

# AAE: An Active Auto-Estimator for Improving Graph Storage

Yu Yan<sup>a</sup>, Man Yang<sup>a</sup>, Hongzhi Wang<sup>a</sup>, Yuzhuo Wang<sup>a</sup>

<sup>a</sup>*Harbin Institute of Technology, Harbin, China*

---

## Abstract

Nowadays, graph becomes an increasingly popular model in many real applications. The efficiency of graph storage is crucial for these applications. Generally speaking, the tune tasks of graph storage rely on the database administrators (DBAs) to find the best graph storage. However, DBAs make the tune decisions by mainly relying on their experiences and intuition. Due to the limitations of DBAs's experiences, the tunes may have an uncertain performance and conduct worse efficiency. In this paper, we observe that an estimator of graph workload has the potential to guarantee the performance of tune operations. Unfortunately, because of the complex characteristics of graph evaluation task, there exists no mature estimator for graph workload. We formulate the evaluation task of graph workload as a classification task and carefully design the feature engineering process, including graph data features, graph workload features and graph storage features. Considering the complex features of graph and the huge time consumption in graph workload execution, it is difficult for the graph workload estimator to obtain enough training set. So, we propose an active auto-estimator (AAE) for the graph workload evaluation by combining the active learning and deep learning. AAE could achieve good evaluation efficiency with limited training set. We test the time efficiency and evaluation accuracy of AAE with two open source graph data, LDBC and Freebase. Experimental results show that our estimator could efficiently complete the graph workload evaluation in milliseconds.

---

\*Corresponding author

---

**keywords:**graph storage, workload evaluation,active learning

## 1. Introduction

Nowadays, graph becomes an increasingly popular model in many real applications [6]. For example, in FaceBook, the relationships among users are stored as property graph. In machine cognitive learning, the knowledge of experience is stored as knowledge graph [2]. Efficient graph storage is the key step for improving the economic income in these applications by accelerating the workload execution.

Generally speaking, the graph storage tuning tasks mainly rely the database administrators (DBAs) in real applications. However, DBAs complete the tuning by utilizing their experiences and intuition, which may conduct performance degradation. For example, a graph has two properties, age and gender. There exist many queries about gender property. DBAs may establish an index in gender according to the frequency. But the gender property only has two kinds of values in which the index is useless. It is unreliable to tune the graph storage only relying on manual work. Importantly, in real applications, any performance degradation will influence the economic income. So, we need an estimator to help DBAs guarantee the efficiency of the storage tunes.

Unfortunately, existing workload estimators [8] [9] [10] [11] [12] which are mainly designed for the sample data models like relational model, key value model and document model is unsuitable for the graph model. For example, the learned model [12] proposed by R.J et al could only process some sample filters, like ‘range’, ‘in’ and etc. while do not support the traversal filters of graph data. Specifically, due to the complexity of graph workload evaluation, there exists no mature work for this task. In this paper, we pay attention to construct an estimator for evaluating the graph workload to guarantee the tune performance in real applications.

Firstly, we conclude three aspect reasons why evaluating graph workload is

Table 1: The Cost of Labeling

database& graph data	Create Operation	Read Operation	Update Operation	Delete Operation	Traversal Operation	Workload Runtime(hours)
New4J+LDBC	38%	15%	2%	13%	32%	11.57
	16%	31%	3%	3%	47%	31.25
	17%	25%	17%	10%	31%	17.36
	9%	17%	9%	10%	31%	23.15
	20%	12%	10%	12%	46%	34.72
New4J+Freebase	20%	18%	13%	11%	28%	18.52
	14%	12%	28%	11%	35%	11.53
	10%	49%	6%	11%	24%	28.33
	12%	23%	4%	33%	28%	57.87
	19%	19%	3%	11%	57%	173.61

difficult as follows:

- Complex Structures** In order to save various entities and relationships, graph data has very complex structures. For example, there exist more than 1000 million edges in LinkedGeoData. And DBpedia has more than 3000 million triples. Big graph has more complex structural features and statistical features [16]. Evaluating the graph workload in such big graph is very difficult.
- Diverse Query Operations** To make full use of graph data information, the query patterns in graph are usually diverse including 'create nodes', 'read nodes', 'update nodes', 'delete nodes' and 'traverse nodes'. And every kind has diverse basic operations [7]. Only the 'read nodes' contains four basic operations seen in Table 2. Evaluate graph workload under diverse query operations is very difficult.

- **Dynamic Workloads** In real applications, the workload is changing with the user’s requests. The dynamic workloads also bring big challenges to the evaluation method. On the one hand, changing workloads require fast evaluation to guarantee the efficient self-tuning. On the other hand, too many factors in workloads, which have great impact on graph workload execution, need to carefully evaluate. It is necessary for the estimator to quickly complete the evaluation task with high accuracy.

Thus, it is necessary to design a special evaluation model for graph workloads. Inspired by the recent machine learning methods [49] [36] in database committee, we firstly observe how to construct an automatic graph workload estimator to carry the above difficulties. We formulate the graph workload evaluation task into a classification task and explore how to utilize deep learning classifiers to complete this evaluation task. The input consists of the graph workload, the graph data, old-graph storage and new-graph storage. The output presents if the new revision of storage structure is better than the old storage structure.

However, deep learning based models need a large number of labeled data. The result of our preliminary investigation in Table 1 shows that it is incredible to obtain enough labeled data in limited time. We can find that the first item in Table 1 would take 11.57 hours in New4J and LDBC dataset. In order to reduce the cost of data labeling, we combine the active learning structure [17] to our estimator and propose an active auto-estimator (AAE) for graph workload to decrease the demand of training sets. Our experimental results prove that only using a bit labeled data can still make good performance. For example, AAE reach 72% accuracy only with 42% training set. And our estimator could complete the evaluation task with milliseconds.

The contributions of this paper are summarized as follows:

- We firstly formulate the definition of graph workload evaluation problem which could help to improve the graph storage and increase the economic efficiency of applications.

- We carefully design the features of graph data, graph workload and graph storage, containing the structural features and statistical features of graph data, various graph query pattern features, basic storage schema and indexes features.
- We design a graph workload estimator, AAE which combines the deep learning and active learning. Our model could process the complex features, including data features, workload features and graph storage features. And AAE could quickly complete the graph workload evaluation with the limited training set.
- We examine our active estimator with two widely used dataset, LDBC and Freebase in microbenchmark [43]. The experimental result shows that AAE could largely reduce the demand of training dataset. And our estimator could complete evaluation with milliseconds.

The structure of this paper is as follows. Section 2 presents some related works about graph storage and graph query acceleration. Section 3 introduces the overview of AAE. Section 4 discuss the feature engineering of three key components of AAE, including graph data, graph workload and graph storage. And Section 5 carefully introduce the active auto-estimator. Finally, we clarify the experimental result in section 6, containing experimental dataset and result analysis.

## 2. related work

Because of the lacks of the graph workload evaluation methods, we only introduce some related works about graph storage and graph workload acceleration.

**Graph Storage** The storage of graph data is divided into native graph database storage [23] and non-native graph database storage [24, 25]. The native graph database implements point-to-point data structure and indexing; the non-native graph database uses other database systems to store graph data

and implement query interfaces. Native graph database systems such as Neo4j, InfiniteGraph, Sparksee [21]. Non-native graph database systems such as OrientDB, ArangoDB and various relational databases. Compared with native storage, non-native storage systems [26, 27] have more mature support in terms of concurrency [22], locks, security, and query optimization.

**Graph workload Acceleration** Recently, researchers mainly optimal the graph query by two key aspects. One is to construct indexes. Compared to other data models, graph data has more complex data structure and diverse query operations. For quickly scan the complex graph, researches index the subgraph [14], path [15] and etc. Different from the sample data model, using the general indexes (like B+ tree, Hash, LSM tree) do not quickly find some complex patterns [30]. Expect the index there also exists some other methods to accelerate graph query. People propose some cache and prefetching methods [18] [31] [32] [33] [34] based on similarity counting. And recent works [33] prefetch graph query by computing the graph edit distance between queries. Totally, the graph database management urgently need the evaluation modules to accelerate the graph query.

### 3. Overview

Traditional estimators could be divided into two main class, single query based [12] and workload based [37]. Considering the demand of real-world settings like FaceBook which could receive millions of queries at one time, we take the total graph workload as the evaluation target. There are two significant influential factors to graph workload performance, graph data and graph storage. In this paper, we formulate these key components of our estimator as follows:

**Graph Data:**  $\mathcal{G} = (N, R, P, T)$ , where  $N$  is the collection of nodes,  $R$  is the collection of edges,  $P$  is the collection of properties, and  $T$  is the property edges.

**Graph Workload:**  $\mathcal{W} = \{q_1, q_2, \dots, q_n\}$ , where  $q_i$  is a graph query.

**Graph Storage:**  $\mathcal{S} = \langle S_S, S_I \rangle$ , where  $S_S$  is the basic storage schema and

the  $S_I$  is the indexes in  $S_S$ .

The total cost of graph workload is denoted as follows.  $w_i$  is the weight of  $q_i$  and  $c_i$  is the cost of  $q_i$ .

$$Cost(G, W, S) = \sum_{i=1}^n w_i * c_i \quad (1)$$

Comparing the workload execution cost of two graph storage plans could help DBAs to stably tune the graph storage. In order to improve graph storage, we formulate the graph workload evaluation task as a classification task.

Definition(the graph workload evaluation task): Give the graph data  $G$  and workload  $W$ , We define the cost a graph workload evaluation task as follows:

$$AAE(G, W, S_{old}, S_{new}) = \begin{cases} 1, & Cost(G, W, S_{old}) > \\ & Cost(G, W, S_{new}), \\ 0, & Cost(G, W, S_{old}) \leq \\ & Cost(G, W, S_{new}). \end{cases} \quad (2)$$

The input of AAE consists of the graph workload, the graph data, the old graph storage and the new graph storage. The output presents if the new graph storage is better than the old one in the current graph data and graph workload. Our AAE could help DBAs to improve the graph storage by quickly evaluating graph workloads.

As figure 1 shows, after receiving the graph storage tune from DBAs, our system implements the evaluation to judge if the tune is positive. If so, this tune from DBAs could be implemented to database. Our estimator could help DBAs to prevent the performance degradation.

#### 4. Feature Engineering

Feature engineering is the key step of constructing the graph workload estimator. Comparing with the evaluation for other sample data models (like relation, key-value), the graph evaluation has more complex characteristics, such as the complex graph structure, diverse query operations and dynamic

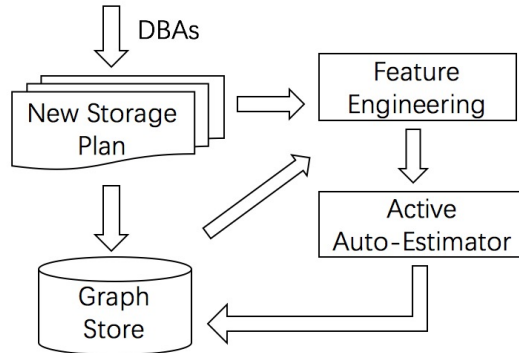


Figure 1: The Workflow of AEE

workloads. In this section, we introduce the feature engineering of constructing graph workload estimator, including the graph data, the graph workload and the graph storage.

**Dataset Feature Extraction:** Generally Speaking, the characteristics of graph data mainly include statistical characteristics and structural characteristics. In this paper, we pay more attention to these two kinds of features for the representation of graph data. Considering the scale and complexity of the graph data, we conclude that some vital statistical features, including data size, number of nodes, number of edges. And some key structural features include the number of nodes or edges’ type, the number of properties and the property values of node or edge. Our estimator which use both statistical features and structural features has more complete representation for graph dataset.

**Workload Feature Extraction:** Graph has various basic query operations. We carefully divide these operations into five types as listed in Table 2, containing create, read, update, delete and traverse. We can find that different from other sample data models (like relation, key-value), graph has more plentiful operations. We conclude these operations by totally analyzing the graph query and existing works [18, 7]. And the rate of these operations in workload has vital influence for graph workload evaluation. For example, in native graph storage, traversal graph queries usually have higher efficiency than read queries.



While in relational graph storage, read queries usually have higher efficiency than traversal queries [6]. By extracting the characteristics of the workload, we can identify the main types of operations, and improve the accuracy of our estimator. Through these workload information, we can have a more detailed judgement of the graph storage. Finally, we utilize the rate of each basic query operation and the frequency of properties as the representation of the graph workloads.

**Storage Solution Features:** In order to query graph data quickly, there are some different kinds of graph store engines, including New4J, Titan, ArangoDB and etc. Also, there are more other features in corresponding engines, such as index support, query optimizer and other feature of distribution design. In this paper, we choose the basic stores and indexes as the total represented features for graph storage solution. Because of the basic stores and indexes information are not numeral features, we utilize the one-hot method [39] to encode them. And one-hot method could largely retain the original characteristics. Specifically, we encode every kind of engines as a one-hot vector. For index information, we consider the indexes on graph node properties and edge properties. We also utilize the one-hot vector to identify if a property has the index.

## 5. Active Auto-Estimator

It is difficult to obtaining a large amount of training set in graph workload evaluation task with limited time. We propose an active-based solution to reduce the time consumption of training model. In this section, we introduce the active auto-estimator module.

### 5.1. Active Structure

Generally speaking, the classifiers based deep learning need a large number of labeled dataset for model training. However, as our preliminary results in Table 1, graph workload execution requires a lot of time and calculation resources, it is incredible for obtaining enough labeled dataset by implement graph queries

Table 2: graph query operations

Type	Query Function	Description
Create	addVertex(G, n)	Create a new node n to graph G
	addEdge(G, e)	Create an edge e from n1 to n2 in Graph G
	addProperty(G.N.n,p)	Create a property to nodes or edges
Read	getCount(G.N)	Count the total number of nodes or edges
	getProperty(G.E)	Get all edge or nodes properties in graph G
	findProperty(G.N, p)	search nodes or edges in G.N according to the property p
	find(G.N, n)	Search nodes or edges in G.N according to the given condition
Update	setProperty(G.N,n,p)	Update the properties of nodes or edges
Delete	removeVertex(G, n)	Delete the node according to the identifier n
	removeEdge(G, e)	Delete the edge according to the identifier e
	removeProperty(G.N,n,p)	Delete the property p of nodes or edges
Traverse	in(G.N, n)	Query all nodes adjacent to the node n with an incoming edge
	out(G.N, n)	Query all nodes adjacent to the node n with an outgoing edge
	all(G.N, n)	Query all nodes adjacent to the node n
	TFilter(count(in(G.N,n)) $\geq$ k)	Filter all the nodes which have at least k-incoming-degree
	allinPath(BFS(G,n))	Search all the nodes reached via breadth-First traversal from node n
	allinPath(BFS(G,n,l))	Same as above, but with edges' labels limitation
	shortPath(G,n1,n2)	Get the shortest path from node n1 to node n2 without weight
shortPath(G,n1,n2,l)	Same as above, but with edges' labels limitation	

---

**Algorithm 1:** Active Learning

---

**input:**  $U$  is the unlabeled dataset  
 $L$  is the labeled dataset  
 $\Theta$  is the parameters of deep classifiers  
 $T$  is the threshold of the estimation

**output:**  $\Theta$  the trained parameters

```
1 def  $AAE(U, L, \Theta, T)$ :  
2 while  $U \neq \emptyset$  do  
3   Sample  $S$  from  $U$   
4    $\Theta = \mathcal{P}(S, \Theta)$   
5    $U = \{u \mid f(\Theta, u) < T, u \in U\}$   
6 end  
7 return  $(\Theta)$ 
```

---

in graph database with certain storage and graph. In order to reduce the cost of obtaining labeled data, we design an active based model to evaluate graph storage. The active learning [46] means training models only with partial dataset. Firstly, we sample some data from the unlabeled data. And next, we implement data training with these sample data. Then, we use the trained model to evaluate unlabeled data. If the accuracy reach the threshold, we stop this interaction. And if the accuracy is not enough, we would sample the poor data point which does not perform well under the trained model to retrain the classifier. Repeat the above process, we could obtain an evaluation model with smaller training set finally.

Our active learning method is shown in Algorithm 1. The input of this algorithm consists of four parts, including the unlabeled dataset  $U$ , the labeled dataset  $L$ , the parameters of deep classifier  $\Theta$  and the threshold of estimation  $T$ . We repeat the sampling and retraining process until that the unlabeled dataset is empty. In every round, we first randomly sample some data from the unlabeled dataset. And then use these data to train deep classifier. Next,

we use the deep classifiers to evaluate the point in  $U$ . If the evaluation result reaches the threshold, we put this point into  $L$ . Finally, we return the parameter of the deep classifier. The deep classifier is the key component of our AAE. We consider to construct our classification based estimator by utilizing popular neural networks.

## 5.2. Deep classifier

Nowadays, deep learning express great learning ability in many fields, such as image classification, NLP and etc. The popular deep model mainly contains two basic kinds, convolutional neural networks(CNN) and recurrent neural networks(RNN). And many traditional networks consist of them, such as VGG [38], RCN [40] and etc. In this paper, considering the different characteristics of different graph dataset, we attempt to employ two kinds of networks to learn the complex graph features introduced in section 4.

### 5.2.1. CNN based classifier

The features of our evaluation task is every complex. For example, there exists significant interaction between the statistical features of graph data and the storage vectors. So, we firstly attempt to use CNN to construct evaluation model. Compared with NNs(Neural Networks), CNN has more powerful capabilities in learning sectional features. The different convolution kernel filters could learn more mature features. And convolutional filters could extract the sectional features of the graph data, graph workload, graph storage solution. Consequently, the multi-layer convolution kernel will learn model detailed knowledge.

Considering that the number of layers has great influence on learning ability. In this paper, we design two CNN models, a simpler one and a deeper one, for researching the effect of network structure in graph workload evaluation. As shown in the Table 3, the first simpler CNN model(SCNN) contains four convolutional layers. And the deeper CNN model(DCNN) contains six convolutional layers. For the SCNN, we first use two convolutional layers with 16 filters to ex-

Table 3: The deep classifier based on CNN

Layer Name	Parameter Value	
Input		ALL
Conv1	filters=16,kernel_size=3	ALL
Conv2	filters=16,kernel_size=3,activation=tanh	ALL
MaxPooling1	pool_size=3	ALL
Conv3	filters=16,kernel_size=3,activation=tanh	ALL
Conv4	filters=16,kernel_size=3,activation=tanh	ALL
MaxPooling2	pool_size=3	ALL
Conv5	filters=16,kernel_size=3,activation=tanh	DCNN
Conv6	filters=16,kernel_size=3,activation=tanh	DCNN
MaxPooling3	pool_size=3	DCNN
Flatten		ALL
Dense	1,activation='sigmoid'	ALL
Output		ALL

tract features in original input vector. And we employ the tanh as the activation function to avoid the gradient disappearance. After two convolutional layers, a maximization pool MaxPooling1D is added to retain the main features. Next, we add two convolutional layers and one pooling layer to enhance the learning ability of our model. Finally, we add the flatten layer to connect the convolutional layer and the fully connected layer. And for the DCNN, we add extra two convolutional layers and pooling layer to improve the capacity of evaluation model.

### 5.2.2. RNN based classifier

For complex characteristics of graph evaluation task, the previous partial input would influence the later output in current partial input. Therefore, we

also consider to utilize RNN to process our evaluation task. Different from traditional NNs, RNNs introduce directional loops, which can deal with the correlation between inputs. The specific manifestation is that the network will memorize the previous information and apply it to the calculation of the current output. However, for traditional RNNs, there are some problems of gradient disappearance and gradient explosion. In order to solve these problems, K et al. proposed GRU[47], a variant of LSTM network[48], is also a type of RNN. And the internal structure of GRU is similar to that of LSTM. So, we utilize the GRU to construct the deep graph estimator.

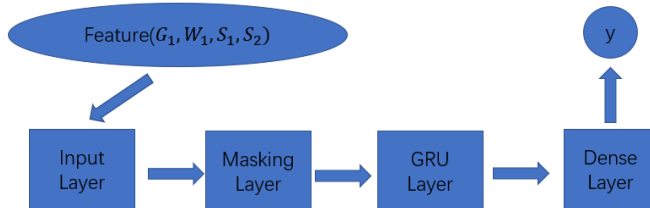


Figure 2: The Deep classifier based on GRU

As shown in Figure 2, The GRU classifier consists of four layers. The input layers is responsible for receiving input vector which is support after feature engineering. Because the different graph data, graph workload and graph storage solution have the different features, constructed feature vectors have variable length. In order to train these vectors in one batch, we implement data preprocessing (fill -1 to vectors) to unify the length of vectors. And correspondingly, we add a mask layer to prevent the influence of the filling operation. Next, we add an GRU layer to learn the evaluation task. Finally, we could obtain the evaluation result after a dense layer.

## 6. Experimental Results

In this section, we express our examination results of our model. First, we introduce the experimental settings in section 6.1, containing dataset and

execution environment. Then, we discuss our experimental results from three aspects, active performance in section 6.2, classifier comparison in section 6.3 and evaluation time cost in section 6.4. The source codes of our experiments are available <sup>1</sup>.

## 6.1. Experimental Settings

### 6.1.1. Dataset

Table 4: Features of datasets

	<b>Freebase- small</b>	<b>Freebase- middle</b>	<b>LDBC</b>
number of nodes	480577	4264156	184328
number of edges	314753	3147537	767894
number of node types	1	1	8
number of edge types	1814	2912	15
number of property types	3	3	62

In this paper, we utilize two open source graph datasets, LDBC and Freebase in micro-benchmark [43] to conduct our experiments. Micro-benchmark consists of seven synthetic and real dataset. The selected two datasets in our study are the Freebase [44](a real dataset) and the LDBC dataset [45](a synthetic dataset). In Freebase, a subgraph (Frb-O), which only has nodes related to organization, business, government, finance, geography, and military are considered, is created. Then, by randomly selecting 0.1 percent and 1 percent edges of this subgraph, two datasets, Freebase-small and Freebase-middle are derived. For the LDBC dataset, the synthetic dataset (LDBC) is generated by using the data generator provided by the Linked Data Benchmark Council. The graph it generates simulates the characteristics of a real social network with a power

---

<sup>1</sup><https://github.com/yangmanST/Graph-Storage/tree/master>

law structure and the characteristics of the real world. Finally, we use different dataset to conduct our experiments, as shown in Table 4.

### 6.1.2. Environment

In this paper, we use two open source graph databases to complete our experiments, as shown in Table 5. And our experiments are conducted on the Ubuntu 18.04 with 8G memory and 80G disk. The deep classifier is trained with CPU and 8G memory.

Table 5: Features of datasets

Name	Version	Storage	Query Language
New4J	1.9	Native Graph	Gremlin
Titan	0.4.5	Columnar store	Gremlin

### 6.2. Active Performance

In this section, we introduce the evaluation accuracy with the different rate of training set in different datasets. Table 6 reports the classification accuracy of different active stage (58%, 49%, 41%) in Freebase-small, Freebase-middle and LDBC. We train our classifiers with learning rate = 0.01 and traditional BP algorithm. After carefully tuning, our active learned estimator achieves 99.85% prediction accuracy in LDBC labeled dataset ( $L$  in algorithm 1) and reach 79.34% prediction accuracy in LDBC unlabeled dataset ( $U$  in algorithm 1) with only 41% training dataset. We can see that our active structure could reduce the demand of training set. The estimator still performs well in the low rate of training set. That is, our active method could support efficient evaluation with lower time consumption.



Table 6: The Classification Accuracy in Different TrainSet

Dataset	Train Set Rate	Test Set	SCNN	DCNN	GRU
Freebase-small	58%	<i>L</i>	94.19%	93.02%	91.86%
		<i>U</i>	68.75%	68.75%	56.25%
Freebase-middle	49%	<i>L</i>	93.25%	82.02%	85.39%
		<i>U</i>	72.73%	63.64%	61.16%
LDBC	41%	<i>L</i>	99.85%	99.85%	90.62%
		<i>U</i>	78.51%	79.34%	68.60%

### 6.3. CNN & GRU

In this paper, we propose three deep classifiers to learning the complex graph workload evaluation task. In this section, we compare the performance of these deep classifiers in two aspects. One is the classification accuracy, and the other is the cross validation in different classifiers.

As shown in Figure 3, we test the prediction accuracy in our three classifiers in unlabeled dataset (*U*). We can find that SCNN is significantly better than other two classifiers (DCNN and GRU). For SCNN, its accuracy is 10.1% higher than DCNN and 12.6% higher than GRU on average. Correspondingly, Figure 4 shows the accuracy in three classifiers with full dataset (*L + U*). The accuracy of SCNN is 7.8% higher than DCNN on average, and 18.3% higher than GRU.

Figure 5 presents the result of cross-validation of the three classifiers on the LDBC. The average accuracy and its standard deviation of SCNN is 89.8% and 4.2%, DCNN is 86.1% and 4.3%, GRU is 89.7% and 4.1%. It's shown that the performance of SCNN is more stable than other two classifiers.

Comprehensively, compared to DCNN and GRU, SCNN which is expert at extracting sample features and mining the correlation of different features is the

best classifier in our graph workload evaluation task.

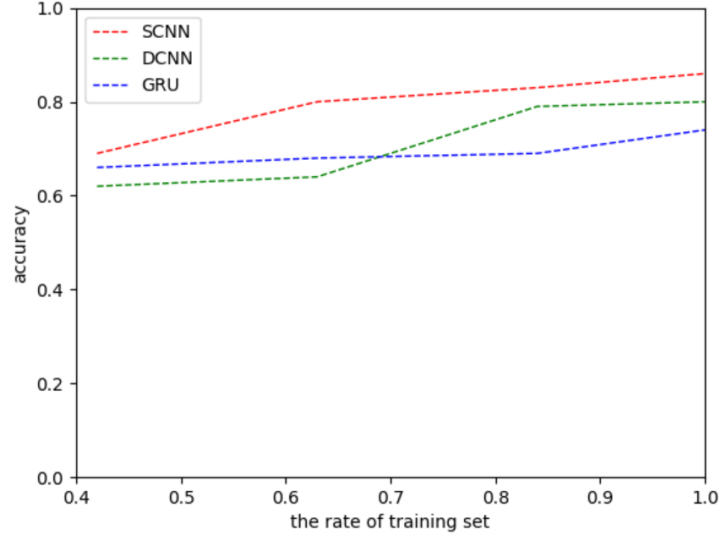


Figure 3: The Accuracy of Unlabeled Dataset (LDBC) on Different classifiers

#### 6.4. Evaluation Time Cost

In order to process the changing graph workload, the evaluation model should be able to quickly process the evaluation requirements. In this section, we discuss the average evaluation time of our classification-based estimator. As shown in Figure 6, in neo4j, we use 50 data to test the average estimation time in different dataset and classifiers. Freebase-small dataset takes separately 0.017s, 0.017s, 0.016s for graph workload evaluation. Freebase-medium dataset takes separately 0.028s, 0.023s, 0.024s. LDBC dataset takes separately 0.049s, 0.047s, 0.046s. Figure 7 shows the average time consumption in Titan. Freebase-small dataset takes separately 0.016s, 0.012s, 0.014s. Freebase-medium dataset takes separately 0.025s, 0.018s, 0.026s. LDBC dataset takes separately 0.048s, 0.037s, 0.052s. We can see that the prediction time is mainly related to the structure of classifiers. Because the runtime of deep learning based model has directly relationship with the number of neurons in classifiers. Our estimator could also quickly predict in any other dataset after some limited training.

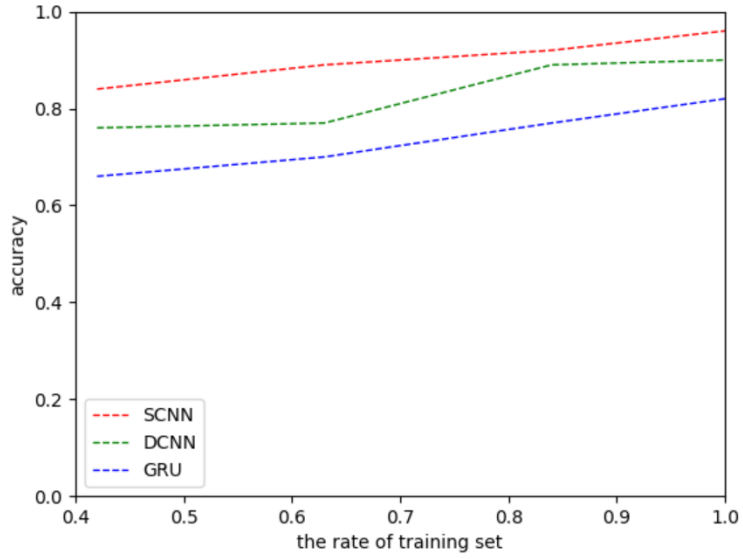


Figure 4: The Accuracy of Full Dataset(LDBC) on Different classifier

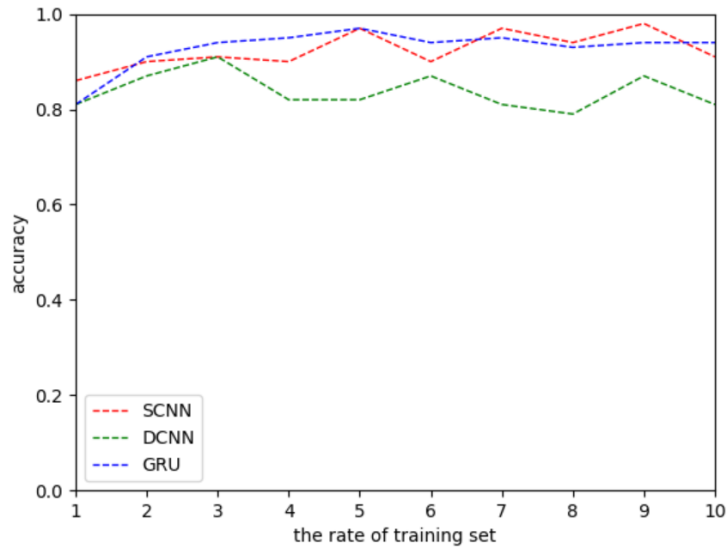


Figure 5: Cross-validation Results of LDBC on Different classifier

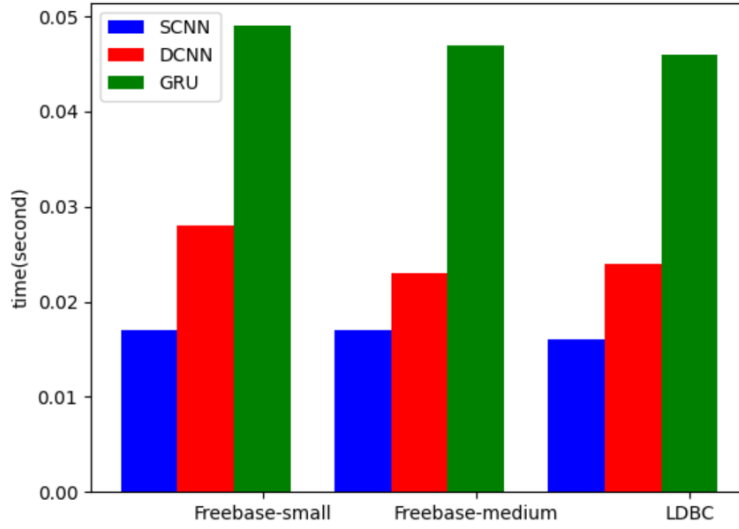


Figure 6: Average Time Cost in Different Dataset(neo4j)

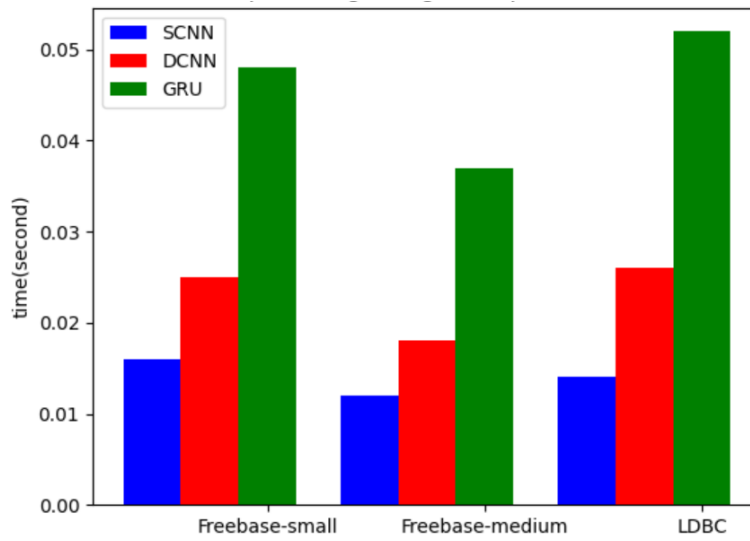


Figure 7: Average Time Cost in Different Dataset(Titan)

## 7. Conclusion

In this paper, we design an active auto-estimator for improving the graph storage solution. We carefully design the feature engineering, including graph data features, graph workload features and graph storage solution. And we clarify the time efficiency and evaluation accuracy in two open source datasets. In the future, we expect to make more improvements in the following two aspects. On the one hand, we want to figure out if the more fine-grained features (such as the graph itself) would performs well than the processed statistical features (such as the number of nodes). After all, it would loss some information after some processing in original dataset. On the other hand, we would consider to research if the new researches about auto-ML would inspire the graph workload evaluation for obtaining a perfect deep classifier.

## References

- [1] Wu W , Yun C , Zhu S , et al. Predicting query execution time: Are optimizer cost models really unusable?[C]// Data Engineering (ICDE), 2013 IEEE 29th International Conference on. IEEE, 2013.
- [2] Wang Y , Ruhe G , Gavrilova M L , et al. A formal knowledge representation system for the cognitive learning engine[C]// Proceedings of the 10th IEEE International Conference on Cognitive Informatics and Cognitive Computing, ICCI\*CC 2011, 18-20 August 2011, Banff, Alberta, Canada. IEEE, 2011.
- [3] Y. Chi, H. J. Moon. iCBS: Incremental costbased scheduling under piecewise linear slas. PVLDB, 4(9):563–574, 2011.
- [4] Kossmann J , Halfpap S , Jankrift M , et al. Magic mirror in my hand, which is the best in the land? An Experimental Evaluation of Index Selection Algorithms[J]. Proceedings of the VLDB Endowment, 2020.

- [5] Das S , Li F , Narasayya V R , et al. Automated Demand-driven Resource Scaling in Relational Database-as-a-Service[C]// the 2016 International Conference. ACM, 2016.
- [6] Z. Qi, H. Wang and H. Zhang, "A Dual-Store Structure for Knowledge Graphs," in IEEE Transactions on Knowledge and Data Engineering, doi: 10.1109/TKDE.2021.3093200.
- [7] M. A. Rodriguez and P. Neubauer, "The graph traversal pattern," Graph Data Management Techniques & Applications, 2012.
- [8] W.-C. Hou and G. Ozsoyoglu. Statistical estimators for aggregate relational algebra queries. ACM Trans. Database Syst., 16:600–654, 1991
- [9] R. J. Lipton, J. F. Naughton, and D. A. Schneider. Practical selectivity estimation through adaptive sampling. In SIGMOD, 1990.
- [10] Y. E. Ioannidis. The history of histograms (abridged). In VLDB, pages 19–30, 2003.
- [11] L. Getoor, B. Taskar, and D. Koller. Selectivity estimation using probabilistic models. In SIGMOD, 2001.
- [12] Ji Sun and Guoliang Li. 2019. An end-to-end learning-based cost estimator. Proc. VLDB Endow. 13, 3 (November 2019), 307–319.
- [13] K. Tzoumas, A. Deshpande, and C. S. Jensen. Lightweight graphical models for selectivity estimation without independence assumptions. PVLDB, 4(11):852–863, 2011.
- [14] D Sofie, Tom M , F Jan, et al. The Index-Based Subgraph Matching Algorithm (ISMA): Fast Subgraph Enumeration in Large Networks Using Optimized Search Trees[J]. Plos One, 2013, 8(4):e61183.
- [15] Kuldeep BR, Kuma PS. Efficient approximate SPARQL querying of Web of linked data. In: Bobillo F, ed. Proc. of the 6th Int'l Workshop on Uncertainty Reasoning for the Semantic Web (URSW 2010).

- [16] Meng W , Fu W , Hao S , et al. Learning on Big Graph: Label Inference and Regularization with Anchor Hierarchy[J]. IEEE Transactions on Knowledge and Data Engineering, 2017, PP(5):1-1.
- [17] Prince M . Does Active Learning Work? A Review of the Research[J]. Journal of Engineering Education, 2004, 93(3).
- [18] Papailiou N , Tsoumakos D , Karras P , et al. Graph-Aware, Workload-Adaptive SPARQL Query Caching[C]// Acm Sigmod International Conference on Management of Data. ACM, 2015.
- [19] Schmid M . On efficiently storing huge property graphs in relational database management systems[C]// iiWAS2019: The 21st International Conference on Information Integration and Web-based Applications&Services. 2019.
- [20] Bornea M A , Dolby J , Kementsietsidis A , et al. Building an efficient RDF store over a relational database[C]// Acm Sigmod International Conference on Management of Data. ACM, 2013.
- [21] N. Mart ´mez-Bazan, V. Munt ´es-Mulero, S. G ´omez-Villamor, J. Nin, M.-A. S ´anchez-Mart ´mez, and J.-L. Larriba-Pey. DEX: High- Performance Exploration on Large Graphs for Information Retrieval. In Proceedings of the 16th Conference on Information and Knowledge Management (CIKM), pages 573–582. ACM, 2007
- [22] Zhao K F,Yu J X.All-in-one:graph processing in RDBMSs revisited[C]//Proceedings of the 2017 ACM SIGMOD International Conference on Management of Data.New York: ACM,2017:1165-1180
- [23] L. Zou, M. T. Ozsul, L. Chen, X. Shen, R. Huang, and D. Zhao, “gstore: a graph-based sparql query engine” The VLDB journal, vol. 23, no. 4, pp. 565–590, 2014.
- [24] W. Sun, A. Fokoue, K. Srinivas, A. Kementsietsidis, G. Hu, andG. Xie, “Sqlgraph: An efficient relational-based property graph store,” in Proceed-

ings of the 2015 ACM SIGMOD International Conference on Management of Data, 2015, pp. 1887–1901.

- [25] Z. Pan and J. Hefflin, “Dldb: Extending relational databases to support semantic web queries,” in Proceedings of the 1st International Workshop on Practical and Scalable Semantic Systems, 2004, pp. 109–113
- [26] ZHANG Xiao, SUN Yiming, WU Xufeng. Research on query-aware relation-graph database adaptive storage technology. *Computer Engineering and Applications*, 2020, 56(17): 100-108.
- [27] Zhixin Qi, Hongzhi Wang, Member, IEEE, and Haoran Zhang A Dual-Store Structure for Knowledge Graphs *Transactions on Knowledge and Data Engineering* 2020
- [28] Chad Vicknair, Michael Macias, Zhendong Zhao, Xiaofei Nan, Yixin Chen, and Dawn Wilkins. A comparison of a graph database and a relational database: a data provenance perspective. In Proceedings of the 48th Annual Southeast Regional Conference, 2010, pages 42:1–42:6.
- [29] Mccoll, Rob & Ediger, David & Poovey, Jason & Campbell, Dan & Bader, David. (2013). A Brief Study of Open Source Graph Databases, 2013.
- [30] M. A. Rodriguez and P. Neubauer, The graph traversal pattern, *Graph Data Management Techniques & Applications*, 2012
- [31] Hogan A, Mellotte M, Powell G, Stampouli D. Towards Fuzzy Query-Relaxation for RDF. *Lecture Notes in Computer Science*. 2012:687-702.
- [32] Elbassuoni S, Ramanath M, Weikum G. Query Relaxation for Entity-Relationship Search. *The Semantic Web: Research and Applications*. 2011:62-76.
- [33] Lorey J, Naumann F. Detecting SPARQL Query Templates for Data Prefetching. *The Semantic Web: Semantics and Big Data*. 2013:124-139.



- [34] Lorey J, Naumann F. Caching and Prefetching Strategies for SPARQL Queries. *Advanced Information Systems Engineering*. 2013:46-65.
- [35] Lan H , Bao Z , Peng Y . An Index Advisor Using Deep Reinforcement Learning[C]// *CIKM '20: The 29th ACM International Conference on Information and Knowledge Management*. ACM, 2020.
- [36] Ji Sun, Guoliang Li. An End-to-End Learning-based Cost Estimator. *PVLDB*, 13(3): 307-319, 2019.
- [37] Kang J , Gaurav, Tan S Y , et al. Efficient Deep Learning Pipelines for Accurate Cost Estimations Over Large Scale Query Workload. *PVLDB*, 2021.
- [38] Simonyan K , Zisserman A . Very Deep Convolutional Networks for Large-Scale Image Recognition[J]. *Computer Science*, 2014.
- [39] Sun, Ji and Li, Guoliang. An end-to-end learning-based cost estimator. *VLDB*, 2021.
- [40] Siwei Lai, Liheng Xu, Kang Liu, and Jun Zhao. 2015. Recurrent convolutional neural networks for text classification. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence (AAAI'15)*. AAAI Press, 2267–2273.
- [41] Agrawal, Smita & Patel, Atul. (2016). A Study on Graph Storage Database of NOSQL. *International Journal on Soft Computing, Artificial Intelligence and Applications*. 5. 33-39. 10.5121/ijscai.2016.5104.
- [42] ZHANG Xiao, SUN Yiming, WU Xufeng. Research on query-aware relation-graph database adaptive storage technology. *Computer Engineering and Applications*, 2020, 56(17):100-108
- [43] Matteo Lissandrini, Martin Brugnara, and Yannis Velegrakis. Beyond macrobenchmarks: microbenchmark-based graph database evaluation. *Proc. VLDB Endow.* 12, 4 (December 2018), 2018.390–403.

- [44] Google.Freebase data dumps. <https://developers.google.com/freebase/data>, 2015
- [45] O. Erling, A. Averbuch, J. Larriba-Pey, H. Chafi, A. Gubichev, A. Prat, M.-D. Pham, and P. Boncz. The ldbc social network benchmark: Interactive workload. In SIGMOD, pages 619–630, 2015.
- [46] Lin Ma, Bailu Ding, Sudipto Das, and Adith Swaminathan . 2020. Active Learning for ML Enhanced Database Systems. In Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data (SIGMOD’20), June 14–19, 2020.
- [47] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014.
- [48] S. Hochreiter and J. Schmidhuber, “Long Short-Term Memory,” in *Neural Computation*, vol. 9, no. 8, pp. 1735-1780, 15 Nov. 1997, doi: 10.1162/neco.1997.9.8.1735.
- [49] Ding, Bailu, et al. ”Ai meets ai: Leveraging query executions to improve index recommendations.” Proceedings of the 2019 International Conference on Management of Data. 2019.