# GPU-Accelerated Join Selectivity Estimation using KDE Models
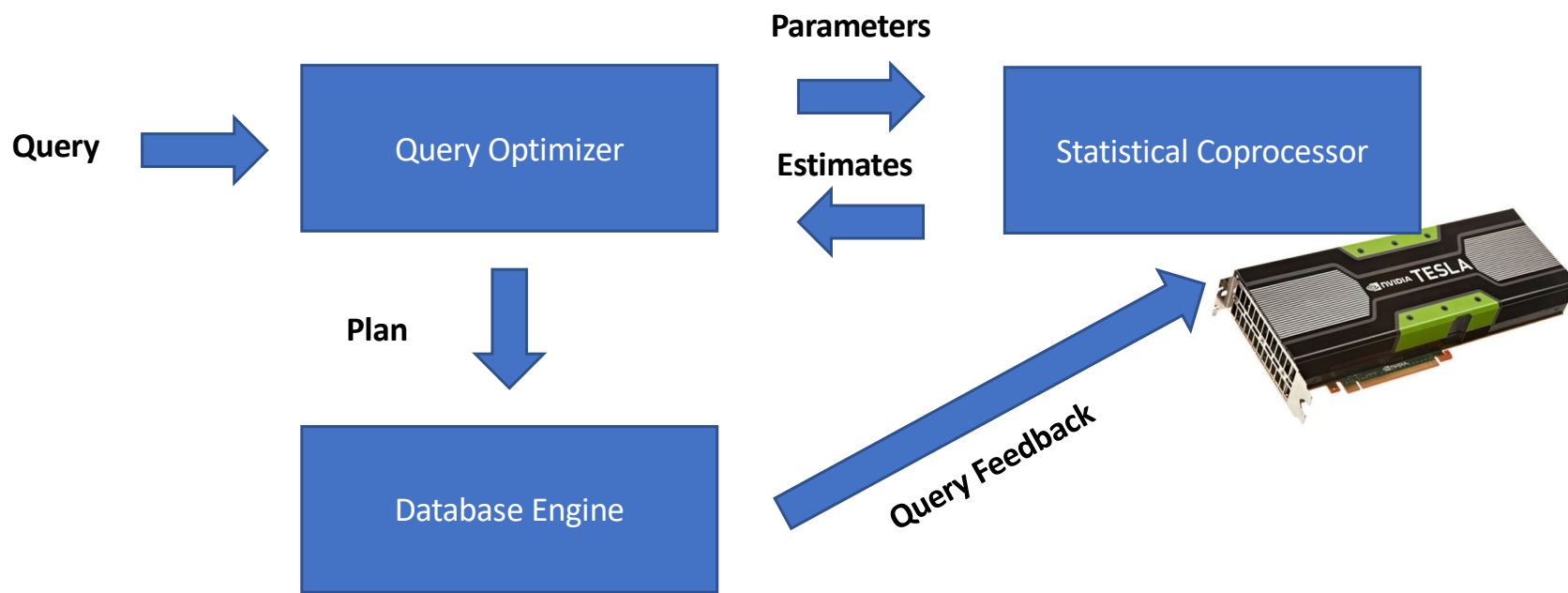
**Paper:**

*Estimating Join Selectivities using Bandwidth-Optimized Kernel Density Models,*

Martin Kiefer, Max Heimel, Sebastian Breß, Volker Markl
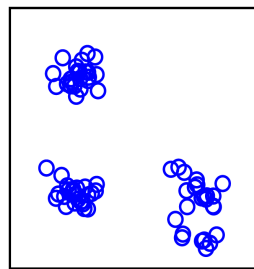
PVLDB, Volume 10 Issue 13, September 2017

Technische
Universität
Berlin

German
Research Center
for Artificial
Intelligence

# GPU-Accelerated Kernel Density Estimation for Join Selectivity Estimation

**Query** →

**Query Optimizer**

**Parameters** →

**Statistical Coprocessor**

← **Estimates**

**Plan** ↓

**Database Engine**

**Query Feedback** →

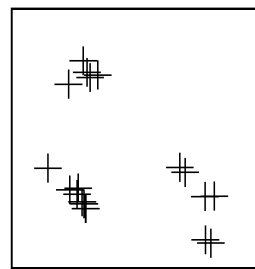# Background: Kernel Density Estimators

Average…        … over the kernel contributions
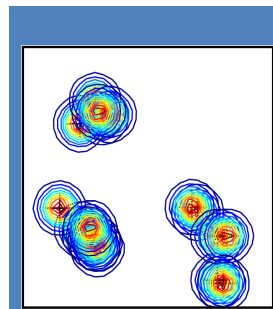
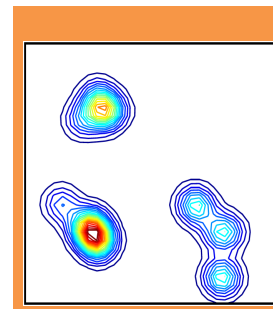$$\hat{P}_H(\vec{x}) = \frac{1}{|S|} \sum_{i=1}^{|S|} K_H(s_i, \vec{x})$$



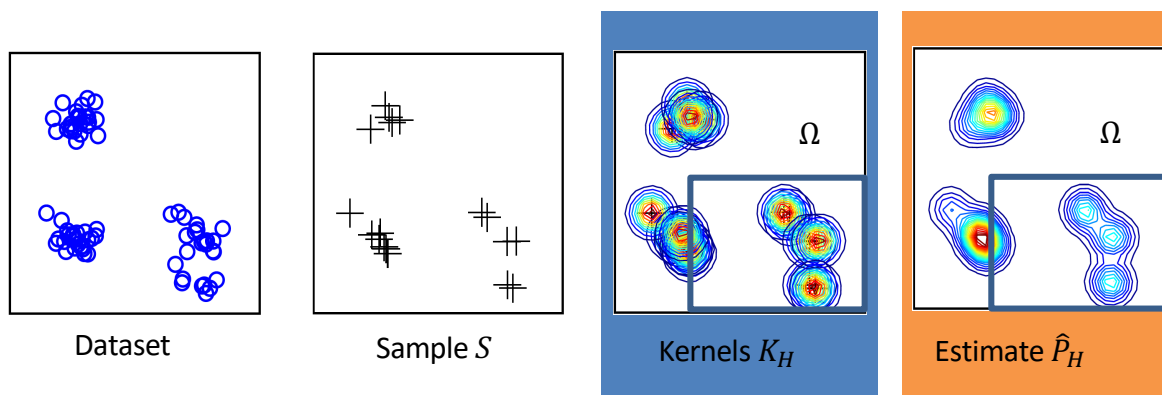Dataset        Sample $S$        Kernels $K_H$        Estimate $\hat{P}_H$
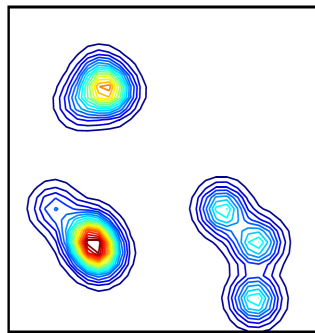
# Background: Kernel Density Estimators

Average...          ... over the kernel contributions

$$\text{sel}(\Omega) = \frac{1}{|S|} \sum_{i=1}^{|S|} \int_{\Omega} K_H(s_i, \vec{x})\, d\vec{x}$$



Dataset          Sample $S$          Kernels $K_H$          Estimate $\hat{P}_H$
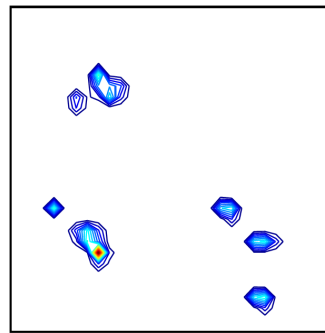
# Background: Kernel Density Estimators for Multi-Dimensional Selectivity Estimation [1]
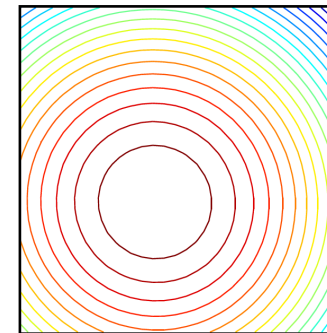
The bandwidth matrix $H$ controls the smoothing applied on the sample



Good fit          Overfit          Underfit

- Range selections over base tables
- Bandwidth optimization based on the estimation error
- Easy model maintenance

[1] Self-Tuning, GPU-Accelerated Kernel Density Models for Multidimensional Selectivity Estimation, SIGMOD'15

# The Problem:
# Multi-Dimensional Join Selectivity Estimation

$$Q = \sigma_{c_1}\left(R_1\right) \bowtie_{R_1.A_1 = R_2.A_1} \sigma_{c_2}\left(R_2\right)$$
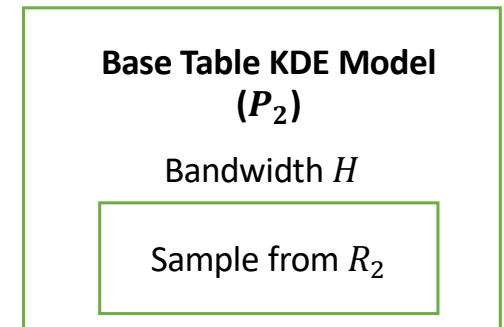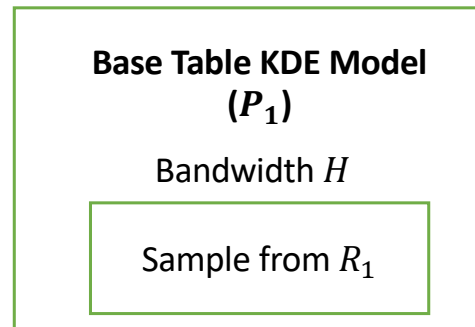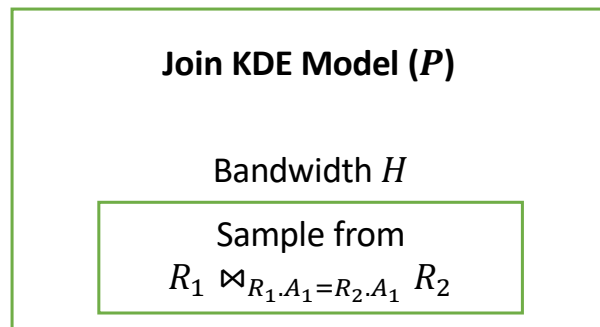
- and generalization to multiple joins
- **Databases:** Independence Assumption
  - Often violated
  - Introduce large errors, potentially bad query plans
- **Research:** Various Methods (e.g. Sampling, Sketches)
- **Our Approach:** Kernel Density Estimators

# Why KDEs for Join Selectivities?

- Multivariate Estimator
- No independence assumption
- Hybrid between samples and histograms
  - Small bandwidth: Sample evaluation
  - Increasing bandwidth: More smoothing, increasing bucket sizes
  - Bandwidth optimization selects proper bandwidth

# The Approach: Join and Base Table Models

$$Q = \sigma_{c_1}(R_1) \bowtie_{R_1.A_1 = R_2.A_1} \sigma_{c_2}(R_2)$$

**Join KDE Model ($P$)**

Bandwidth $H$

Sample from
$R_1 \bowtie_{R_1.A_1 = R_2.A_1} R_2$

**Base Table KDE Model ($P_1$)**

Bandwidth $H$

Sample from $R_1$

**Base Table KDE Model ($P_2$)**

Bandwidth $H$

Sample from $R_2$

**Compute:** $P(c_1 \wedge c_2)$

**Compute:** $\sum_{v \in A} P_1(A_1 = v \wedge c_1) \cdot P_2(A_2 = v \wedge c_2)$

Easy to evaluate, better estimates

Support for base table and join selectivities
Easy to construct and to maintain

# Table Model: Computation Components

$$Q = \sigma_{c_1}\left(R_1\right) \bowtie_{R_1.A_1 = R_2.A_1} \sigma_{c_2}\left(R_2\right)$$

Selectivity:

$$\frac{1}{s_1 \cdot s_2} \underbrace{\sum_{i=1,j=1}^{s_1,s_2}}_{\substack{\text{Sum over cross} \\ \text{product of two} \\ \text{samples}}} \underbrace{\hat{p}_1^{(i)}\left(c_1\right) \cdot \hat{p}_2^{(j)}\left(c_2\right)}_{\substack{\textbf{Invariant Contributions:} \\ \text{Contribution of sample} \\ \text{points wrt. selection} \\ \text{predicate}}} \cdot \underbrace{\hat{J}_{i,j}}_{\substack{\textbf{Cross Contribution:} \\ \text{Distance function on join} \\ \text{attributes of sample points}}}$$

48

# Table Model: Sample Pruning

$$\frac{1}{s_1 \cdot s_2} \sum_{i=1,j=1}^{s_1,s_2} \boxed{\hat{p}_1^{(i)}(c_1)} \cdot \hat{p}_2^{(j)}(c_2) \cdot \hat{J}_{i,j}$$
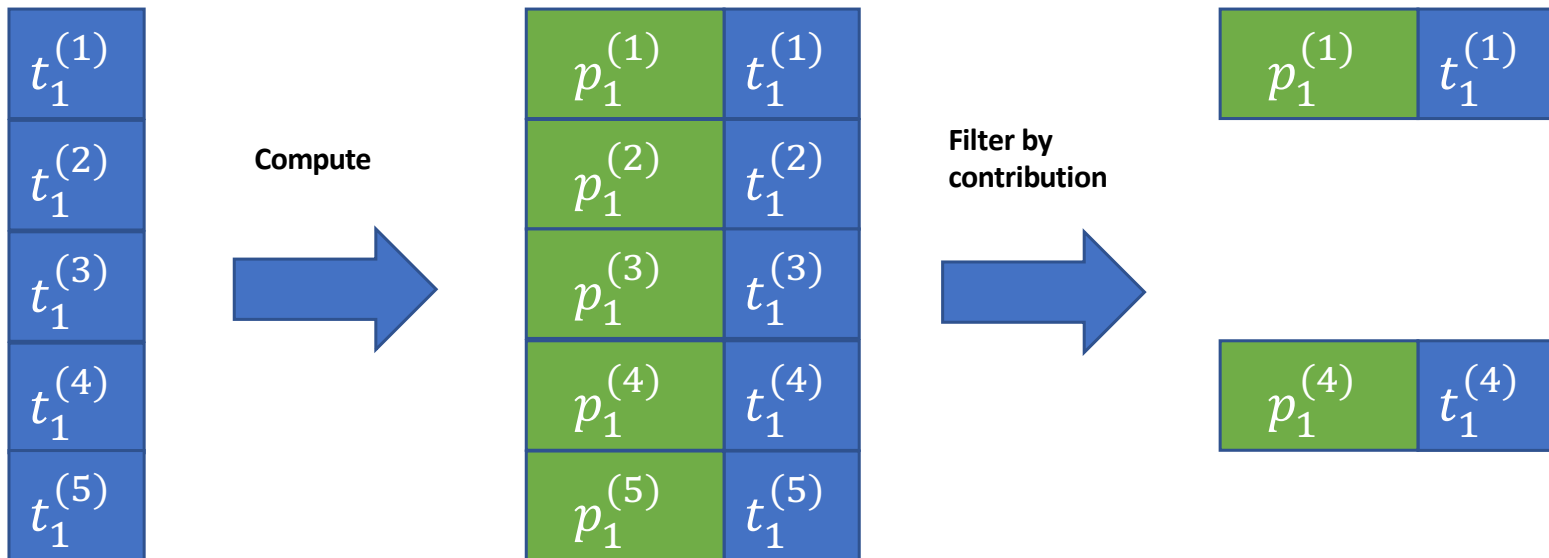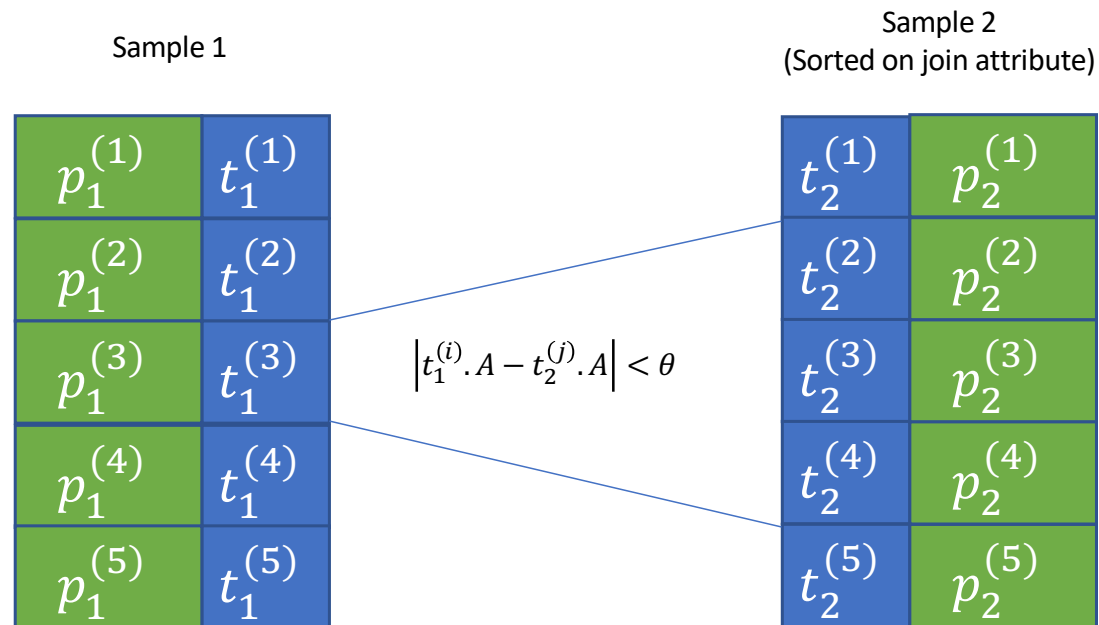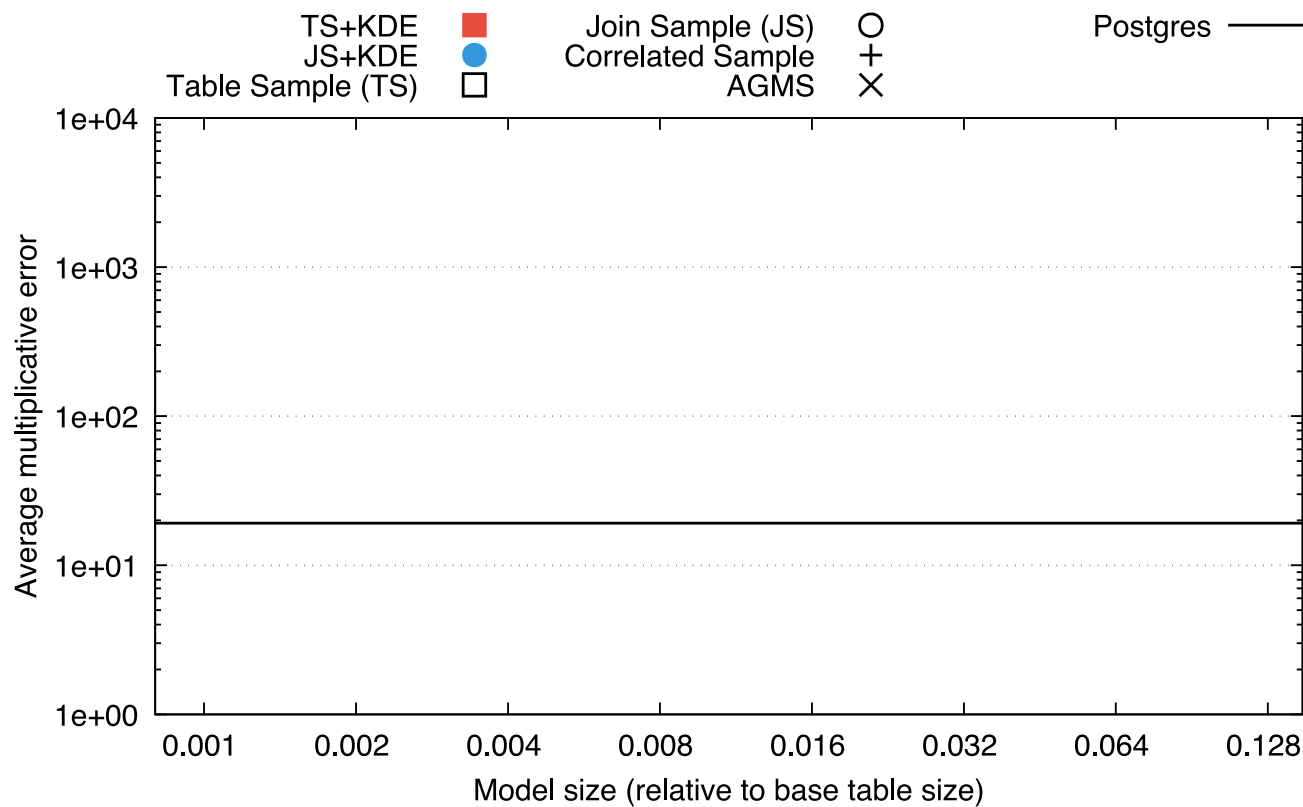
**Sample 1**

# Table Model: Cross Pruning

$$\frac{1}{s_1 \cdot s_2} \sum_{i=1,j=1}^{s_1,s_2} \hat{p}_1^{(i)}(c_1) \cdot \hat{p}_2^{(j)}(c_2) \cdot \boxed{\hat{J}_{i,j}}$$
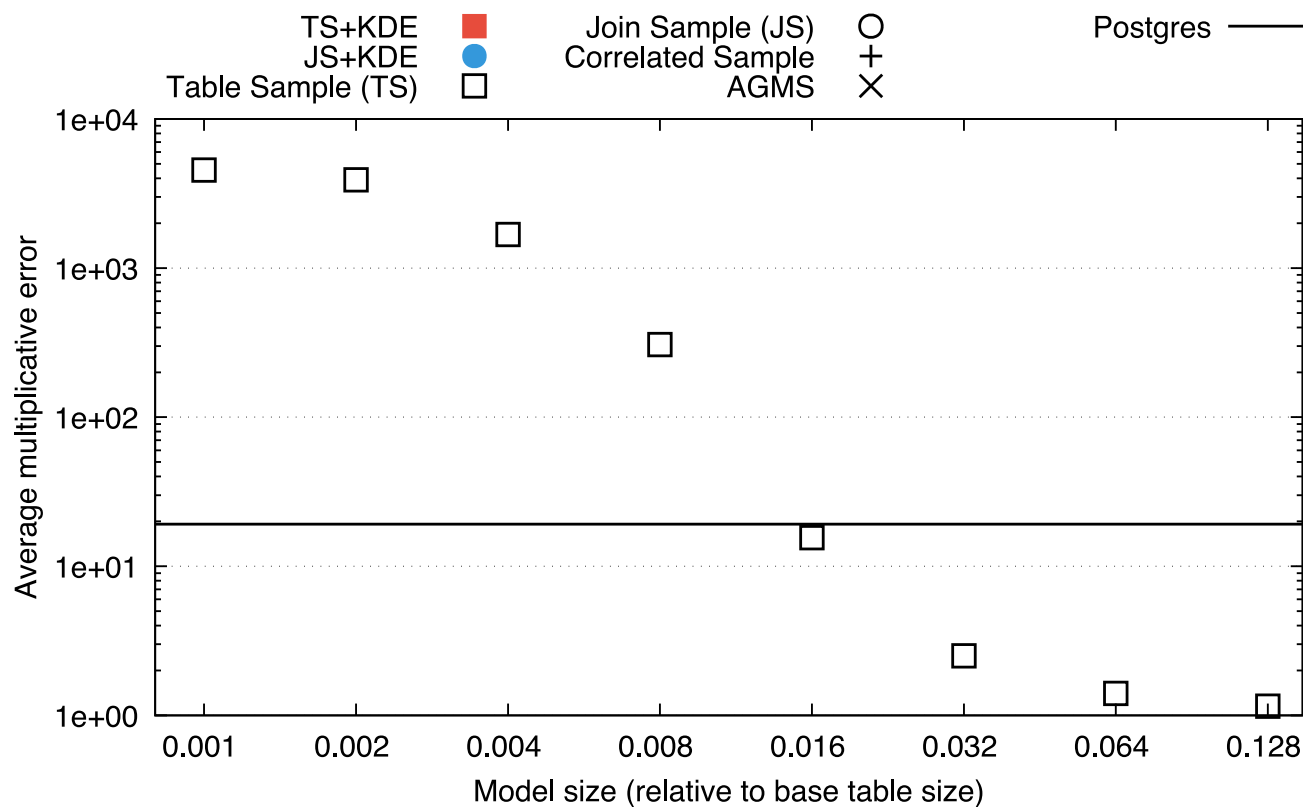
Sample 1

Sample 2
(Sorted on join attribute)



$$\left| t_1^{(i)}.A - t_2^{(j)}.A \right| < \theta$$

50

# Evaluation: Scaling the Model Size (Postgres)
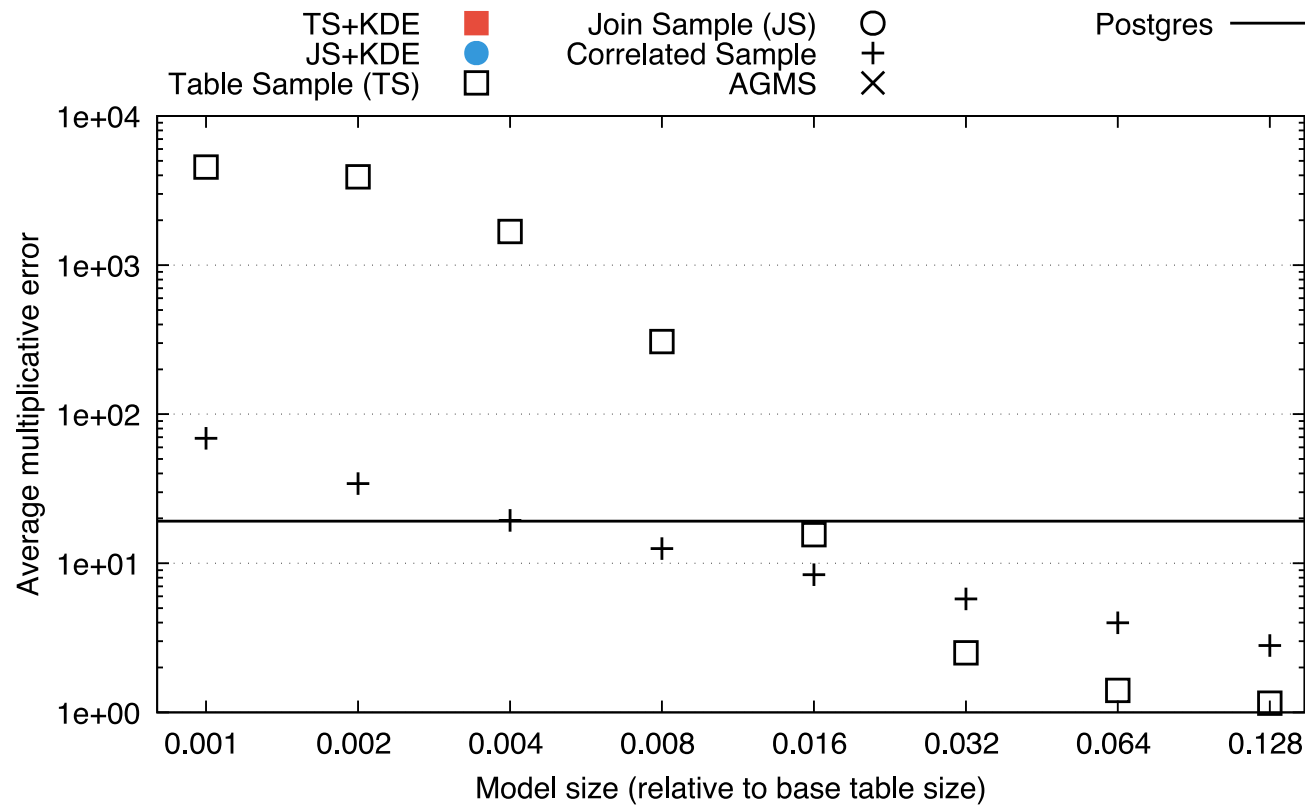
**Dataset:** DMV
**Query:** Q1U

# Evaluation: Scaling the Model Size (Table Sample)
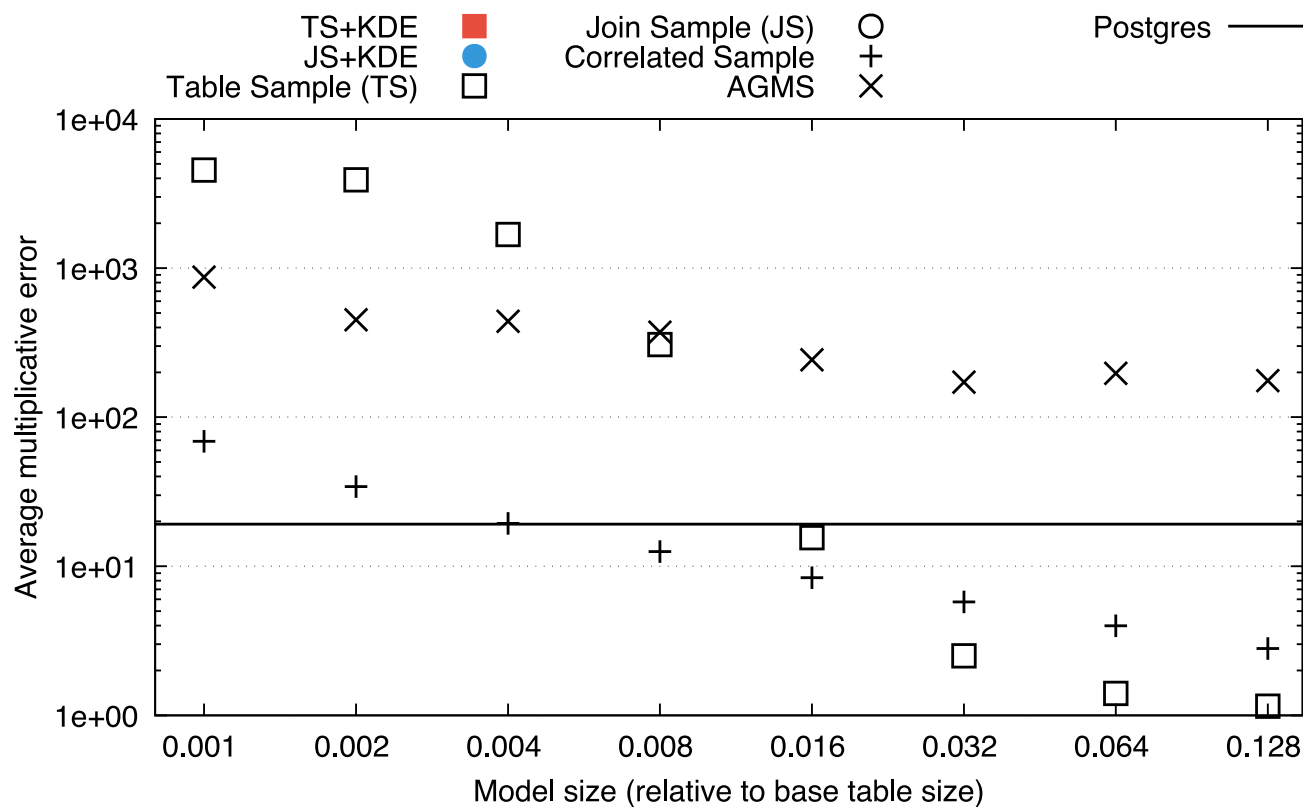
**Dataset:** DMV
**Query:** Q1U

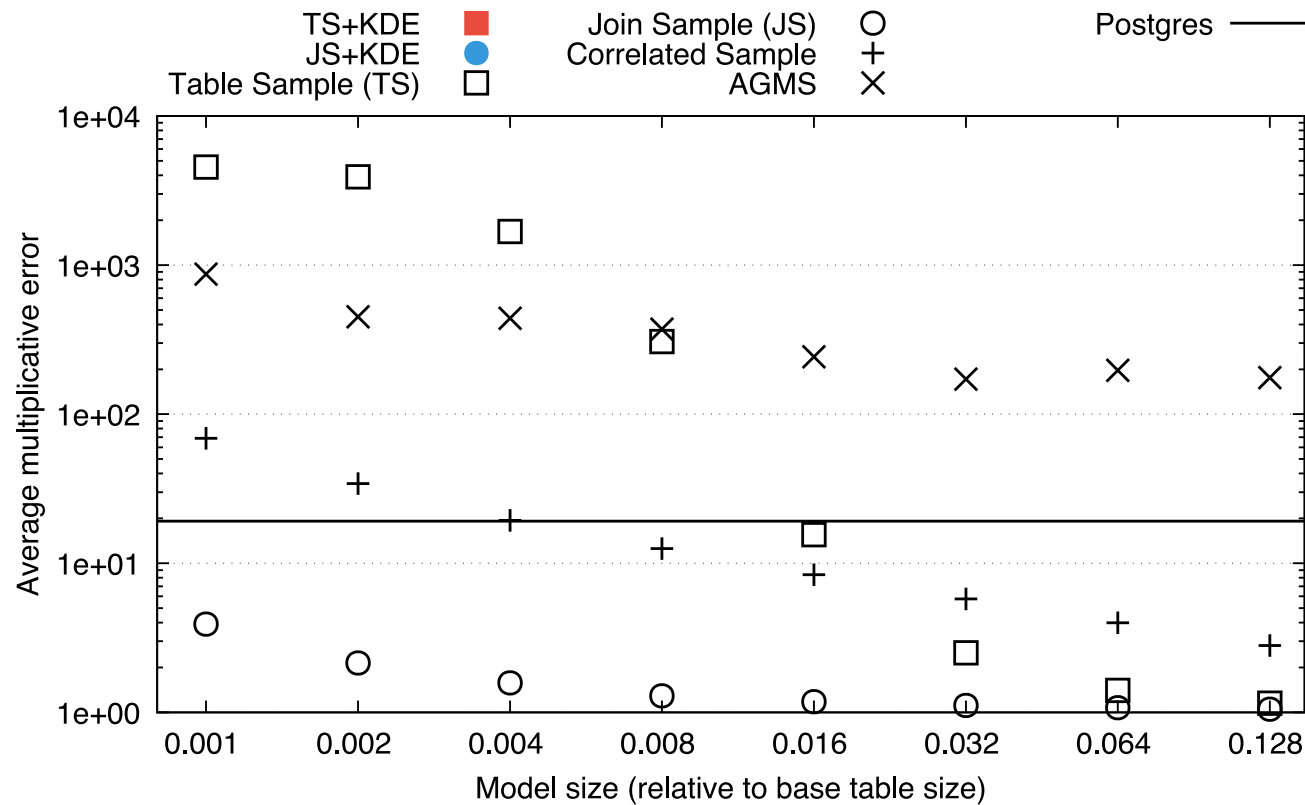# Evaluation: Scaling the Model Size (Correlated Sample)

**Dataset:** DMV
**Query:** Q1U

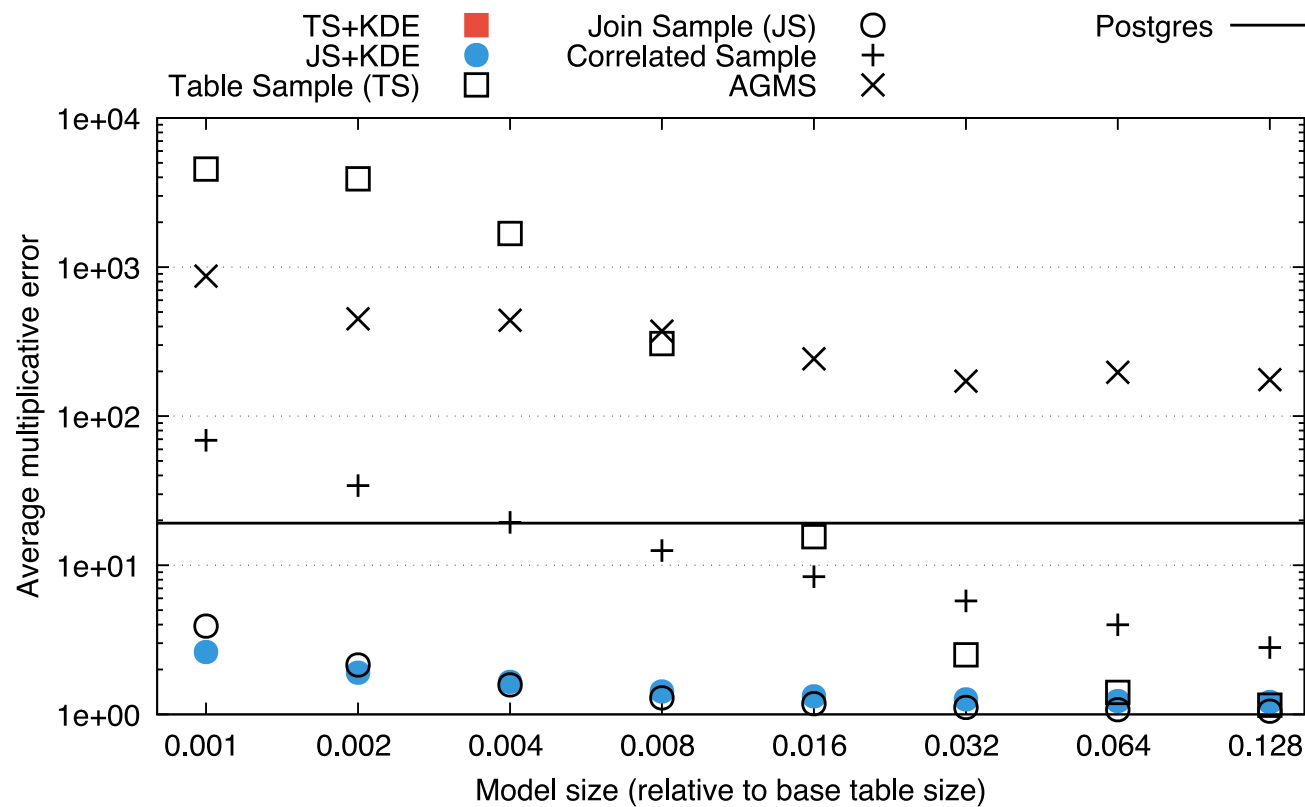# Evaluation: Scaling the Model Size (AGMS Sketch)

**Dataset:** DMV
**Query:** Q1U

# Evaluation: Scaling the Model Size
# (Join Sample)
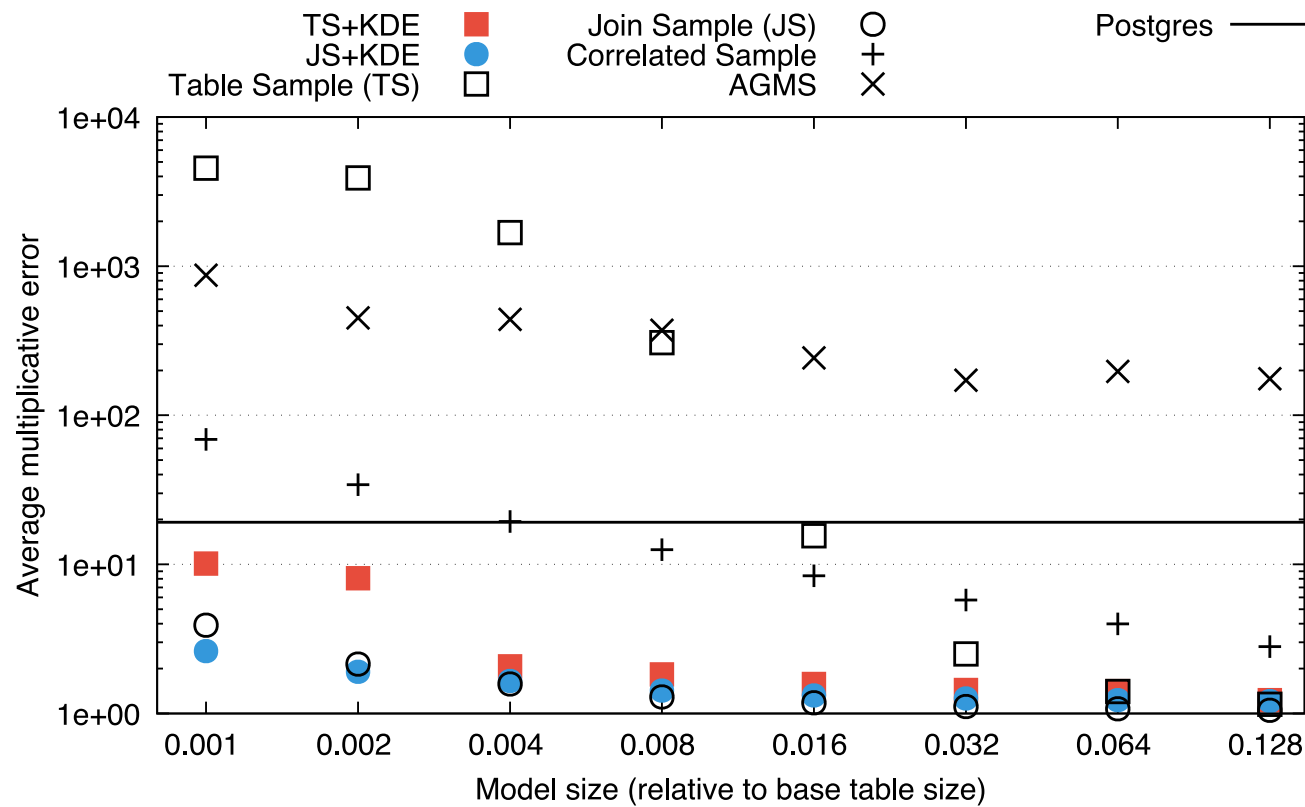
**Dataset:** DMV
**Query:** Q1U

# Evaluation: Scaling the Model Size
# (Join Sample + KDE)
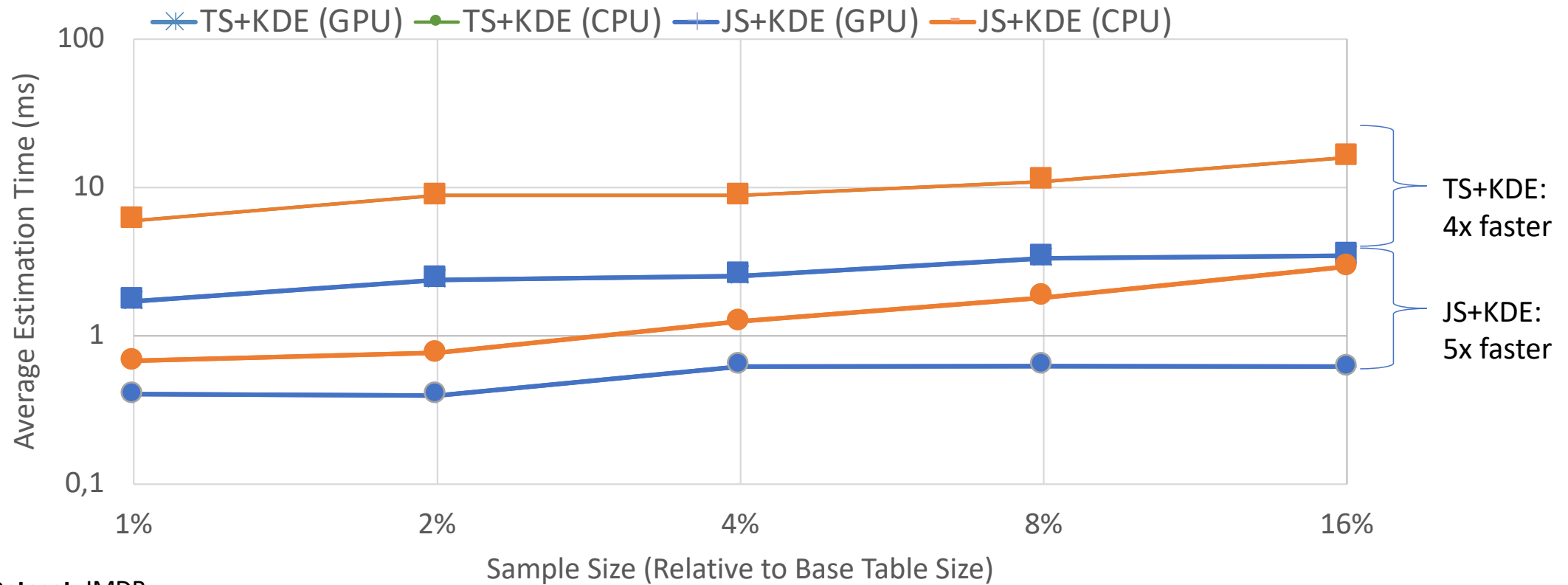
**Dataset:** DMV
**Query:** Q1U

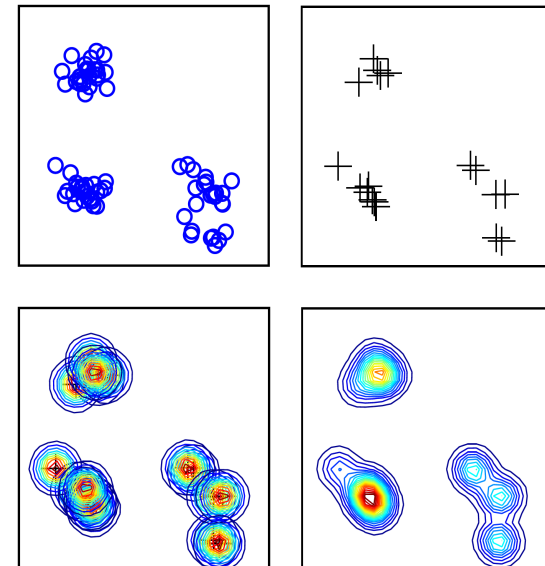# Evaluation: Scaling the Model Size
# (Table Sample + KDE)

**Dataset:** DMV
**Query:** Q1U

# Runtime: CPU vs GPU



**Dataset:** IMDB
**Workload:** Q1U
**GPU:** Tesla V100
**CPU:** Intel Xeon Gold 5115

# Conclusion

- KDE models for join selectivity estimation

- "Getting most out of your sample"

- Based on join or base table KDE models

- Learning hybrid between histograms and samples

- GPU-acceleration possible

- Experiments, data, and code online



github.com/martinkiefer/join-kde

"Estimating Join Selectivities using Bandwidth-Optimized Kernel Density Models", PVLDB 17

# Estimating Join Selectivities using Bandwidth-Optimized Kernel Density Models

**Martin Kiefer,** Max Heimel, Sebastian Breß, Volker Markl

**Further Publications on GPU-Accelerated Kernel Density Estimation:**

Self-Tuning, GPU-Accelerated Kernel Density Models for Multidimensional Selectivity Estimation

SIGMOD 2015

Demonstrating Transfer-Efficient Sample Maintenance on Graphics Cards

EDBT 2015

Technische Universität Berlin

German Research Center for Artificial Intelligence