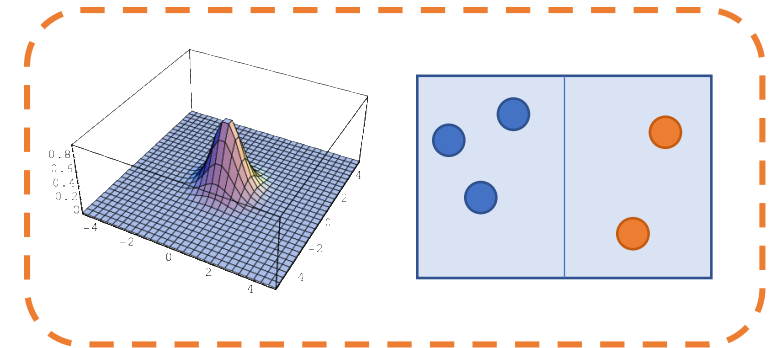# Scalable Kernel Density Classification via Threshold-Based Pruning
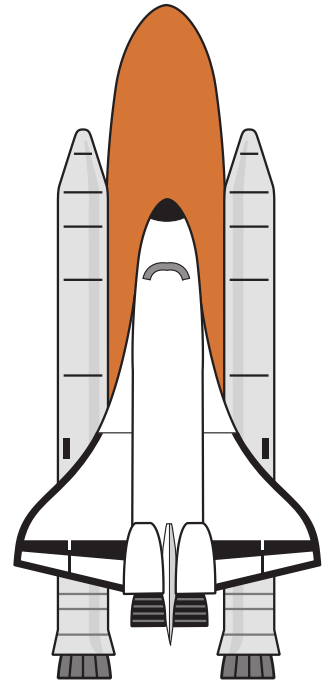
Edward Gan & Peter Bailis

# MacroBase: Analytics on Fast Streams

- Increasing Streaming Data
  - Manufacturing, Sensors, Mobile
  - Multi-dimensional + Latent anomalies

- Running in production
  - see CIDR17, SIGMOD17

- End-to-end operator cascades for:
  - Feature Transformation
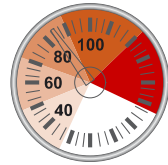  - Statistical Classification
  - Data Summarization

# Example: Space Shuttle Sensors

8 Sensors Total

"Fuel Flow"

"Flight Speed"
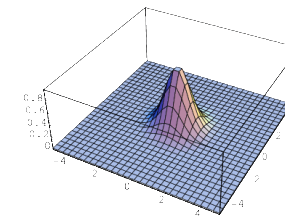
| Speed | Flow | Status |
|-------|------|-----------|
| 28 | 27 | Fpv Close |
| 34 | 43 | High |
| 52 | 30 | Rad Flow |
| 28 | 40 | Rad Flow |
| ... | ... | |

[UCI Repository]

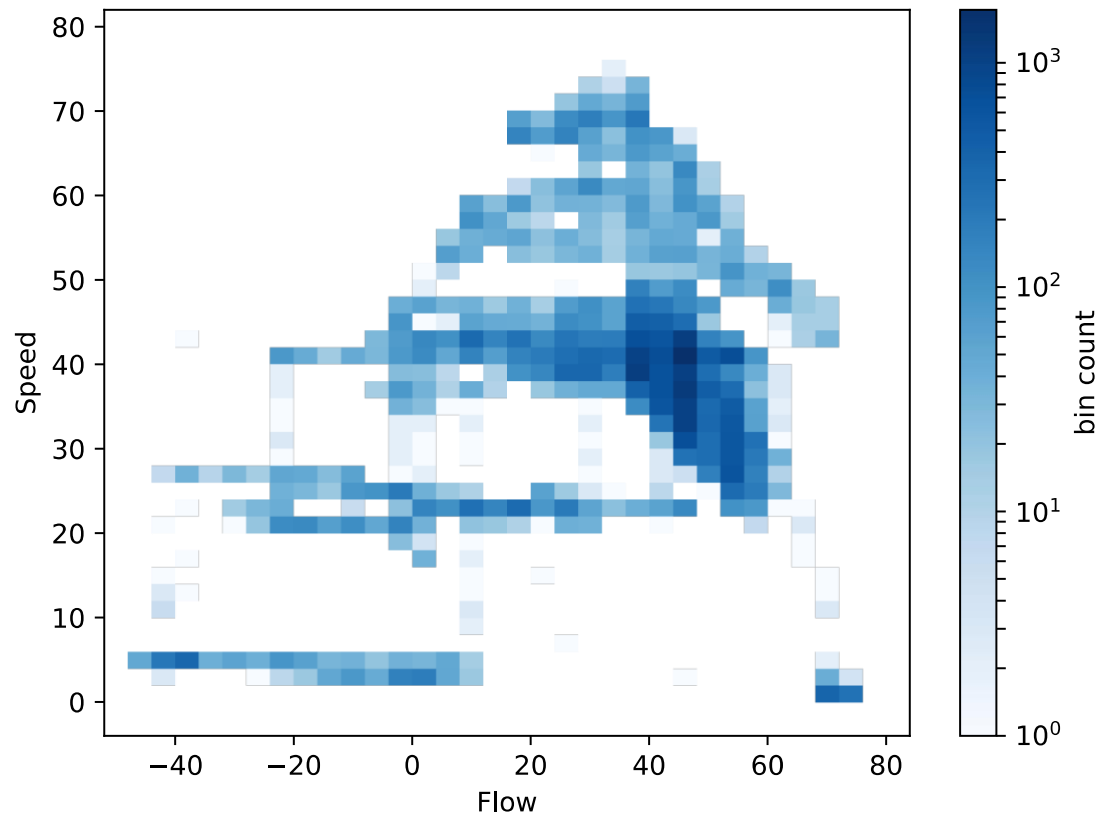End-Goal: Explain anomalous speed / flow measurements.

Problem: Model distribution of speed / flow measurements.
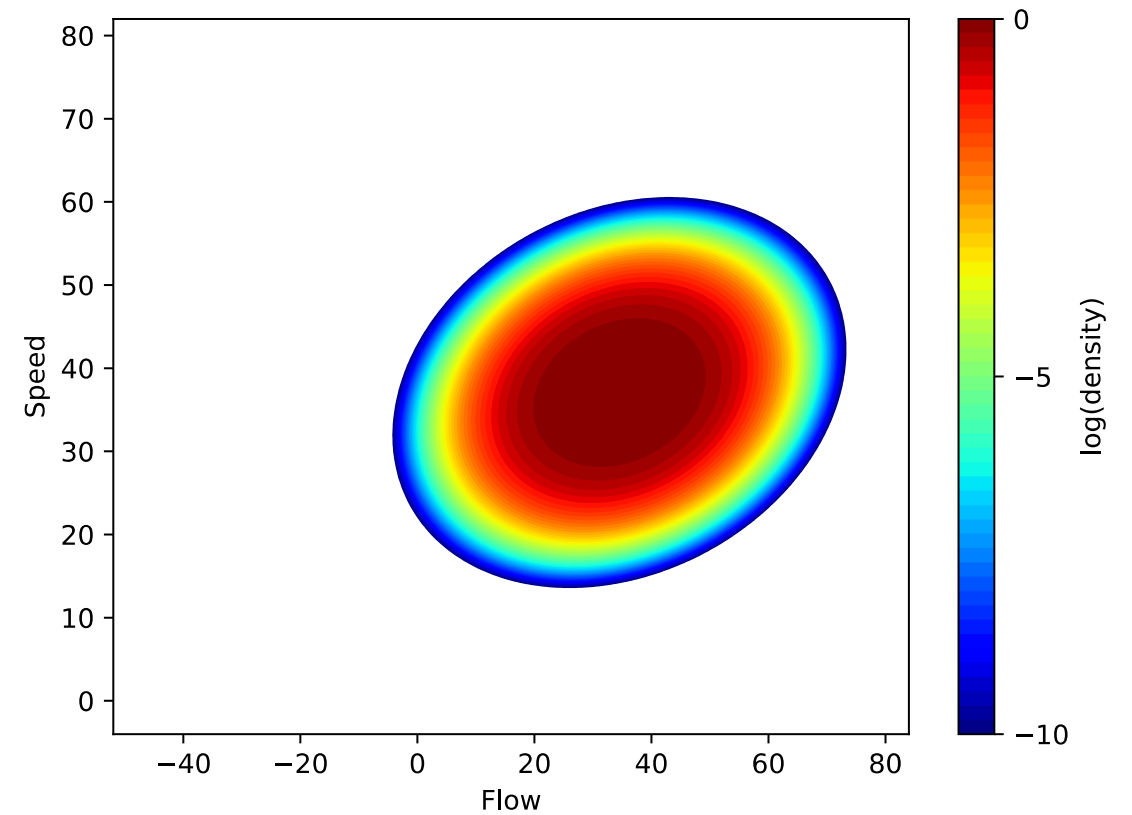
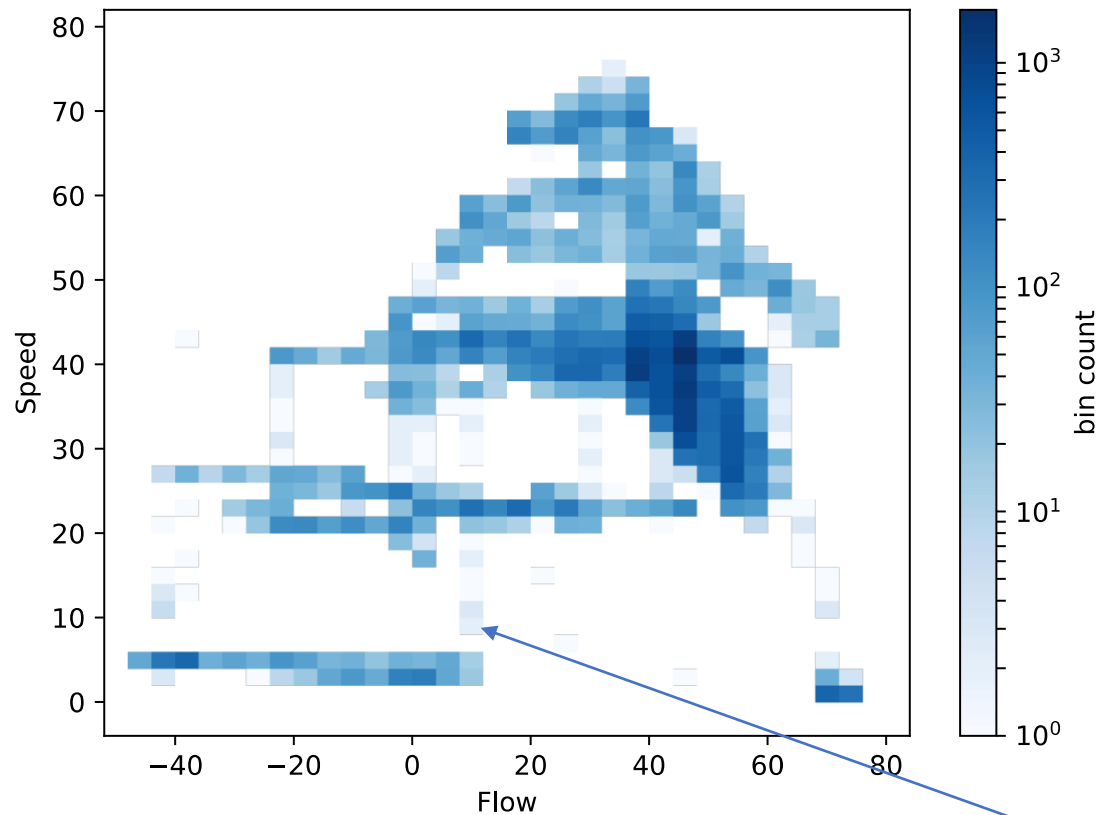# Difficulties in Data Modelling

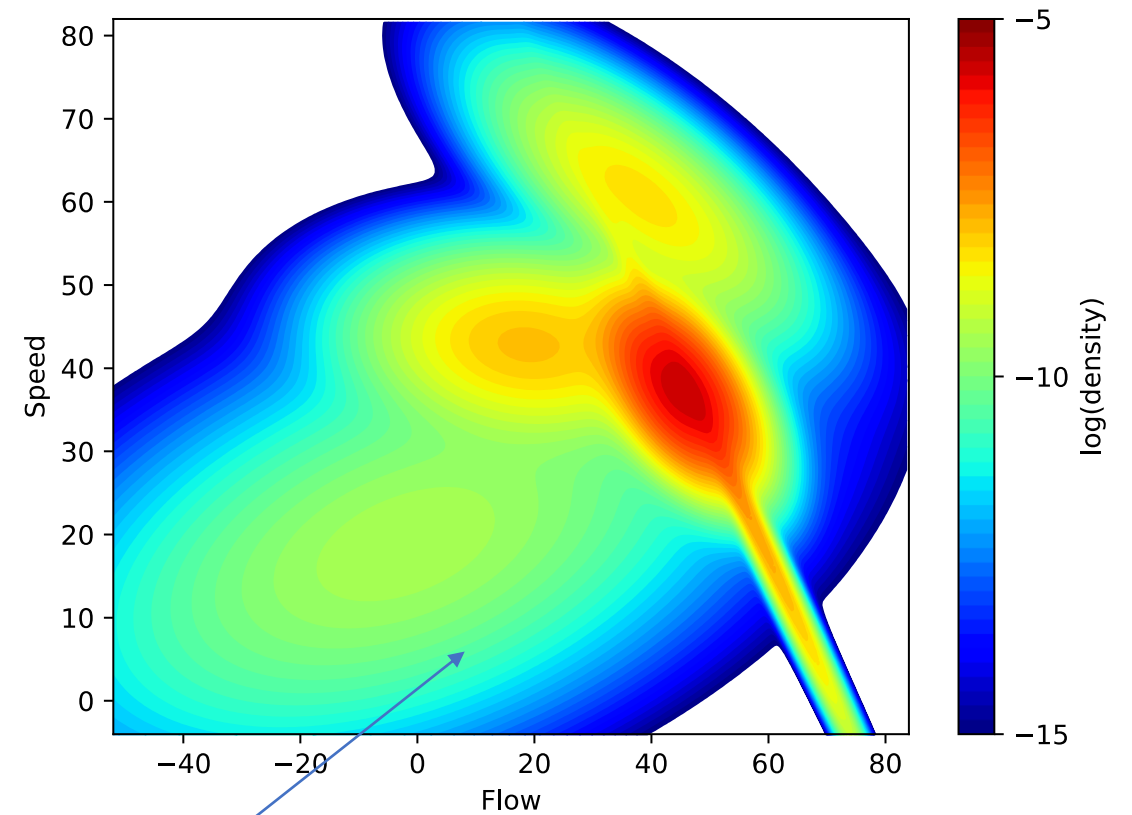Data Histogram

Gaussian Model



Poor Fit

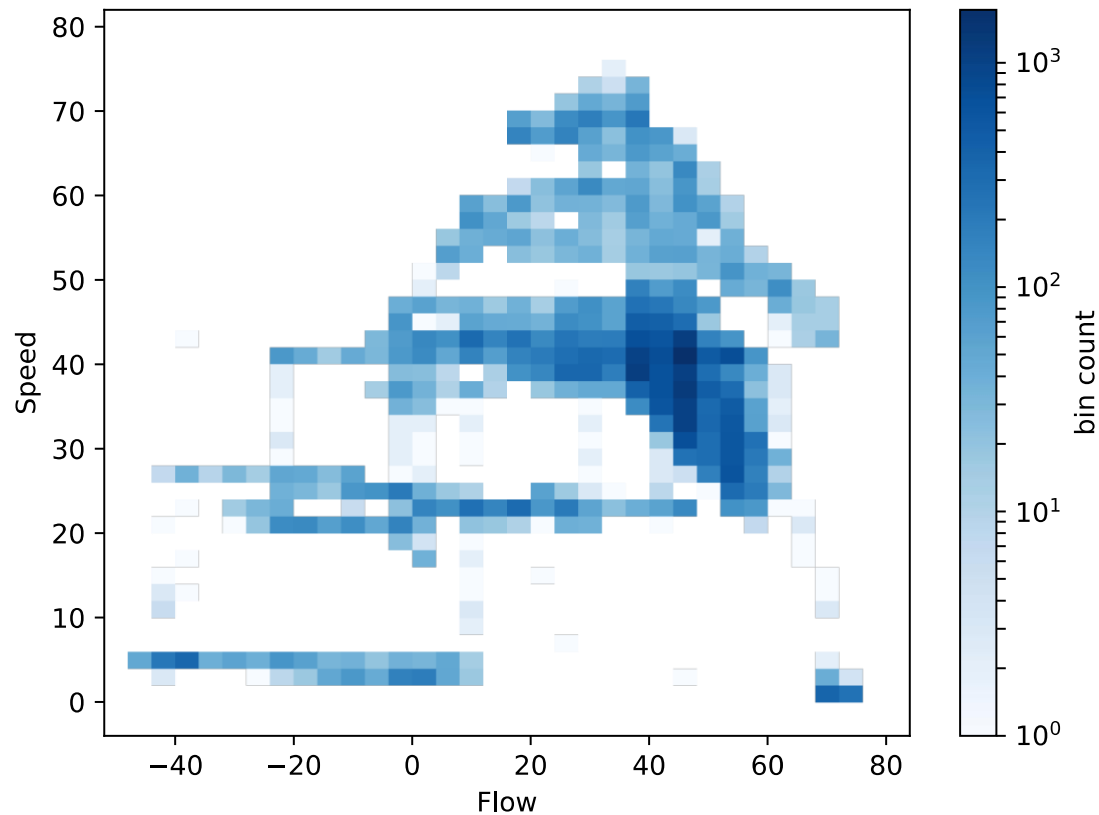# Difficulties in Data Modelling
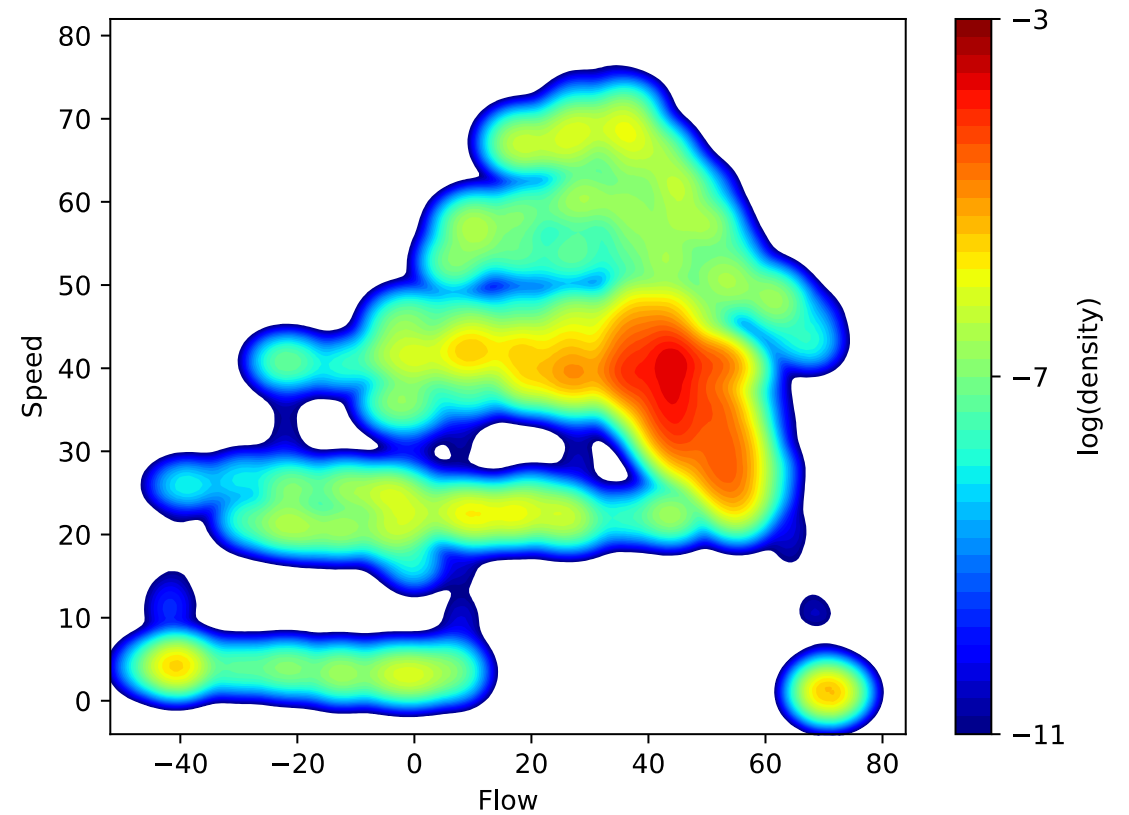


Data Histogram

Mixture of Gaussians

Inaccurate: Gaps not captured

# Kernel Density Estimation (KDE)



Data Histogram

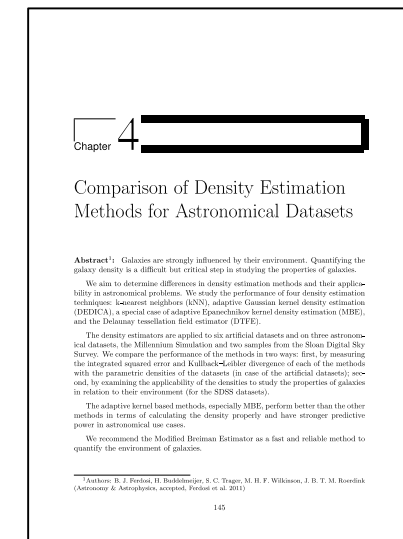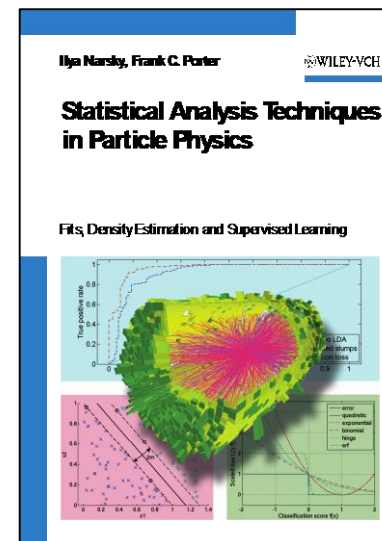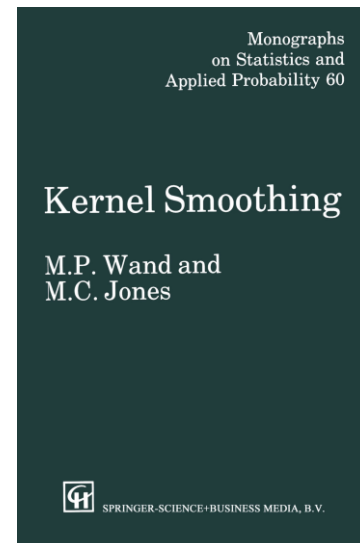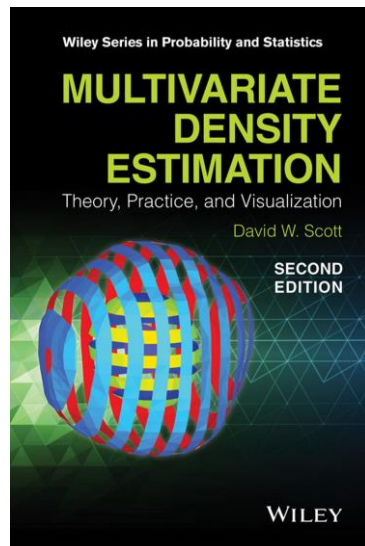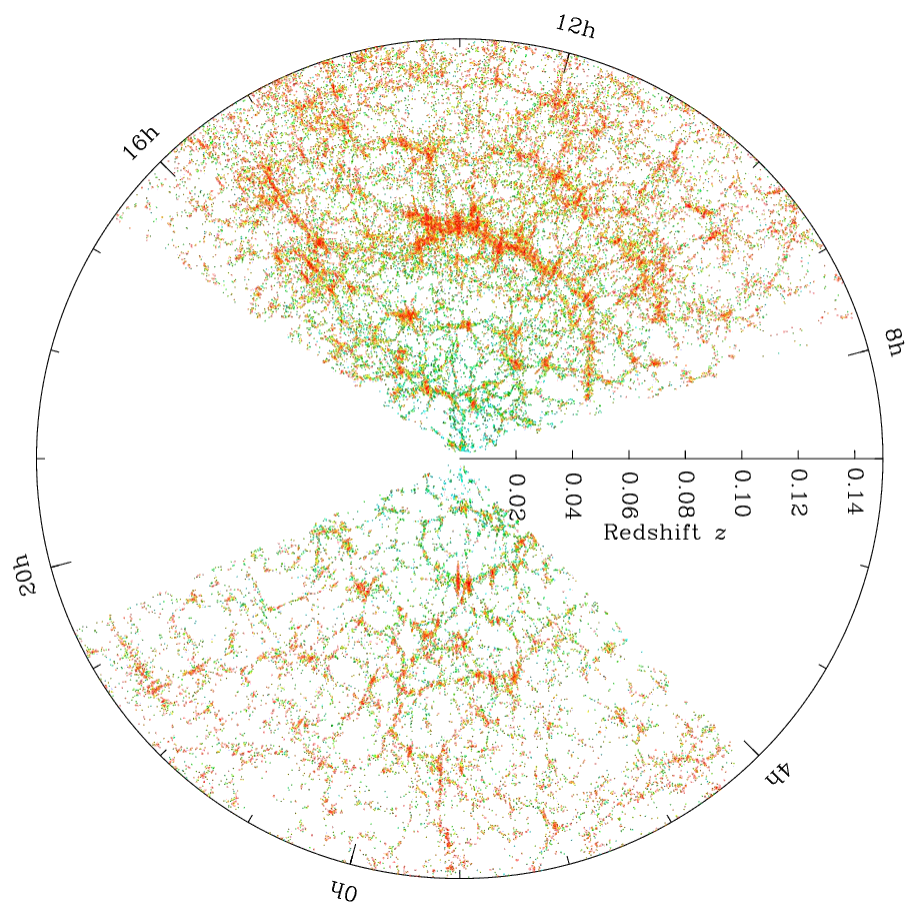Kernel Density Estimate

Much better fit

# KDE: Statistical Gold Standard

- Guaranteed to converge to the underlying distribution
- Provides normalized, true probability densities
- Few assumptions about shape of distribution: inferred from data

# KDE Usage



Galaxy Mass Distribution
[Sloan Digital Sky Survey]



Distribution of Bowhead Whales
[L.T. Quackenbush et al, *Arctic* 2010]

# KDE Definition



Each point in dataset contributes a *kernel*

Kernel: localized Gaussian "bump"

Kernels summed up to form estimate

Mixture of N Gaussians: N is the dataset size

$$f(x) = \frac{1}{n} \sum_{x_i \in \text{Data}} K(x - x_i)$$

# Problem: KDE does not scale

$$f(x) = \frac{1}{n} \sum_{x_i \in \text{Data}} K(x - x_i)$$
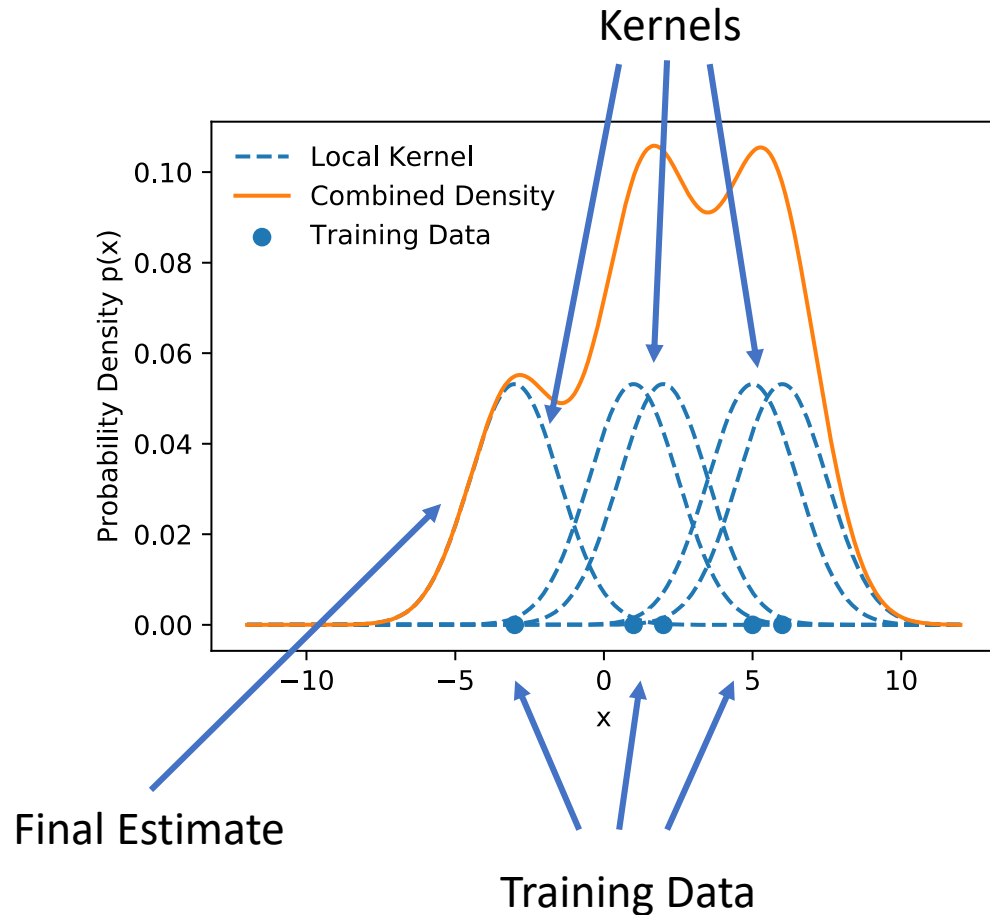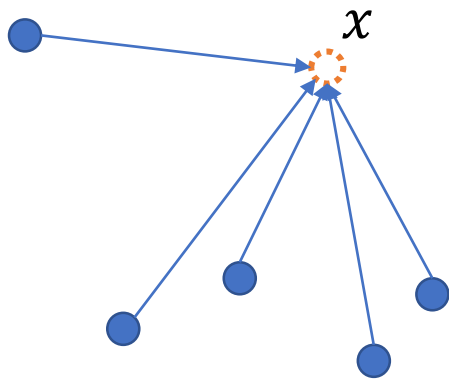
$O(n)$ to compute single density $f(x)$

$O(n^2)$ to compute all densities in data

2 hours to compute on 1M points
on 2.9Ghz Core i5
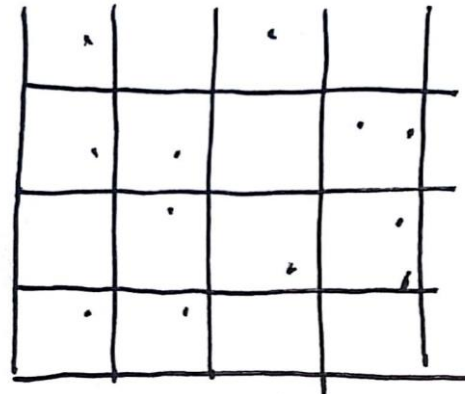
$x$

Training Data

How can we speed this up?

# Strawman Optimization: Histograms



Training Dataset      Binned Counts      Grid computation

Benefit: Runtime depends on grid size rather than N

Problem: Bin explosion in high dimensions

[Wand, *J. of Computational and Graphical Statics* 1994]

# Stepping Back: What users need

Anomaly Explanation

```
SELECT flight_mode FROM shuttle_sensors
WHERE kde(flow,speed) < threshold
```



Hypothesis Testing

```
SELECT color FROM galaxies
WHERE kde(x,y,z) < threshold
```

# From Estimation to Classification

```
SELECT flight_mode FROM shuttle_sensors
WHERE kde(flow,speed) < Threshold
```



Training Data



KDE Model



Classification

# End to End Query

Training Data
$$\mathbb{R}^{d \times n}$$

Kernel Density Estimation

$x$

Densities
$$\mathbb{R}^{n}$$

Threshold Filter

$$\begin{cases} \text{High } \textit{if } f(x) \geq t \\ \text{Low } \textit{if } f(x) < t \end{cases}$$

Classification
$$\{\text{Hi,Low}\}^{n}$$

Unnecessary for final output

# End to End Query

Kernel Density Estimation

Training Data

$$\mathbb{R}^{d \times n}$$

$x$

Threshold Filter

$$\begin{cases} \text{High } if \ f(x) \geq t \\ \text{Low } if \ f(x) < t \end{cases}$$

Classification

$$\{\text{Hi,Low}\}^{n}$$

# Recap

- KDE can model complex distributions

- Problem: KDE scales quadratically with dataset size

- Real Usage: KDE + Predicates = Kernel Density Classification

- Idea: Apply Predicate Pushdown to KDE

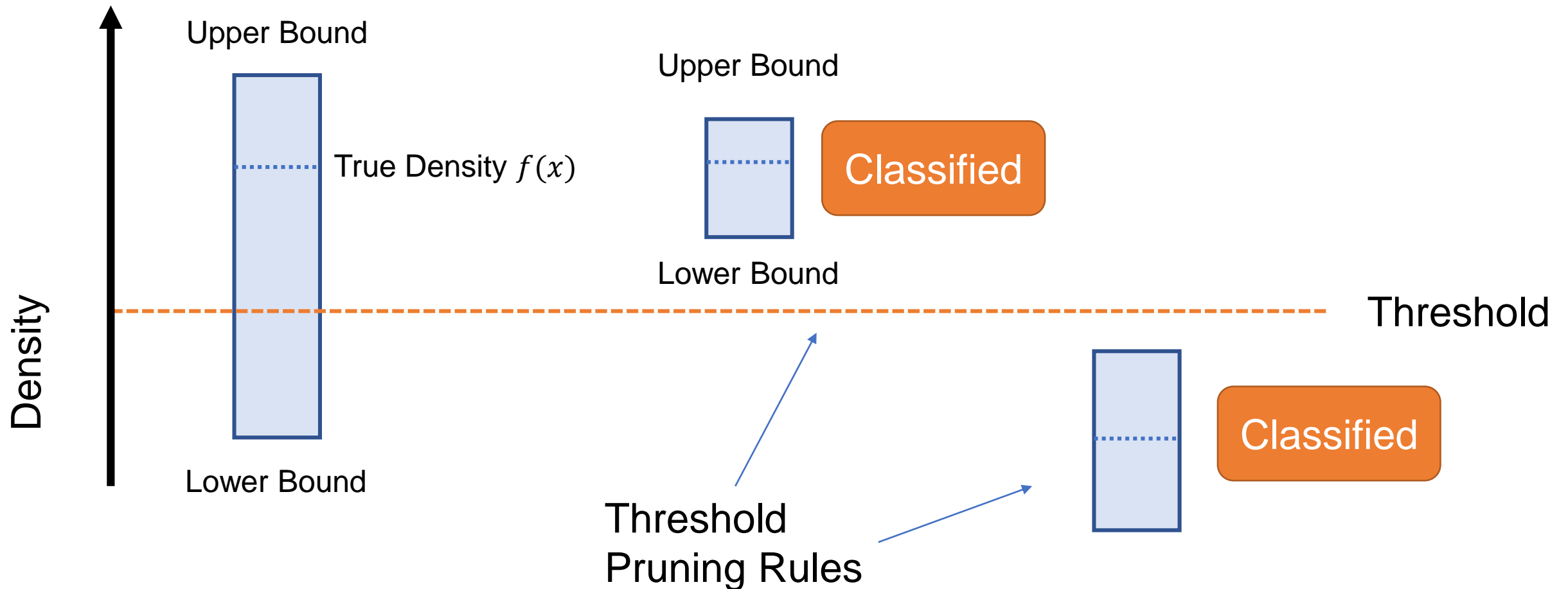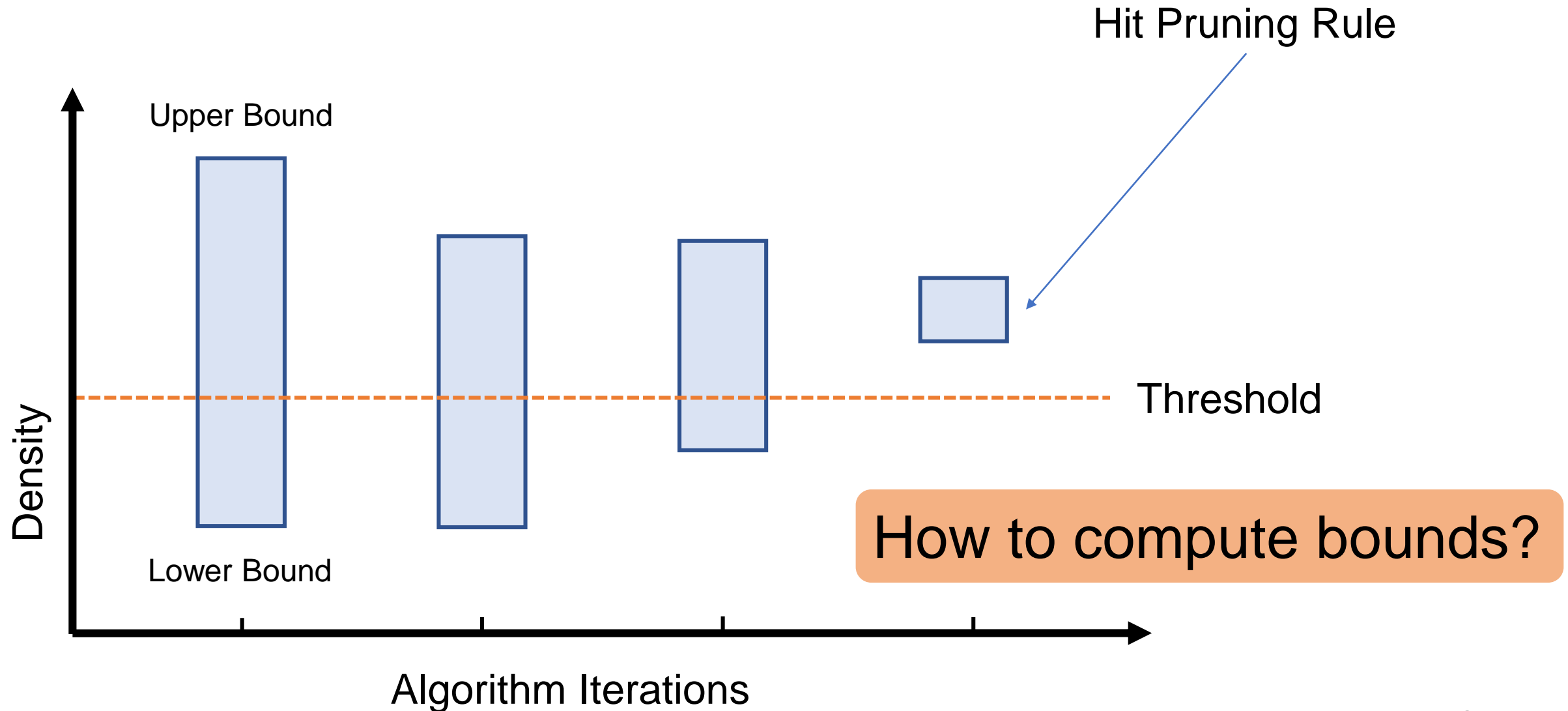# tkdc Algorithm Overview

1. Pick a threshold

2. Repeat: Calculate bounds on point density
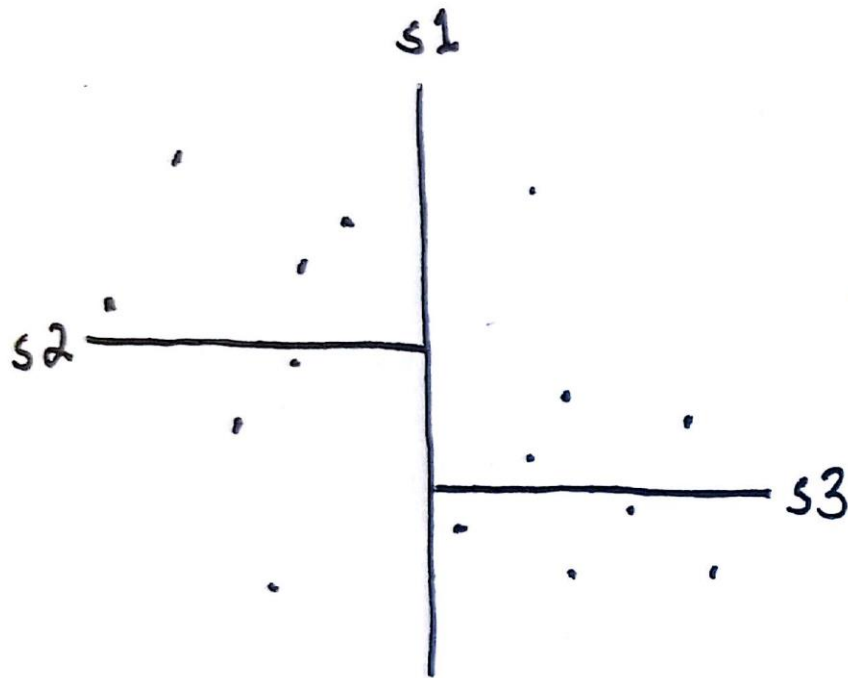
3. Stop when we can make a classification

# Classifying the density based on bounds
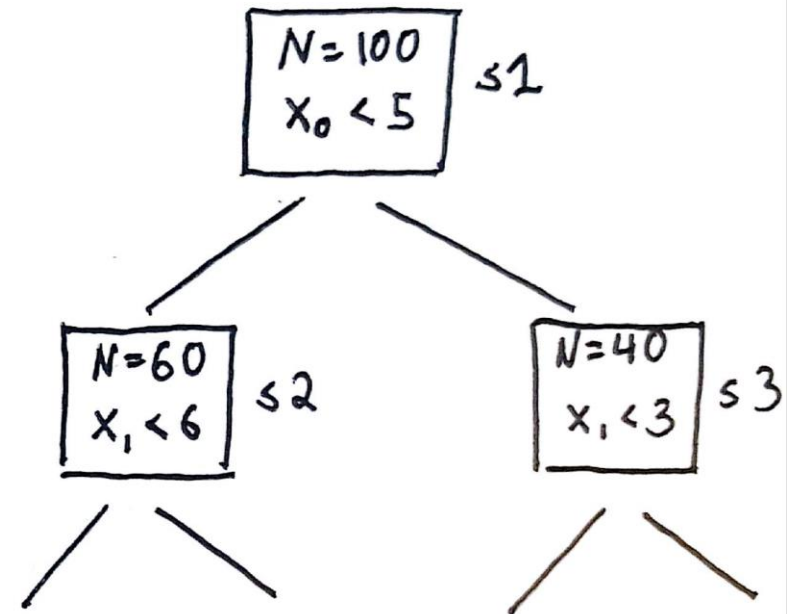
# Iterative Refinement

# k-d tree Spatial Indices
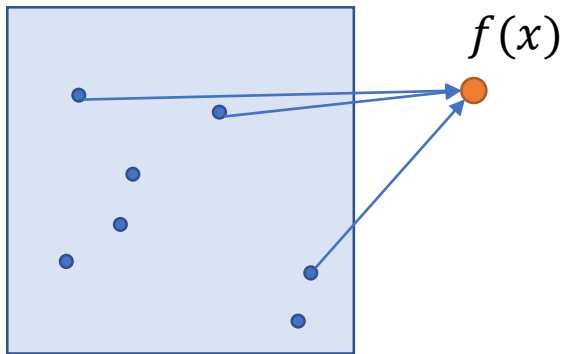


Divide N-dimensional space
1 axis at a time

Nodes for each Region
Track # of points + bounding box

# Bounding the densities
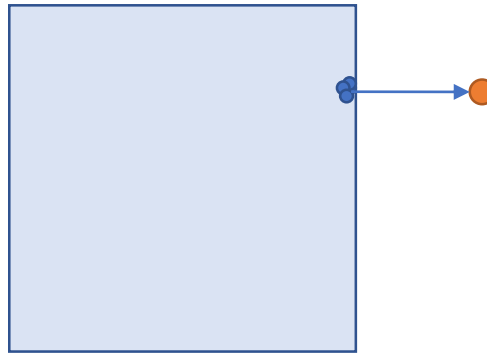
Given from k-d tree: Bounding Box, # Points Contained

Total contribution from a region can be bounded



Total Contribution | Maximum Contribution | Minimum Contribution

[Gray & Moore, *ICDM* 2003]

# Iterative Refinement

Initial Estimate       Step 1       Step 2       Step 3



$f(x)$

k-d tree root node

split

split

split

**Priority Queue**: Split nodes with largest uncertainty first

# tkdc Algorithm Overview

1. Pick a threshold
   - User-Specified
   - Automatically Inferred

2. Calculate bounds on a density
   - k-d tree bounding boxes

3. Refine the bounds until we can classify
   - Priority-queue guided region splitting

# Automatic Threshold Selection

- Probability Densities hard to work with:
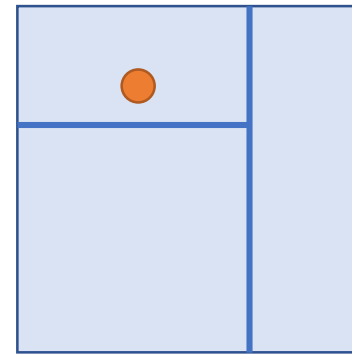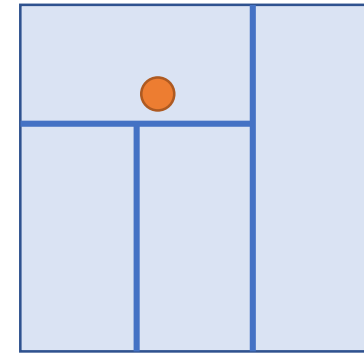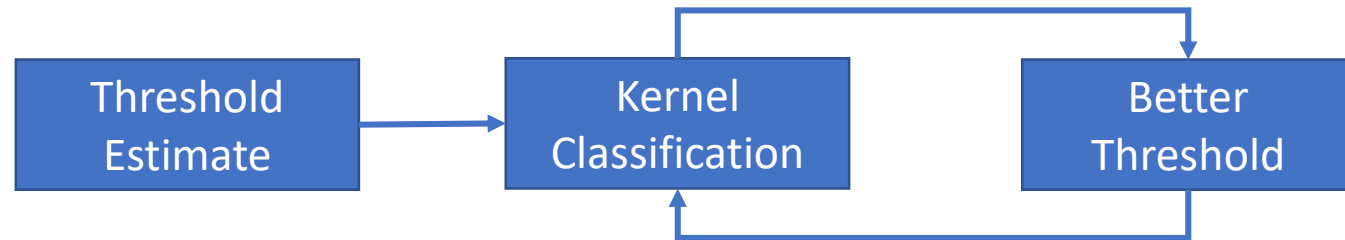  - Unpredictable
  - Huge range of magnitudes
- Good Default: capture a set % of the data

```
SELECT Quantile(kde(A,B), 1%) from shuttle_sensors
```

- Bootstrapping

| Threshold Estimate | → | Kernel Classification | ⇄ | Better Threshold |

- Classification for computing thresholds
  - See paper for details

# tkdc Complete Algorithm

- Pick a threshold
  - Inferred given desired % level


- Calculate bounds on a density
  - k-d tree bounding boxes


- Refine the bounds until we can make classification
  - Priority-queue guided region splitting

# Theorem: Expected Runtime

$n$ number of training points

$d$ dimensionality of data

$$\text{Runtime} = O\left(n^{\frac{d-1}{d}}\right)$$

$$\text{Naive} = O\left(n\right)$$

100 million data points, 2-dimensions $\dfrac{100M}{(100M)^{\frac{1}{2}}} \approx 10{,}000\text{x}$

100 million data points, 8-dimensions $\dfrac{100M}{(100M)^{\frac{7}{8}}} \approx 10\text{x}$

# Runtime in practice: Experimental Setup
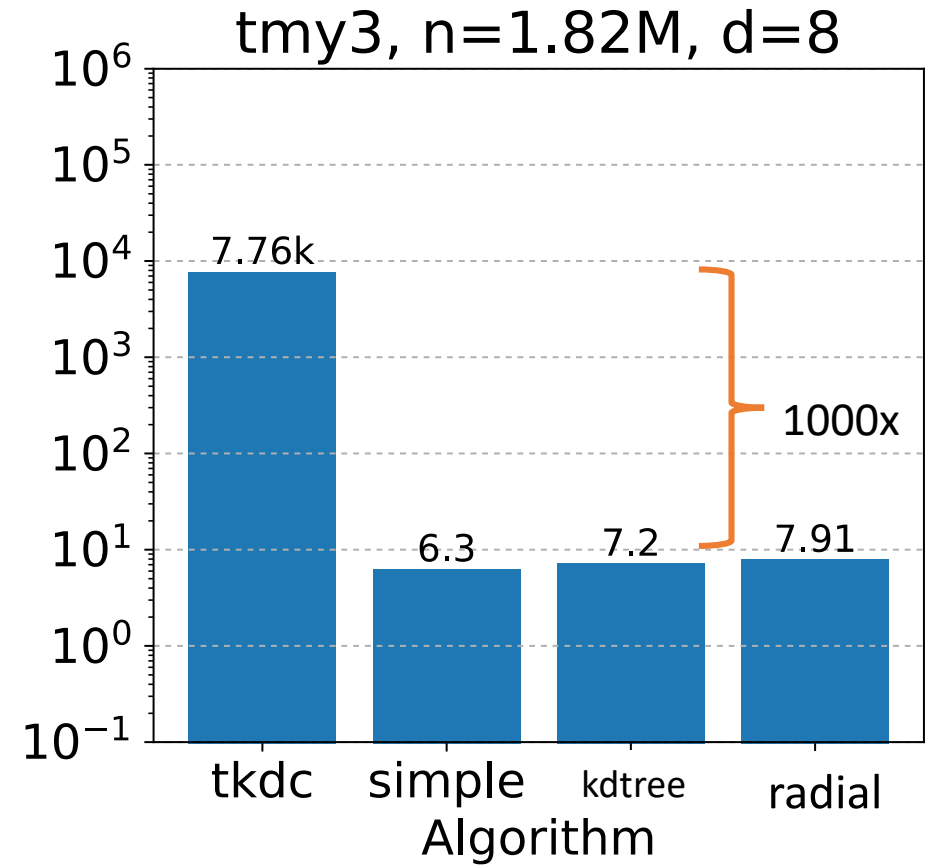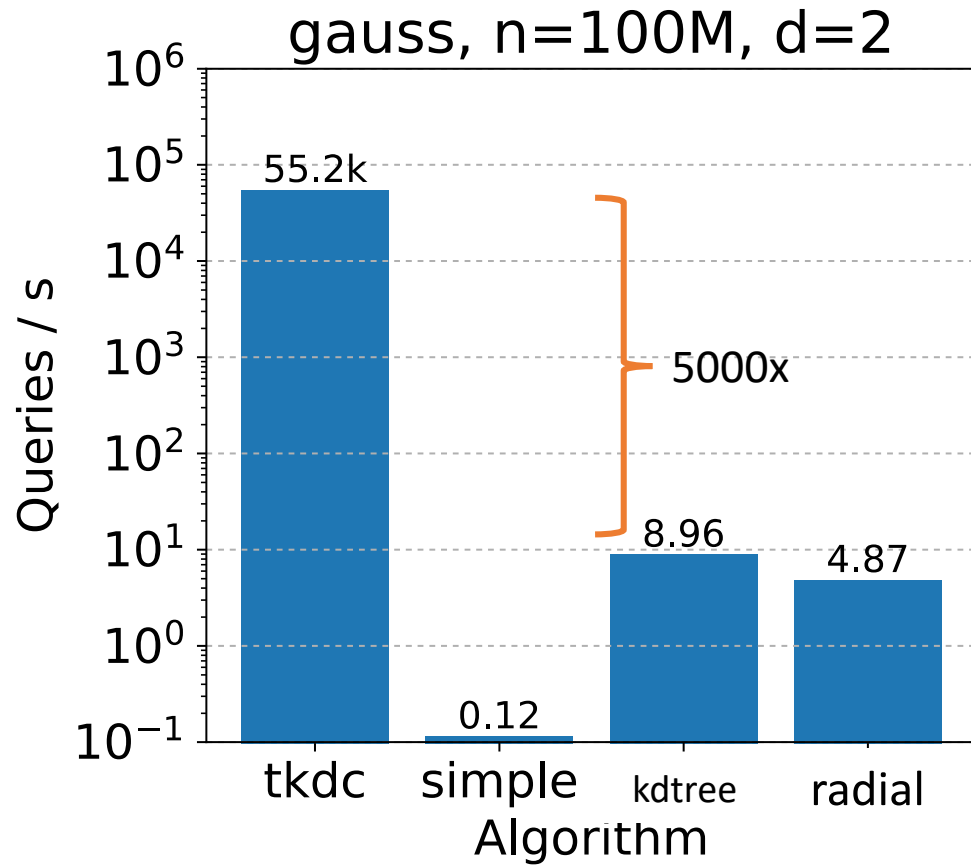
Single Threaded, In-memory

Total Time = Training Time + Threshold Estimation + Classify All

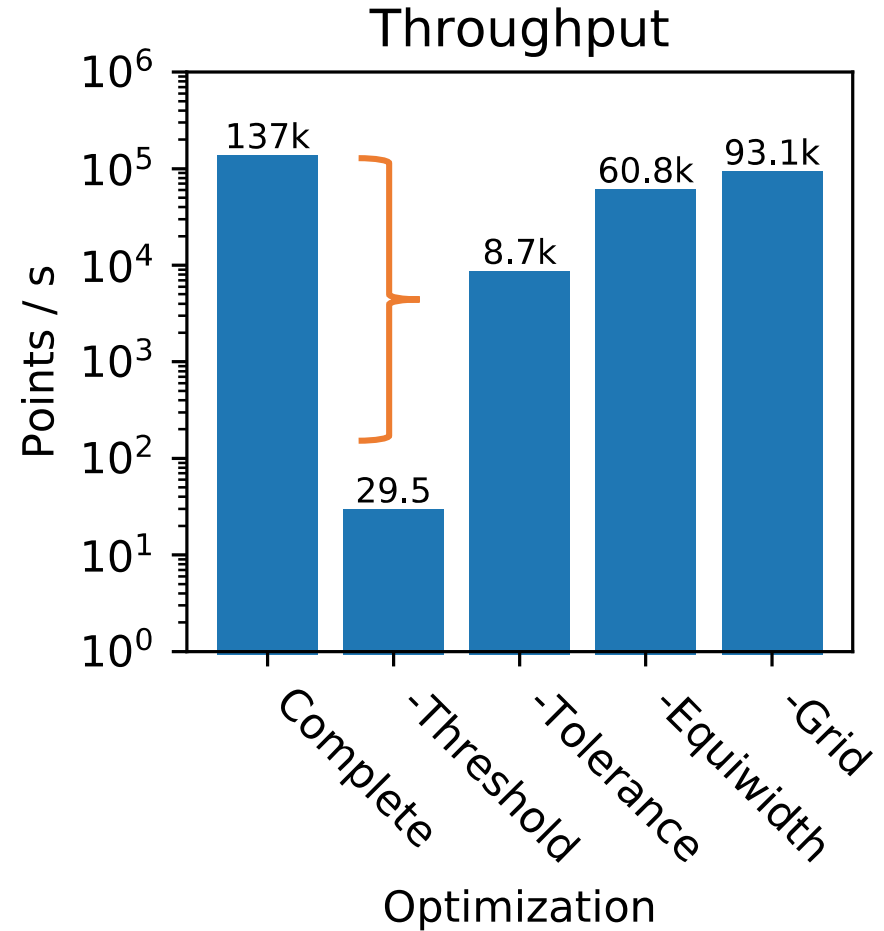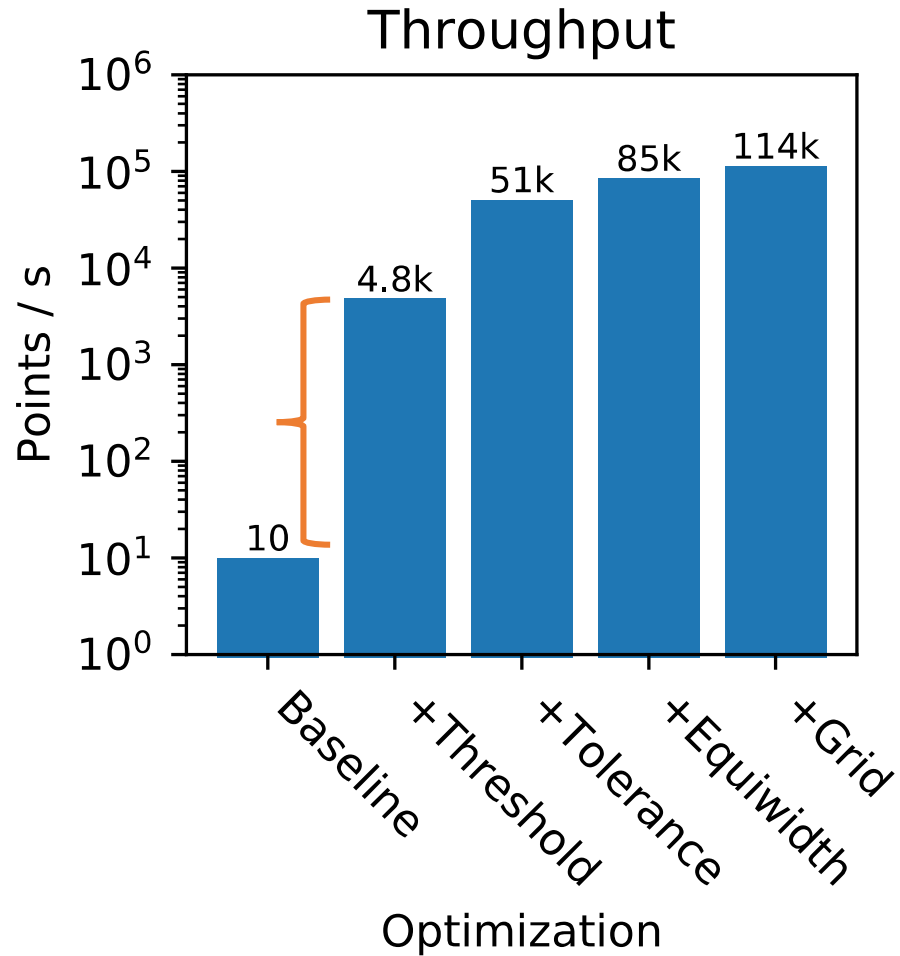Threshold = 1% classification rate

Baselines:
- simple:     naïve for loop over all points
- kdtree:     k-d tree approximate density estimation, no threshold
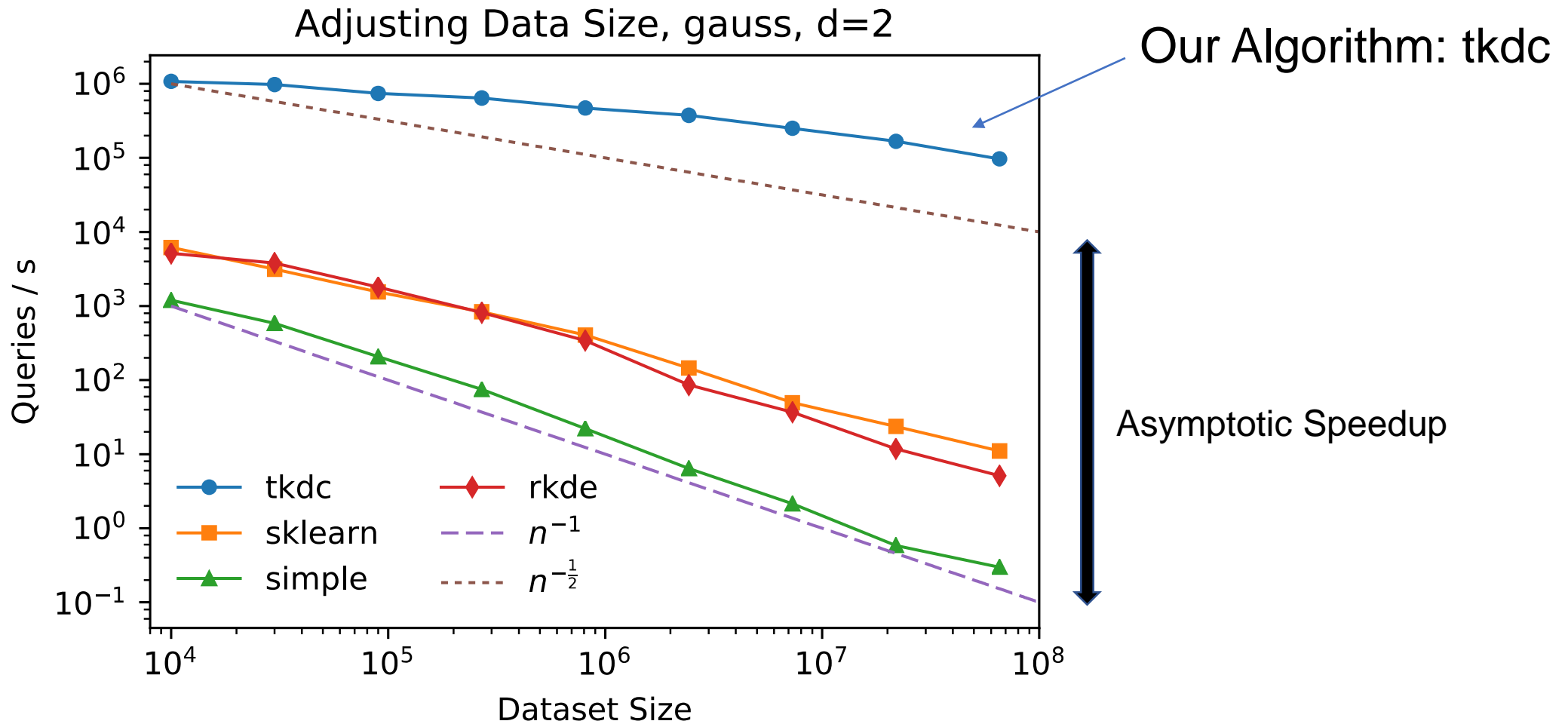- radial:     iterates through points, pruning > certain radius

# KDE Performance Improvement

# Threshold Pruning Contribution

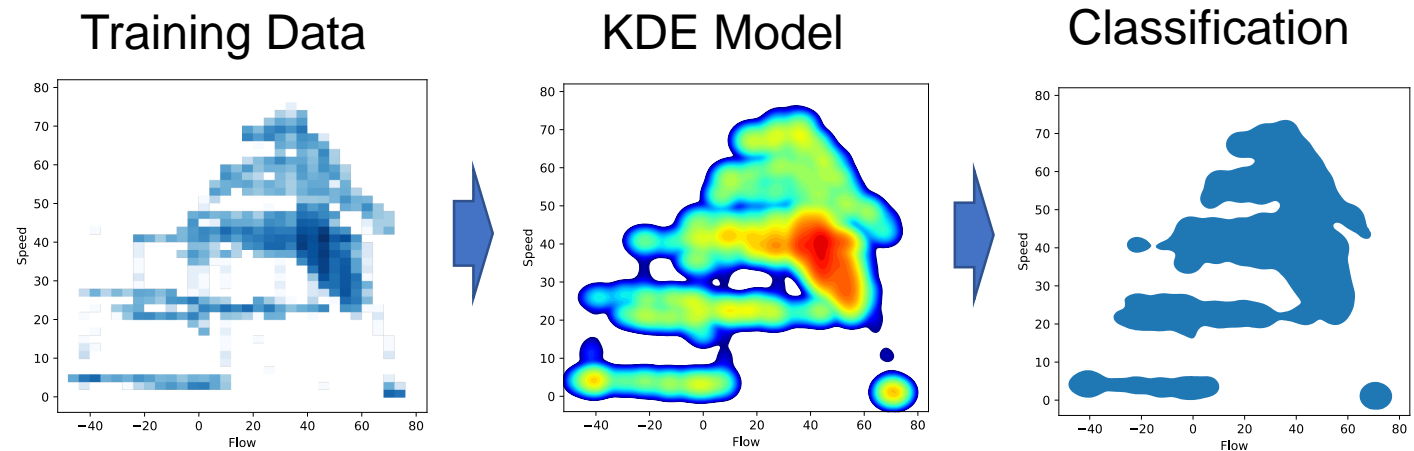# tkdc scales well with dataset size



Adjusting Data Size, gauss, d=2

Our Algorithm: tkdc

Asymptotic Speedup

Legend: tkdc, sklearn, simple, rkde, $n^{-1}$, $n^{-\frac{1}{2}}$

Y-axis: Queries / s

X-axis: Dataset Size

# Conclusion

KDE:
Powerful & Expensive

Real Queries:
MacroBase

Systems Techniques:

https://github.com/stanford-futuredata/tKDC

```
SELECT flight_mode FROM shuttle_sensors
WHERE kde(flow, speed) < Threshold
```



Training Data          KDE Model          Classification

Predicate Pushdown, k-d tree indices:

1000x, Asymptotic Speedups