



Xian, Shaanxi, China  
SIGMOD/PODS 2021

# Learned Cardinality Estimation for Similarity Queries

Ji Sun<sup>+</sup>, Guoliang Li<sup>+</sup>, Nan Tang<sup>\*</sup>

<sup>+</sup>Tsinghua University, <sup>\*</sup>QCRI

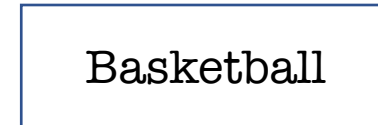
# Problem Statement

- **Cardinality for Similarity Search.** Number of objects in  $D$  whose distances to a query  $q$  are not greater than a distance threshold  $\tau$ .
- **Cardinality for Similarity Join.** Total number of pairs  $(q, p)$  whose distance between  $q \in Q$  and  $p \in D$  is not greater than  $\tau$ .

# Problem Statement



Card=20



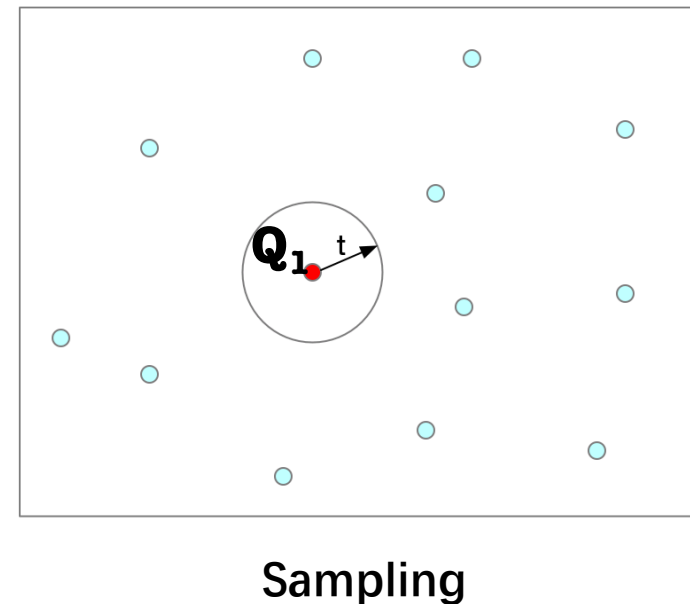
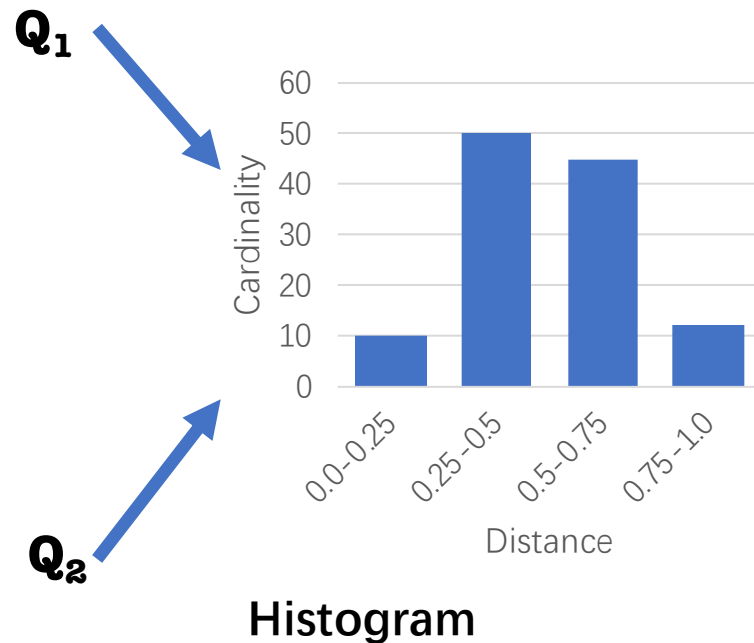
Card=9



Volleyball	Football	Tennis
Soccer	Softball	Lacrosse
Baseball	Women_Basketball	basketball

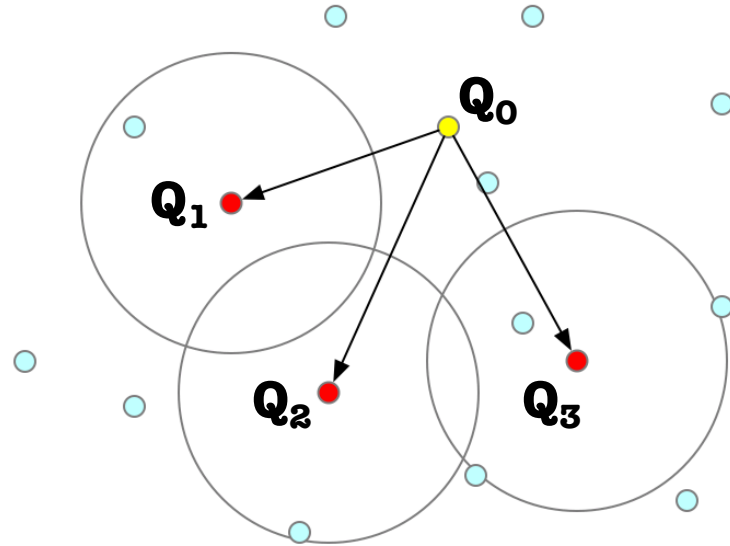
# Related Work

- Cardinality Estimator for Exact Queries
  - **Histogram**: Relative distance is defined on the given query.
  - **Sampling**: 0-tuple problem for high dimensionality.
  - **Data Model**: Hard to fit the sparse continuous data.

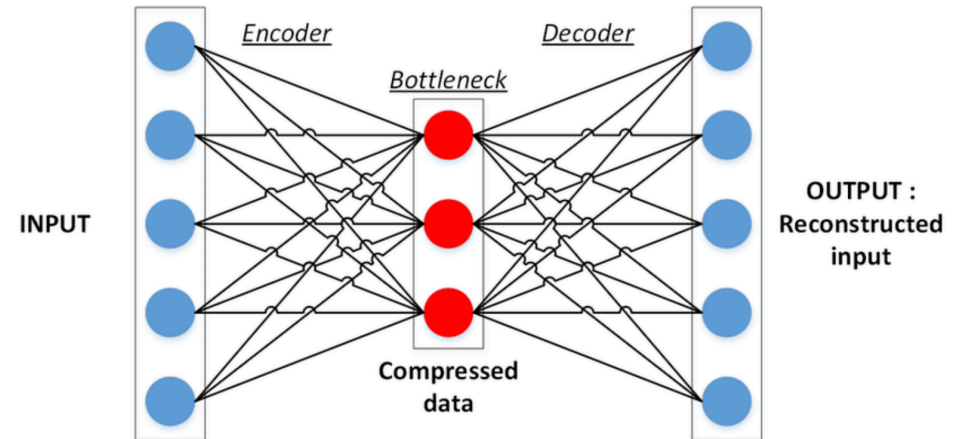


# Related Work

- Cardinality Estimator for Similarity Queries
  - **KDE**: 0-tuple problem for high dimensionality.
  - **Linear Mixture Model**: Less powerful for high dimensional data.
  - **VAE**: Low dimensional embedding is not a distance-aware representation.

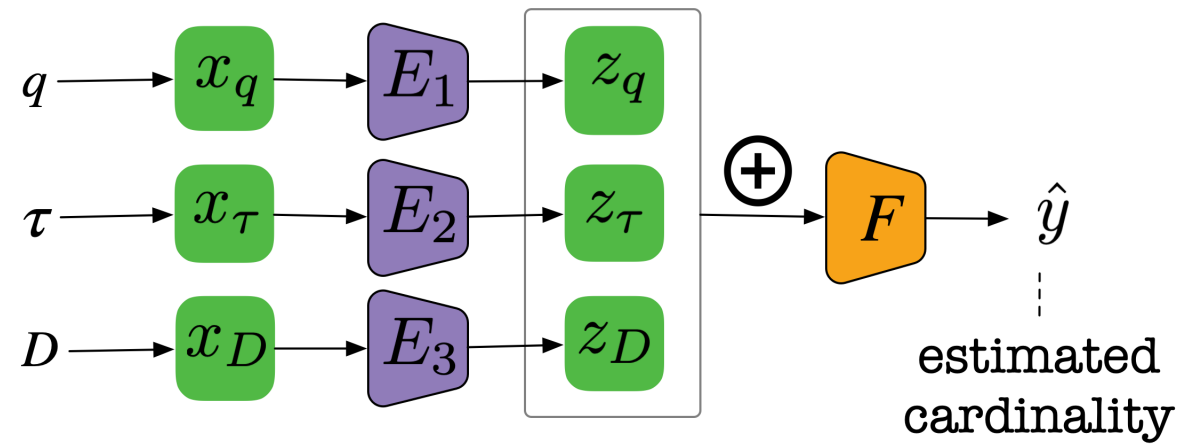


Mixture Model



Variational AutoEncoder(VAE)

# Basic Model



$\mathbf{q}$ : query vector

$\boldsymbol{\tau}$ : distance threshold

$\mathbf{D}$ : data sample

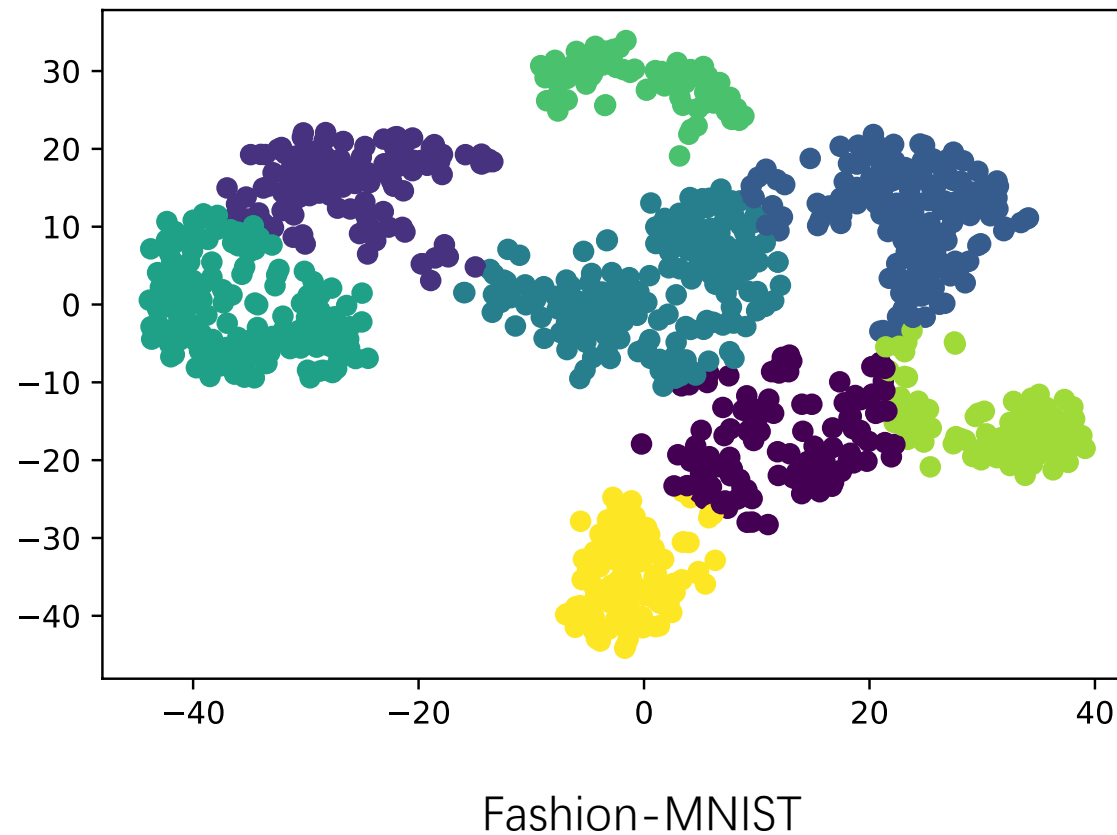
$\mathbf{E1}$ ,  $\mathbf{E2}$ ,  $\mathbf{E3}$ ,  $\mathbf{F}$ : Neural Networks

# Observations & Opportunities

- Vectors far from the query can be ignored.
- The distance of two vectors is related to sum of distances on vector segments.
  - $\text{Hamming}(\text{Sigmod}, \text{Sigkdd}) = \text{Hamming}(\text{Sig}, \text{Sig}) + \text{Hamming}(\text{mod}, \text{kdd})$
  - ...

# Observations & Opportunities

- Clustering





# Observations & Opportunities

- Distance Decomposition

**L<sub>m</sub>-distance**

$$\begin{aligned} \text{dis}_{L_m}(\mathbf{u}, \mathbf{v}) &= \sqrt[m]{\sum_{j=1}^d (\mathbf{u}[j] - \mathbf{v}[j])^m} \\ &= \sqrt[m]{\sum_{i=1}^n \sum_{j=|\mathbf{u}^{(i)}|-(i-1)}^{|\mathbf{u}^{(i)}|} (\mathbf{u}[j] - \mathbf{v}[j])^m} = \sqrt[m]{\sum_{i=1}^n (\text{dis}_{L_m}(\mathbf{u}^{(i)}, \mathbf{v}^{(i)}))^m} \end{aligned}$$

**Cosine distance**

$$\begin{aligned} \text{dis}_{\text{cos}}(\mathbf{u}, \mathbf{v}) &= 1 - \frac{\mathbf{u} \cdot \mathbf{v}}{|\mathbf{u}| \cdot |\mathbf{v}|} = \frac{|\mathbf{u}| \cdot |\mathbf{v}| - \mathbf{u} \cdot \mathbf{v}}{|\mathbf{u}| \cdot |\mathbf{v}|} \\ &= \frac{\mathbf{u}^2 + \mathbf{v}^2 - 2\mathbf{u} \cdot \mathbf{v}}{2|\mathbf{u}| \cdot |\mathbf{v}|} = \frac{\text{dis}_{L_2}(\mathbf{u}, \mathbf{v})}{2} \end{aligned}$$

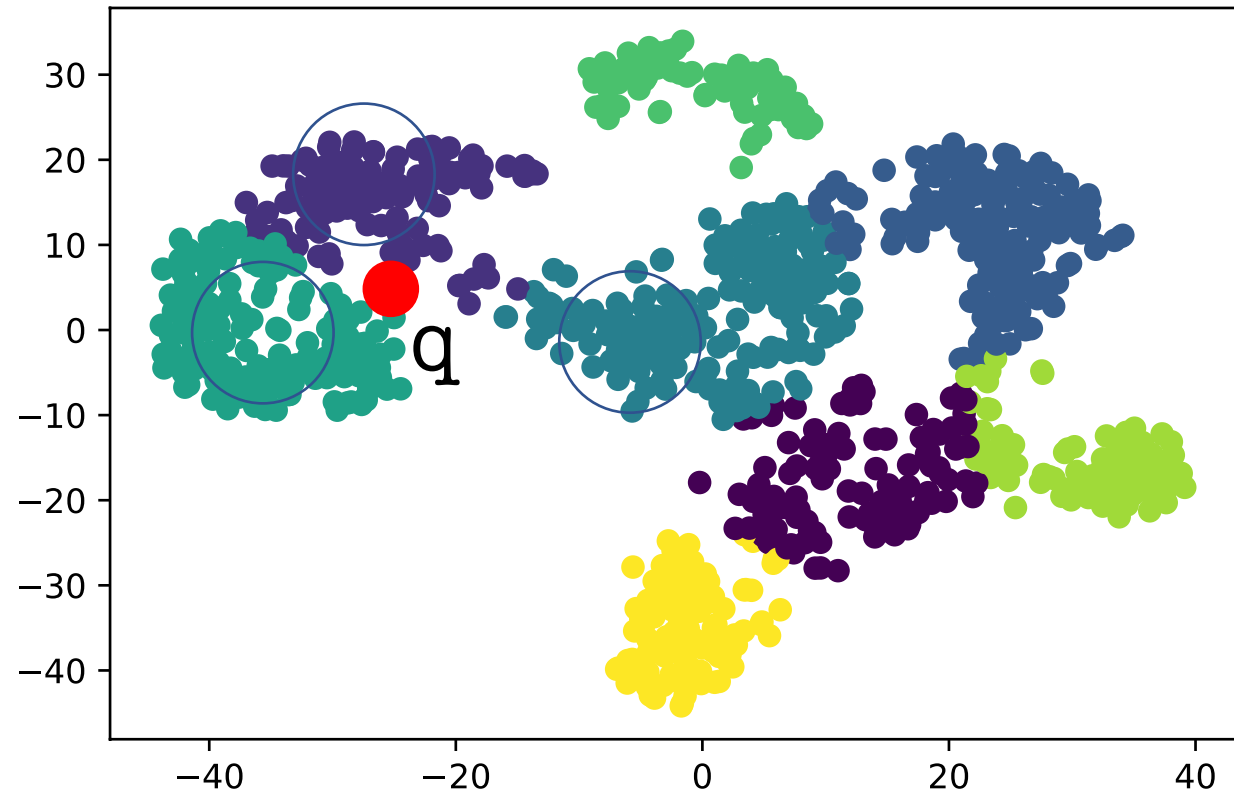
**Angular distance**

$$\text{dis}_{\text{angular}}(\mathbf{u}, \mathbf{v}) = \frac{\arccos \text{dis}_{\text{cos}}(\mathbf{u}, \mathbf{v})}{\pi}$$

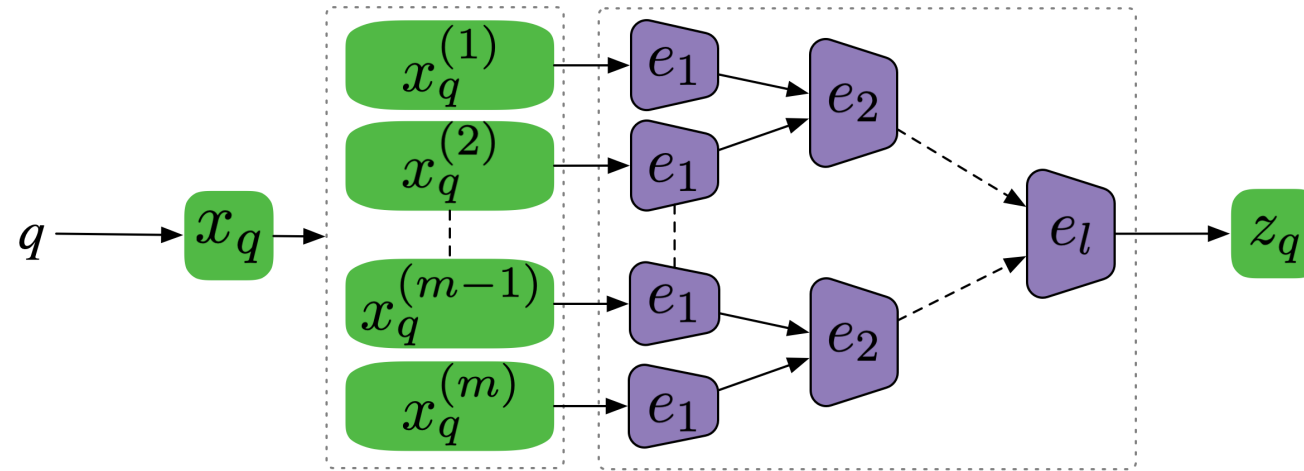
**Hamming distance**

$$\begin{aligned} \text{dis}_{\text{ham}}(\mathbf{u}, \mathbf{v}) &= \sum_{j=1}^d \text{equal}(\mathbf{u}[j], \mathbf{v}[j]) \\ &= \sum_{i=1}^n \sum_{j=|\mathbf{u}^{(i)}|-(i-1)}^{|\mathbf{u}^{(i)}|} \text{equal}(\mathbf{u}[j], \mathbf{v}[j]) = \sum_{i=1}^n \text{dis}_{\text{ham}}(\mathbf{u}^{(i)}, \mathbf{v}^{(i)}) \end{aligned}$$

# Data Segmentation



# Query Segmentation



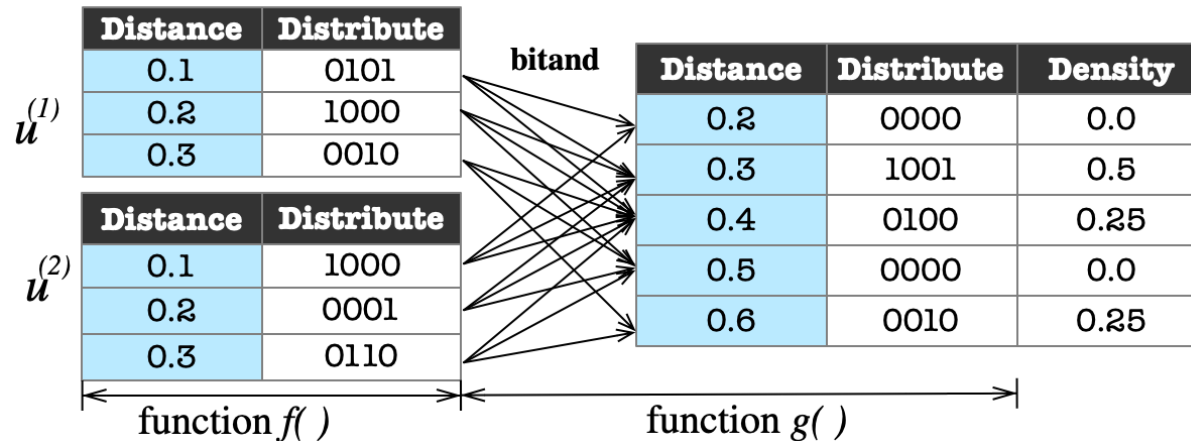
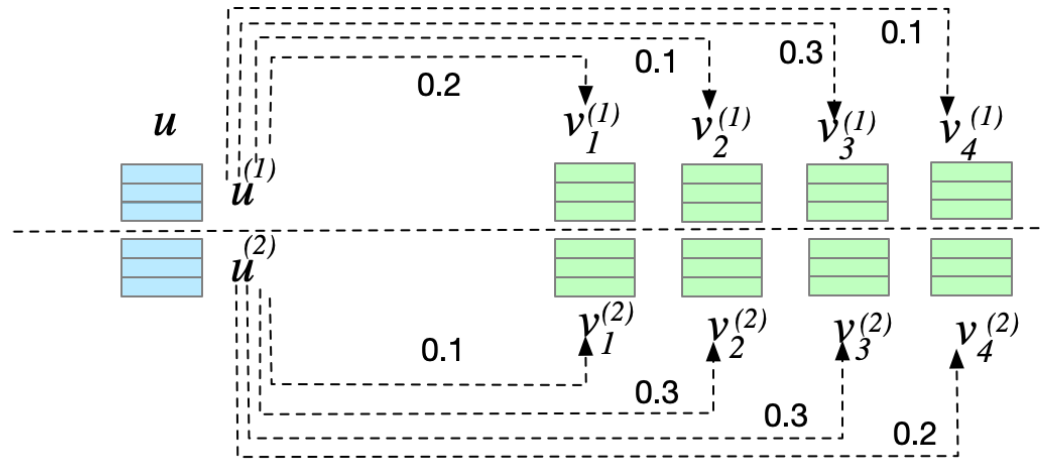
$\mathbf{q}$ : query

$\mathbf{x}_q$ : input vector

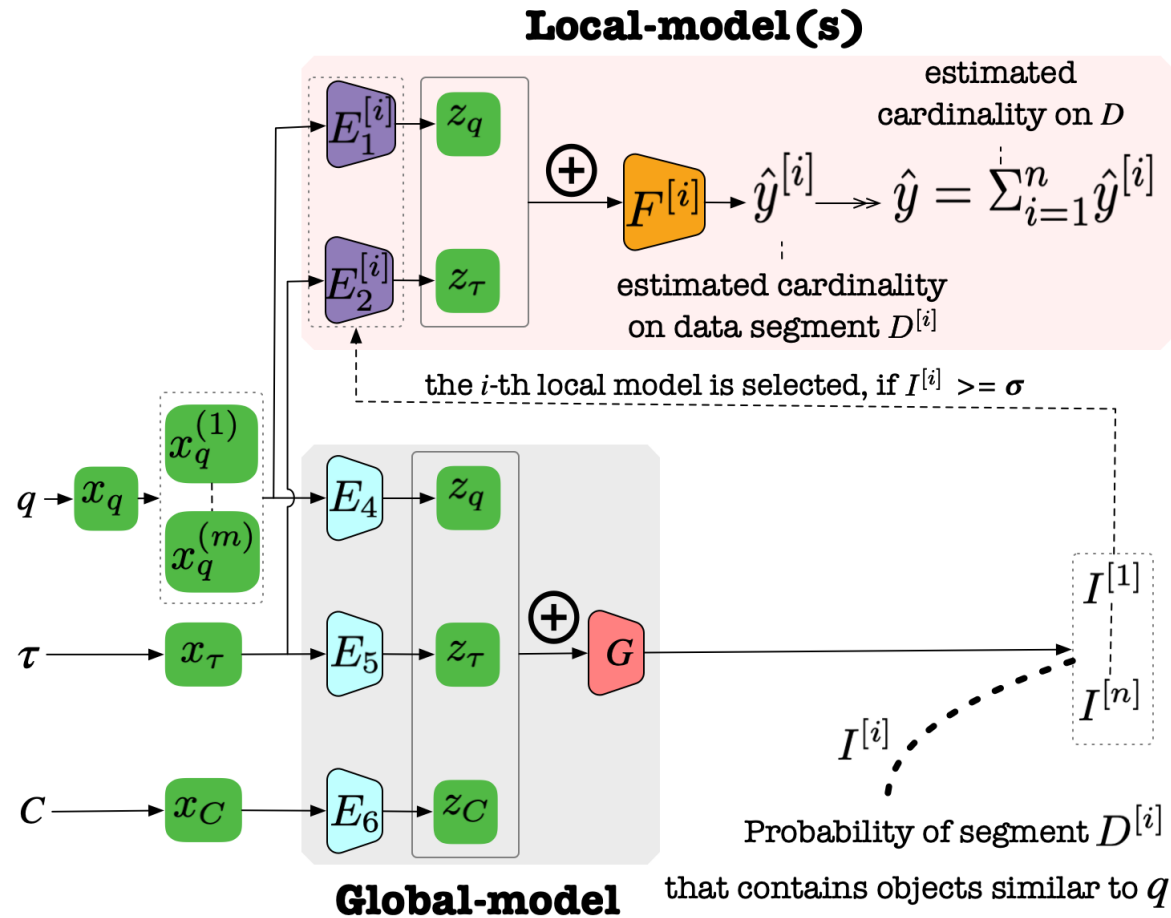
$x_q^{(i)}$ : the  $i$ -th segment of input vector

$\mathbf{e1}$ ,  $\mathbf{e2}$ ,  $\mathbf{el}$ : Neural Networks

# Query Segmentation



# Global-Local Model



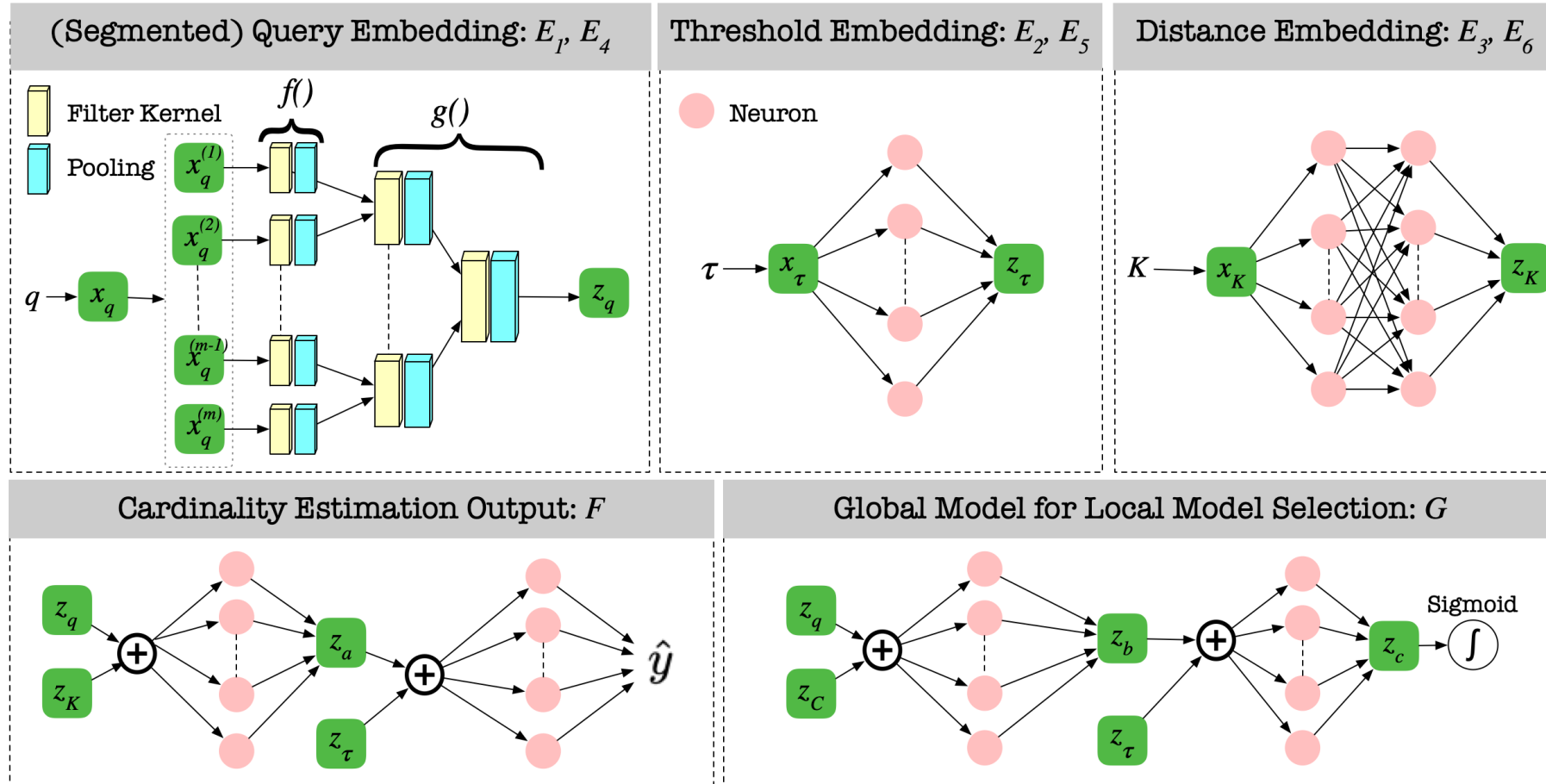
## Global Loss Function

$$\epsilon^{\{j\}[i]} = \frac{\text{card}^{\{j\}[i]} - \min_i \text{card}^{\{j\}[i]}}{\max_i \text{card}^{\{j\}[i]} - \min_i \text{card}^{\{j\}[i]}}$$

$$\mathcal{L}(\theta) = \frac{1}{n \times B_S} \sum_{i=1}^n \sum_{j=1}^{B_S} R^{\{j\}[i]} \log(I^{\{j\}[i]})(1 + \epsilon^{\{j\}[i]}) + (1 - R^{\{j\}[i]}) \log(1 - I^{\{j\}[i]})$$

$$\mathcal{J}(\theta) = - \frac{1}{n \times B_S} \sum_{i=1}^n \sum_{j=1}^{B_S} R^{\{j\}[i]} \log(I^{\{j\}[i]})(1 + \epsilon^{\{j\}[i]}) + (1 - R^{\{j\}[i]}) \log(1 - I^{\{j\}[i]})$$

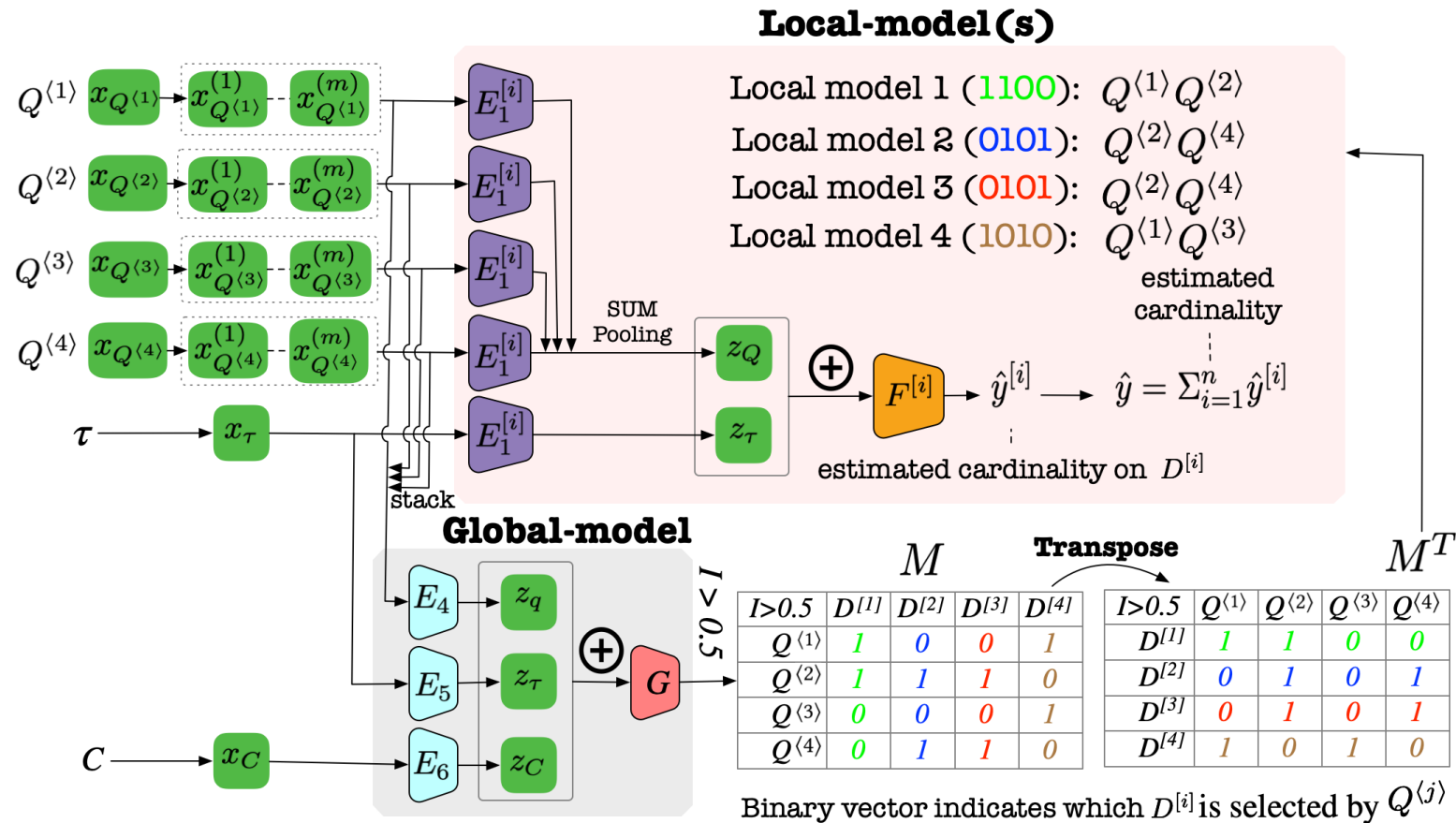
# Implementation Details



# Cardinality Estimation for Similarity Joins

- What if the query is a set of vectors (Joins) ?
  - $Q = \{q_1, q_2, q_3, q_4, \dots\}$
- A naïve way is to estimate for each vector and sum them up.
  - $\text{Card}(Q) = \text{Card}(q_1) + \text{Card}(q_2) + \text{Card}(q_3) + \text{Card}(q_4)$
  - Low efficiency

# Cardinality Estimation for Similarity Joins





# Experiments

- Datasets

<b>Dataset</b>	<b>Dimension</b>	<b>#Data</b>	<b>#Training</b>	<b>#Testing</b>	<b>Metric</b>	$\tau_{max}$
BMS	512	515,597	8,000	2,000	Jaccard	0.50
GloVe300	300	1,917,494	8,000	2,000	Angular	0.60
ImageNET	64	1,431,167	8,000	2,000	Hamming	0.90
Aminer	2,943	1,712,433	4,000	1,000	Edit	0.05
YouTube	1,770	346,194	2,400	600	Euclidean	0.15
DBLP	5,373	1,000,000	2,400	600	Edit	0.20

# Experiments

- Methods

<b>id</b>	<b>Method</b>	<b>Embed</b>	<b>Auto-tuning</b>	<b>Framework</b>	<b>Opt</b>	<b>Data Segment</b>
1	QES	CNN	No	Local	Select	No
2	Local+	CNN	Yes	Local	Select	Yes
3	GL-MLP	MLP	No	Global-Local	Select	Yes
4	GL-CNN	CNN	No	Global-Local	Select	Yes
5	GL+	CNN	Yes	Global-Local	Select	Yes
6	CardNet	VAE	No	Local	Select	No
7	Sampling	-	No	-	Select	No
8	Kernel-based	-	No	-	Select	No
9	MLP	MLP	No	Local	Select	No
10	SimSelect	-	-	-	Select	-
11	CNNJoin	CNN	No	Local	Join	No
12	GLJoin	MLP	No	Global-Local	Join	Yes
13	GLJoin+	CNN	Yes	Global-Local	Join	Yes

# Experiments

- Query
  - Vectors: 80% training, 20% testing
  - Threshold: selectivity lower than 1%
  - Join Size: [1-100) training, [50-100), [100-150), [150,200) testing
- Environment
  - Intel(R) Xeon(R) CPU E5-2630v4@2.20GHz
  - 128 Gigabytes memory
  - PyTorch 1.0.1

# Experiments (Accuracy)

- Cardinality Estimation for Similarity Search

Dataset	Method	Mean	Median	90th	95th	99th	Max
BMS	GL+	2.34	1.09	2.47	4.32	19.7	111
	Local+	2.37	1.05	2.51	4.36	18.4	98.3
	Sampling (10%)	5.18	1.83	11.2	17.4	55.0	165
	GL-CNN	3.50	2.42	8.21	10.6	15.7	291
	GL-MLP	4.41	3.02	9.78	12.8	19.7	439
	QES	7.27	5.05	16.5	21.6	32.2	644
	CardNet	12.4	5.16	31.3	48.8	99.1	335
	MLP	11.2	8.03	36.8	47.7	71.0	700
	Kernel-based	12.8	8.81	29.7	39.2	59.5	135
	Sampling (equal)	12.3	7.0	31.0	41.0	74.0	111
	Sampling (1%)	19.6	13.0	55.0	66.9	74.0	200

Dataset	Method	Mean	Median	90th	95th	99th	Max
Aminer	GL+	1.54	1.07	2.05	2.98	7.79	152
	Local+	1.61	1.12	2.36	3.01	6.46	321
	Sampling (10%)	2.41	1.72	3.90	5.26	14.2	31.0
	GL-CNN	1.83	1.27	4.21	5.39	8.38	154
	GL-MLP	3.09	2.14	7.10	9.18	14.2	290
	QES	5.22	3.63	11.9	15.4	24.4	541
	CardNet	5.45	2.05	7.59	12.9	43.1	3526
	MLP	8.39	5.80	19.4	25.1	38.6	780
	Kernel-based	9.85	6.91	22.6	28.7	44.6	117
	Sampling (equal)	66.5	42.0	182	245	245	245
	Sampling (1%)	19.5	4.20	56.0	75.0	136	245

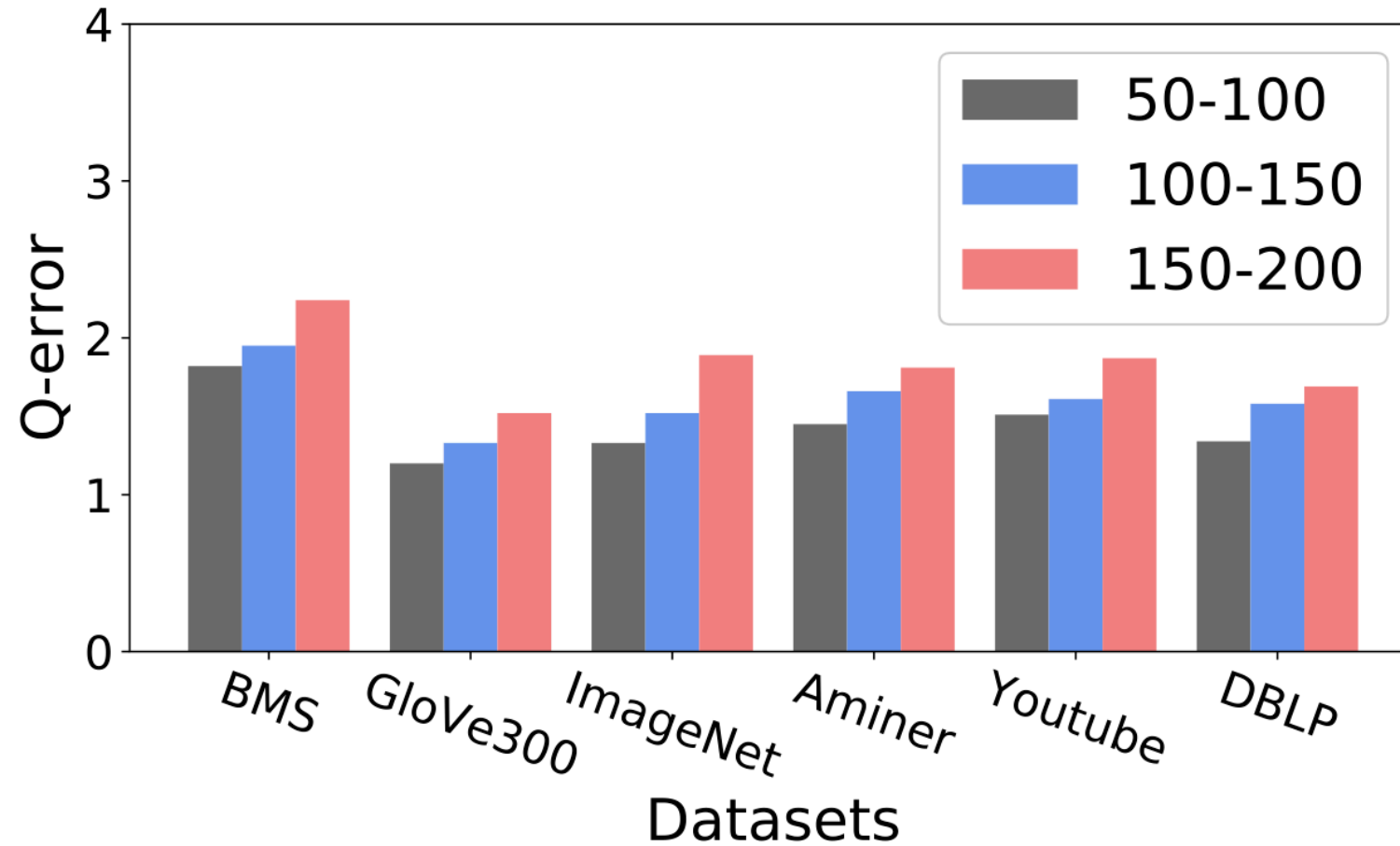
# Experiments (Accuracy)

- Cardinality Estimation for Similarity Join

Dataset	Method	Mean	Median	90th	95th	99th	Max	Dataset	Method	Mean	Median	90th	95th	99th	Max
BMS	GLJoin+	<b>1.87</b>	<b>1.31</b>	<b>4.31</b>	<b>5.51</b>	<b>8.55</b>	174	Aminer	GLJoin+	<b>1.42</b>	<b>1.08</b>	3.26	4.16	6.26	121
	GL+	2.01	1.36	4.59	6.12	9.34	205		GL+	1.70	1.18	3.95	5.10	7.94	171
	Sampling (10%)	3.99	2.18	8.46	13.5	23.1	<b>37.0</b>		Sampling (10%)	2.06	1.90	<b>2.90</b>	<b>3.35</b>	<b>4.57</b>	<b>5.12</b>
	GLJoin	2.51	1.72	5.78	7.56	11.5	265		GLJoin	2.02	1.40	4.66	5.94	9.25	193
	CNNJoin	5.63	3.90	12.9	16.9	26.2	508		CNNJoin	6.58	4.67	15.2	19.6	30.5	788
	CardNet	8.35	5.88	19.1	25.2	37.2	857		CardNet	5.16	3.55	11.7	15.2	24.3	766
	Sampling (equal)	19.3	2.50	15.2	40.9	302	451		Sampling (equal)	124	7.77	371	501	909	1221
	Sampling (1%)	144	3.86	451	800	1505	2701		Sampling (1%)	5.96	1.94	3.98	5.21	86.2	151

# Experiments (Accuracy)

- Cardinality Estimation for Similarity Join



# Experiments (Efficiency)

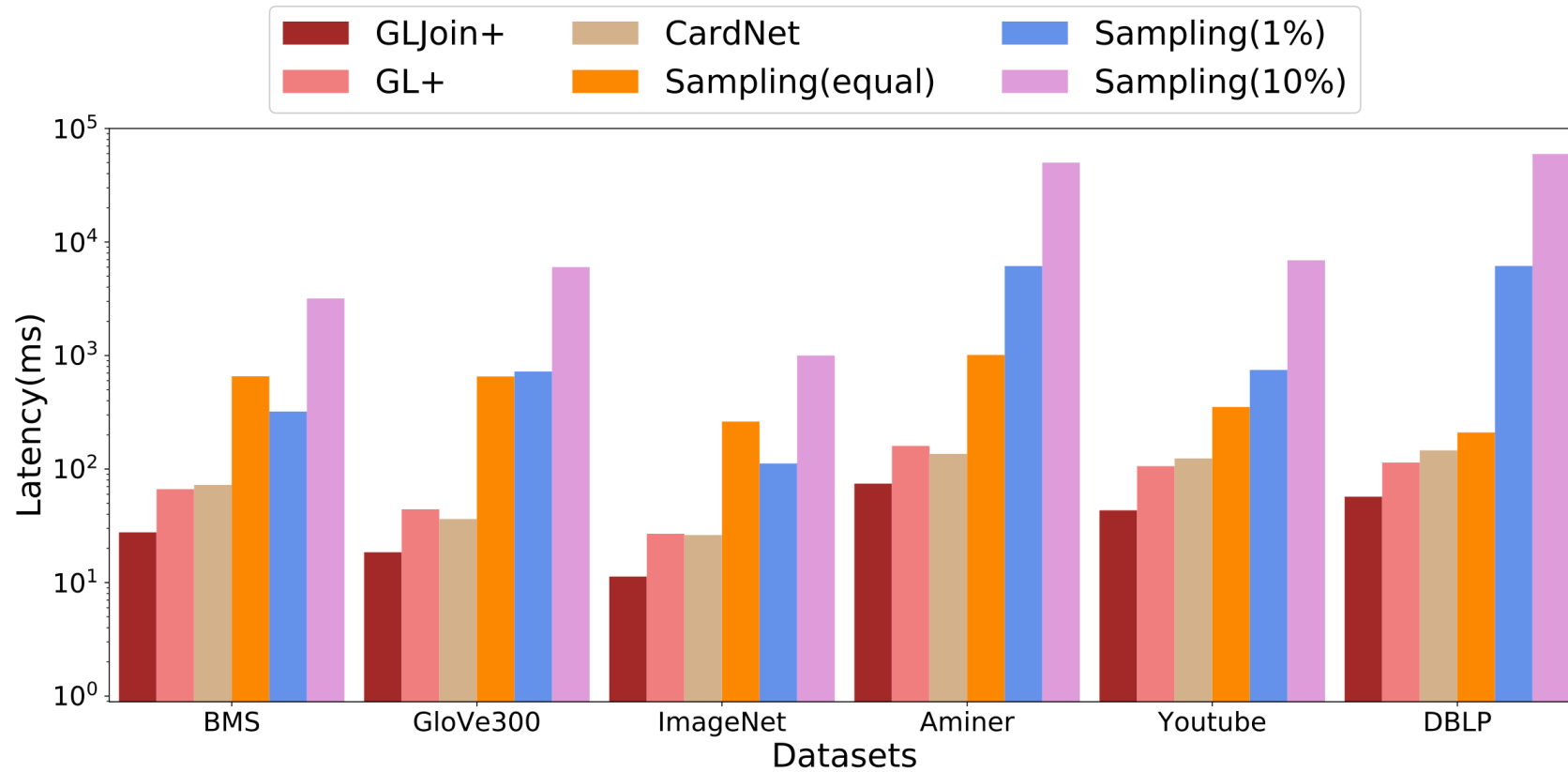
Model	BMS	GloVe300	ImageNET	Aminer	Youtube	DBLP
Sampling (1%)	12.7	27.7	3.66	243	24.5	239
MLP	4.11	3.09	3.21	9.01	8.23	15.3
QES	0.25	0.17	0.18	0.41	0.35	0.58
CardNet	38.8	35.3	16.2	54.5	52.8	55.1
GL-MLP	111	106	101	176	171	203
GL-CNN	29.2	21.3	7.32	35.6	32.1	55.6
GL+	28.3	22.1	7.51	34.2	30.7	50.1
GLJoin+	30.1	21.5	9.04	35.9	31.8	59.1

**Model Size (MB)**

Model	BMS	GloVe300	ImageNET	Aminer	Youtube	DBLP
SimSelect	3.96	12.1	5.22	5.87	12.5	18.6
Kernel-based	10.3	15.1	6.43	125	21.3	138
Sampling (10%)	30.9	70.1	10.5	587	69.5	598
Sampling (equal)	6.78	6.77	2.31	9.56	3.26	2.55
Sampling (1%)	3.21	7.23	1.12	61.4	7.46	61.5
CardNet	0.36	0.18	0.13	0.68	0.62	0.73
Local+	1.46	1.12	0.79	5.12	2.55	3.24
GL-MLP	0.51	0.65	0.28	3.43	2.35	3.69
GL-CNN	0.35	0.21	0.15	0.81	0.49	0.55
GL+	0.33	0.22	0.13	0.80	0.53	0.57
MLP	0.14	0.11	0.046	0.18	0.15	0.27
QES	0.015	0.012	0.007	0.042	0.021	0.032

**Estimation Efficiency for Similarity Search (Milliseconds)**

# Experiments (Efficiency)



Estimation Efficiency for Similarity Join (Milliseconds)



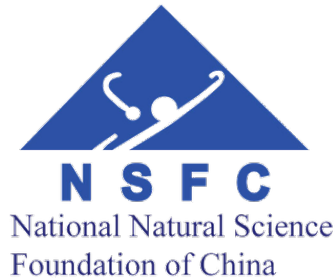
# Conclusion

**sun-j16@mails.tsinghua.edu.cn**

## **We make following contributions:**

- We propose a basic neural network model for cardinality estimation of similarity queries.
- We propose Query segmentation and Data segmentation to improve performance of model.
- We extend model to support similarity join.
- We conduct Comprehensive experiments on real datasets.

## **We special thank to:**



北京信息科学与技术  
国家研究中心  
BEIJING NATIONAL RESEARCH CENTER  
FOR INFORMATION SCIENCE AND TECHNOLOGY

