memory size of 1 MB, and the number of hash functions was set to $d = 7$. The edges were streamed at the maximum throughput of each sketch. Therefore this experiment gives an idea about the average insertion rate of edges for each sketch. A minor drawback of kMatrix is that it takes some time for its initialization stage. However, this initialization time becomes negligible compared to the advantage that kMatrix receives over time due to its faster streaming rate. In both Fig. 5a and Fig. 5c, it has managed to outperform other sketching techniques by a significant margin. In Fig. 5b, kMatrix has shown comparable performance to gMatrix. With the increase of the data contained within the sketch, the number of hash collisions in TCM and gMatrix has grown over time, increasing the computational cost of inserting a new edge. kMatrix has maintained a relatively lower build time as a result of its lower number of hash collisions due to the sketch partitioning before inserting the edges.

B. Edge queries

This experiment investigates how accurately the kMatrix can answer the edge queries after the summarization process. For this, we let our datastream get summarized into the sketch and then queried the frequency of different edges chosen at random. The experiment was repeated for each sketch for the sizes, 200 KB, 300 KB, 400 KB and 512 KB while keeping the number of hash functions at $d = 7$. We have used average relative error and the number of effective queries as the evaluation matrices for this experiment.

1) *Average Relative Error*: kMatrix showed significantly low ARE than all the other sketches for the three datasets we chose. The reason is that kMatrix can maintain frequency uniformity within each partition, making kMatrix relatively more immune to hash collisions than TCM and gMatrix. It is clear from the experimental evidence shown in Fig. 6 that kMatrix vastly outperforms the other state of the art sketching techniques. The superiority of our solution is more apparent when the allocated memory is low.

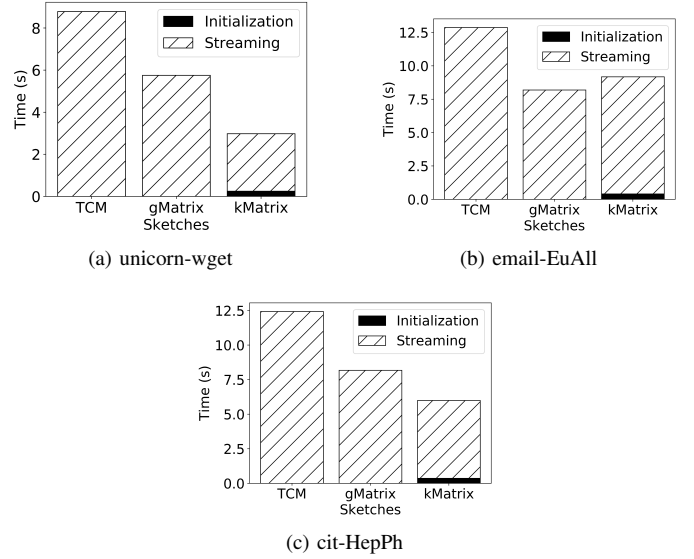


Fig. 5. Build-time

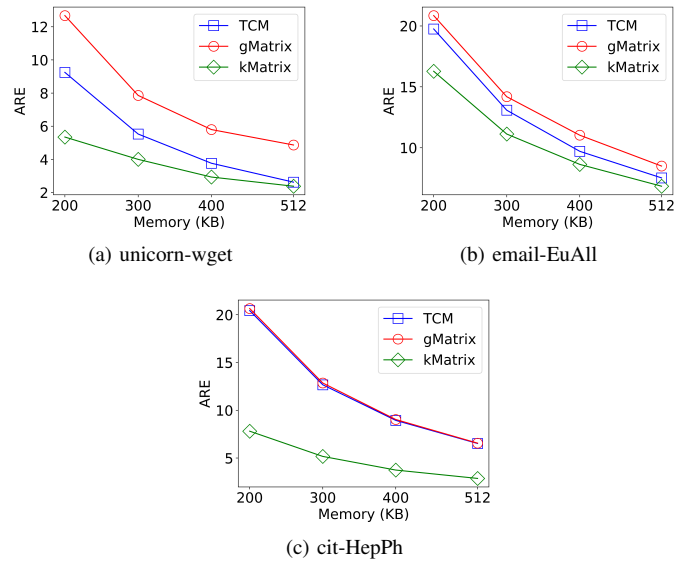


Fig. 6. Average relative error

2) *Number of Effective Queries*: The number of effective queries for each sketch was calculated by querying the sketches against 10,000 edges chosen through reservoir sampling from the original dataset. kMatrix has surpassed the accuracy of both TCM and gMatrix for all the scenarios that we have tested. The results for cit-HepPh in Fig. 7c shows that kMatrix has been able to effectively answer a significantly larger number of queries where the other sketches failed due to hash collisions. This is due to the sketch partitioning process where kMatrix try to minimize the hash collisions in contrast to TCM and gMatrix.

VI. FUTURE WORK

There are multiple aspects such as sliding windows and data partitioning across machines, that should be considered

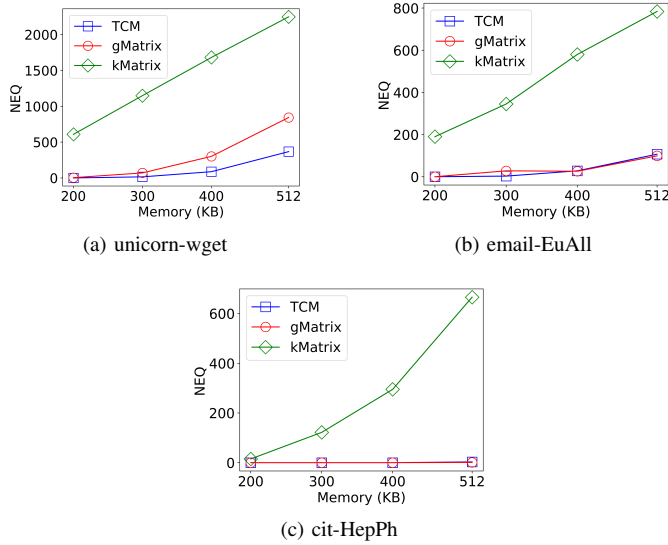


Fig. 7. Number of effective queries

before the kMatrix sketch be used in a practical application. In addition, we have to further test the performance of kMatrix concerning other criteria such as heavy node/edge queries. The test suit's functionalities can be extended and improved upon as a benchmarking tool for testing the graph summarization sketches.

VII. CONCLUSION

kMatrix is a new streaming graph summarization technique proposed through this research. It can answer queries with a much lower average relative error with the same amount of memory compared to the existing state-of-the-art sketching techniques, TCM and gMatrix. Decreasing the error of the queries answered by the summarized sketches will have a significant impact on all the application scenarios that require the use of graph summarization. The dramatic decrease in the error of the queries shown by the kMatrix makes it a better replacement for any current application domain where TCM or gMatrix has been utilized. We have benchmarked kMatrix using three datasets in different application domains to test out its performance. We believe that the experimental results show the superiority of the proposed solution in comparison to the existing steaming graph summarization techniques.

REFERENCES

- [1] S. Brin and L. Page, "The anatomy of a large-scale hypertextual Web search engine," *Computer Networks and ISDN Systems*, vol. 30, pp. 107–117, Apr. 1998.
- [2] L. Page, "The PageRank Citation Ranking: Bringing Order to the Web,"
- [3] A. Ching, S. Edunov, M. Kabiljo, D. Logothetis, and S. Muthukrishnan, "One trillion edges: graph processing at Facebook-scale," *Proceedings of the VLDB Endowment*, vol. 8, pp. 1804–1815, Aug. 2015.
- [4] A. McGregor, "Graph stream algorithms: a survey," *ACM SIGMOD Record*, vol. 43, pp. 9–20, May 2014.
- [5] M. S. Kumarage, Y. Horawalavithana, and D. Ranasinghe, "An efficient query platform for streaming and dynamic natural graphs," in *2017 IEEE International Conference on Industrial and Information Systems (ICIIS)*, (Peradeniya), pp. 1–6, IEEE, Dec. 2017.
- [6] Y. Liu, T. Safavi, A. Dighe, and D. Koutra, "Graph Summarization Methods and Applications: A Survey," *ACM Computing Surveys*, vol. 51, pp. 1–34, June 2018.
- [7] H. Seo, K. Park, Y. Han, H. Kim, M. Umair, K. U. Khan, and Y.-K. Lee, "An effective graph summarization and compression technique for a large-scaled graph," *The Journal of Supercomputing*, Jan. 2018.
- [8] C. Dunne and B. Shneiderman, "Motif simplification: improving network visualization readability with fan, connector, and clique glyphs," in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems - CHI '13*, (Paris, France), p. 3247, ACM Press, 2013.
- [9] L. Jin and D. Koutra, "ECO : Comparative Visualization of Time-Evolving Network Summaries," p. 8.
- [10] N. Zhang, Y. Tian, and J. M. Patel, "Discovery-driven graph summarization," in *2010 IEEE 26th International Conference on Data Engineering (ICDE 2010)*, (Long Beach, CA, USA), pp. 880–891, IEEE, 2010.
- [11] M. Shooran, A. Thomo, and J. H. Weber-Jahnke, "Zero-knowledge private graph summarization," in *2013 IEEE International Conference on Big Data*, (Silicon Valley, CA, USA), pp. 597–605, IEEE, Oct. 2013.
- [12] R. Cilibrasi and P. M. B. Vitanyi, "Clustering by Compression," *IEEE TRANSACTIONS ON INFORMATION THEORY*, vol. 51, no. 4, p. 21, 2005.
- [13] M. van Leeuwen, J. Vreeken, and A. Siebes, "Compression Picks Item Sets That Matter," in *Knowledge Discovery in Databases: PKDD 2006* (D. Hutchison, T. Kanade, J. Kittler, J. M. Kleinberg, F. Mattern, J. C. Mitchell, M. Naor, O. Nierstrasz, C. Pandu Rangan, B. Steffen, M. Sudan, D. Terzopoulos, D. Tygar, M. Y. Vardi, G. Weikum, J. Fürnkranz, T. Scheffer, and M. Spiliopoulou, eds.), vol. 4213, pp. 585–592, Berlin, Heidelberg: Springer Berlin Heidelberg, 2006.
- [14] D. Chakrabarti, D. S. Modha, S. Papadimitriou, and C. Faloutsos, "Fully Automatic Cross-associations," p. 12.
- [15] K. Smets and J. Vreeken, "The Odd One Out: Identifying and Characterising Anomalies," in *Proceedings of the 2011 SIAM International Conference on Data Mining*, pp. 804–815, Society for Industrial and Applied Mathematics, Apr. 2011.
- [16] L. Akoglu, D. H. Chau, U. Kang, D. Koutra, and C. Faloutsos, "OPAvion: mining and visualization in large graphs," in *Proceedings of the 2012 international conference on Management of Data - SIGMOD '12*, (Scottsdale, Arizona, USA), p. 717, ACM Press, 2012.
- [17] M. Mampaey, N. Tatti, and J. Vreeken, "Tell me what i need to know: succinctly summarizing data with itemsets," in *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD '11*, (San Diego, California, USA), p. 573, ACM Press, 2011.
- [18] B. A. Prakash, J. Vreeken, and C. Faloutsos, "Spotting Culprits in Epidemics: How Many and Which Ones?," in *2012 IEEE 12th International Conference on Data Mining*, (Brussels, Belgium), pp. 11–20, IEEE, Dec. 2012.
- [19] G. Cormode and S. Muthukrishnan, "An Improved Data Stream Summary: The Count-Min Sketch and its Applications," p. 18, 2003.
- [20] P. Zhao, C. C. Aggarwal, and M. Wang, "gSketch: on query estimation in graph streams," *Proceedings of the VLDB Endowment*, vol. 5, pp. 193–204, Nov. 2011.
- [21] N. Tang, Q. Chen, and P. Mitra, "Graph Stream Summarization: From Big Bang to Big Crunch," in *Proceedings of the 2016 International Conference on Management of Data - SIGMOD '16*, (San Francisco, California, USA), pp. 1481–1496, ACM Press, 2016.
- [22] A. Khan and C. Aggarwal, "Query-friendly compression of graph streams," in *2016 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, (San Francisco, CA, USA), pp. 130–137, IEEE, Aug. 2016.
- [23] X. Han, "Hour-Long Wget Benign Dataset (Base Graph)," 2018.
- [24] J. Leskovec, J. Kleinberg, and C. Faloutsos, "Graph evolution: Densefication and shrinking diameters," *ACM Transactions on Knowledge Discovery from Data*, vol. 1, pp. 2–es, Mar. 2007.
- [25] "SNAP: Network datasets: EU email network."
- [26] J. Gehrke, P. Ginsparg, and J. Kleinberg, "Overview of the 2003 KDD Cup," *ACM SIGKDD Explorations Newsletter*, vol. 5, p. 149, Dec. 2003.