# Finding Probabilistic Nearest Neighbors for Query Objects with Imprecise Locations
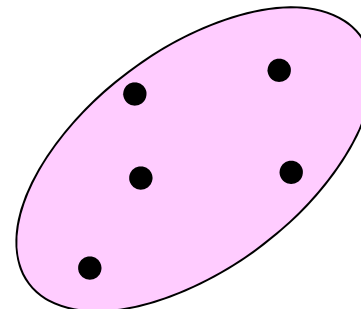
Yuichi Iijima and Yoshiharu Ishikawa
Nagoya University, Japan

# Outline

- <span style="color:red">Background and Problem Formulation</span>
- Related Work
- Query Processing Strategies
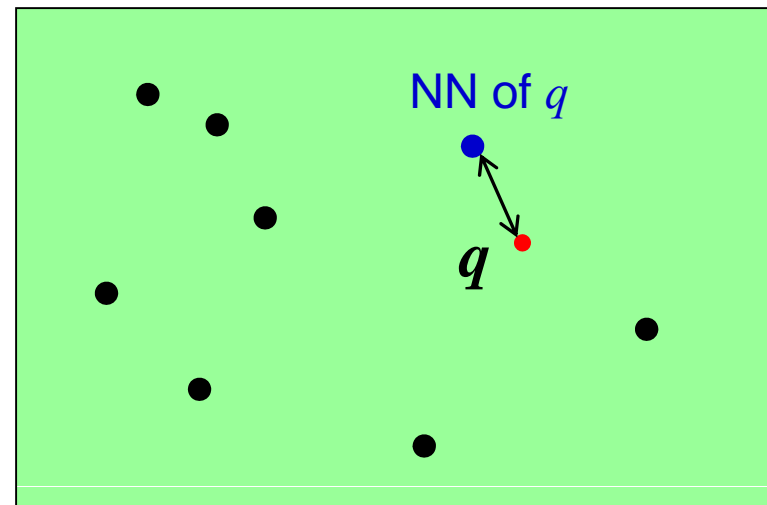- Experimental Results
- Conclusions

# Imprecise Location Information

- ## Sensor Environments
  - Measurement Noise
  - Frequent updates may not be possible
    - GPS-based positioning consumes batteries

- ## Robotics
  - Localization using sensing and movement histories
  - Probabilistic approach has vagueness
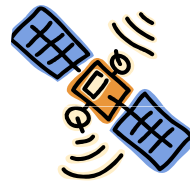
- ## Privacy
  - Location Anonymity

# Nearest Neighbor Queries

- ## Nearest Neighbor Queries
  - Example: Find the closest bus stop
  - Traditional problem in spatial databases
    - Efficient query processing using spatial indices
    - Extensible to multi-dimensional cases (e.g., image retrieval)

- What happen if the location of query object is uncertain?



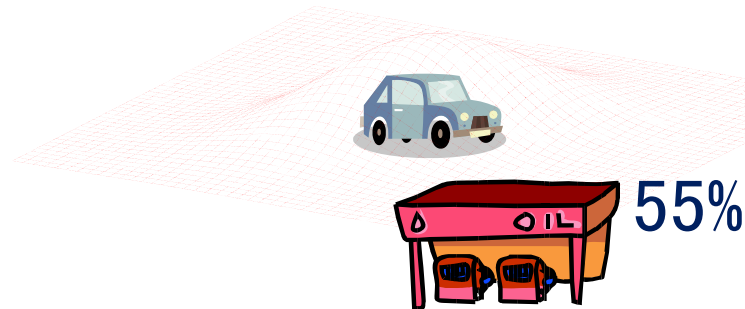NN of $q$

$q$

# Example Scenario (1)

- Query: Find the nearest gas station

0%

35%

10%

55%

Location estimation based on noisy GPS data

Nearest gas station depends on the possible car locations
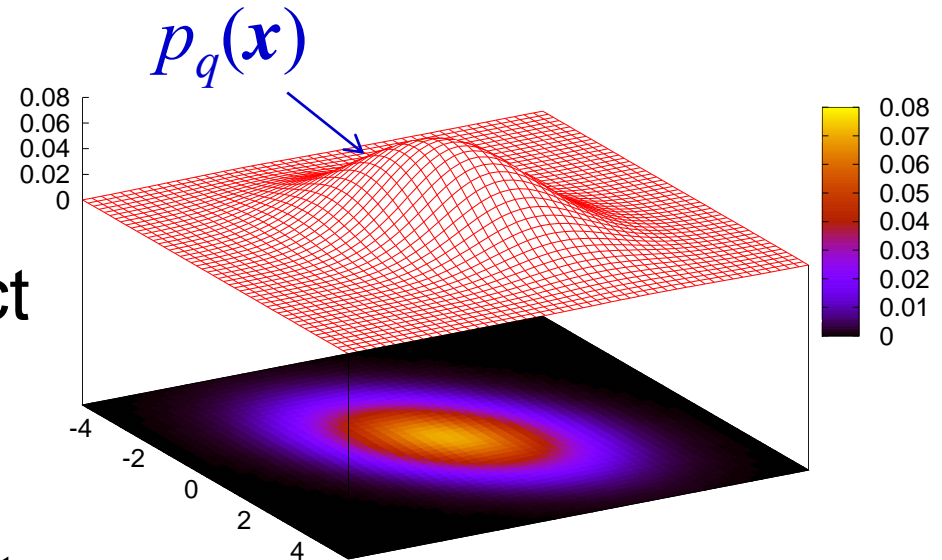
NN objects are defined in a probabilistic way

# Example Scenario (2)

- **Mobile Robotics**
  - Location of the robot is estimated based on movement histories and sensor data
  - Measurements are noisy
  - Localization based on <span style="color:red">probabilistic modeling</span>
    - Kalman filter, particle filter, etc.
  - Estimated location is given as a <span style="color:blue">probabilistic density function (PDF)</span>
    - PDF changes on each estimation

- PNNQ for short

- Assumptions

  - Location of query object $q$ is specified as a <span style="color:red">Gaussian distribution</span>

  - Target data: static points

- Gaussian Distribution

$p_q(x)$

$$p_q(x) = \frac{1}{(2\pi)^{d/2}\,|\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2}(x-q)^t\,\boldsymbol{\Sigma}^{-1}(x-q)\right]$$

  - $\boldsymbol{\Sigma}$: Covariance matrix

- Definition

$$\mathrm{Pr}_{NN}(q, o) = \mathrm{Pr}(\forall o' \in O,\ o' \neq o,\ \|x - o\| \leq \|x - o'\|)$$

$$PNNQ(q, \theta) = \{o \mid o \in O,\ \mathrm{Pr}_{NN}(q, o) \geq \theta\}$$

- Find objects which satisfy the condition

  – The probability that the object is the nearest neighbor of $q$ is greater than or equal to $\theta$

7

# Outline

- Background and Problem Formulation
- Related Work
- Query Processing Strategies
- Experimental Results
- Conclusions

# Related Work

- Query processing methods for uncertain (location) data
  - Cheng, Prabhakar, et al. (SIGMOD'03, VLDB'04, …)
  - Tao et al. (VLDB'05, TODS'07)
  - Consider arbitrary PDFs or uniform PDFs
  - Most of the case, target objects are imprecise
- Research related to Gaussian distribution
  - Gauss-tree [Böhm et al., ICDE'06]
    - Target objects are based on Gaussian distributions
- Our former work
  - Ishikawa, Iijima, Yu (ICDE'09): Probabilistic range queries

# Outline

- Background and Problem Formulation
- Related Work
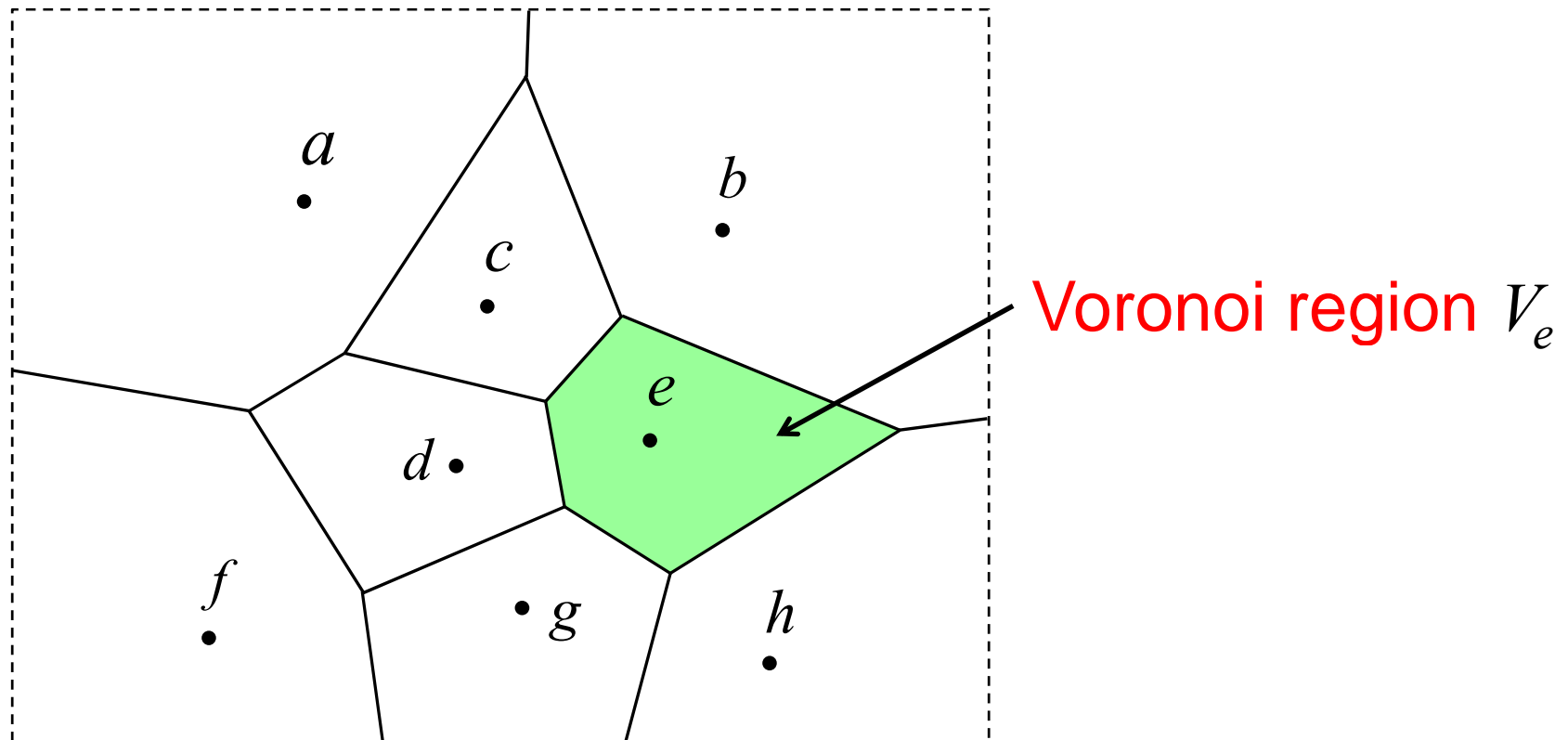- Query Processing Strategies
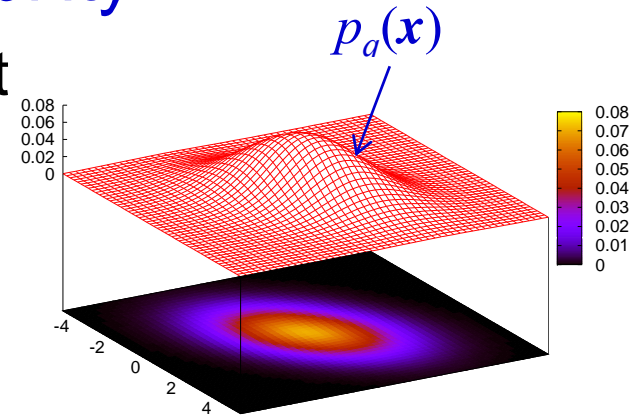- Experimental Results
- Conclusions

# Naïve Approach (1)

- Use of Voronoi Diagram
  - Well-known method for standard (non-imprecise) nearest neighbor queries



Voronoi region $V_e$

# Naïve Approach (2)

- $\text{Pr}_{NN}(q, o)$: Nearest neighbor probability
  - Probability that object $o$ is the nearest neighbor of query object $q$
  - It can be computed by integrating the probability density function $p_q(\boldsymbol{x})$ over Voronoi region $V_o$



$p_a(\boldsymbol{x})$

$$\text{Pr}_{NN}(q, o) = \int_{V_o} p_q(\boldsymbol{x}) d\boldsymbol{x}$$

- Problem
  - Need to consider all target objects
  - Numerical integration (Monte Carlo method) is quite costly!



12

# Our Approach

- **Outline of processing**
  1. **Filtering**
     - Prune non-candidate objects whose $\text{Pr}_{NN}$ are obviously smaller than the threshold $\theta$
     - Low-cost filtering conditions
  2. **Numerical integration** for the remaining candidate objects

- **Two strategies**
  - $\theta$-region-based Approach
  - SES-based Approach
    - SES: Smallest Enclosing Sphere

- $\theta$-region

  - Similar concepts are often used in query processing for uncertain spatial databases

  - Definition: Ellipsoidal region for which the result of the integration becomes $1 - 2\theta$:
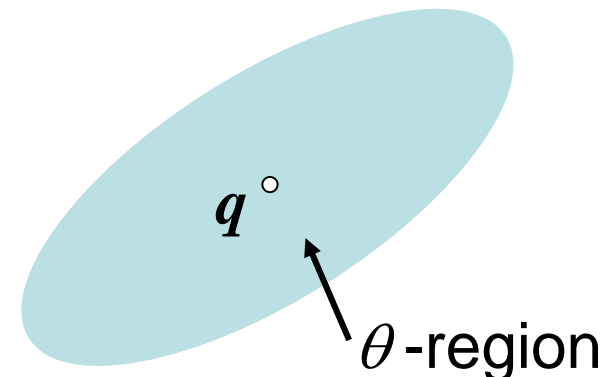
$$\int_{(x-q)^t \Sigma^{-1}(x-q) \leq r_\theta^2} p_q(x)dx = 1 - 2\theta$$

The ellipsoidal region

$$(x - q)^t \Sigma^{-1} (x - q) \leq r_\theta^2$$

is the $\theta$-region

- Example: $\theta$ is specified as 1%, we consider 98% $\theta$-region

$q^{\circ}$

$\theta$-region

- **Query Processing**
  1. *$\theta$-region* for the query is computed at first
     - $\theta$-region can be derived using $r_\theta$-table:
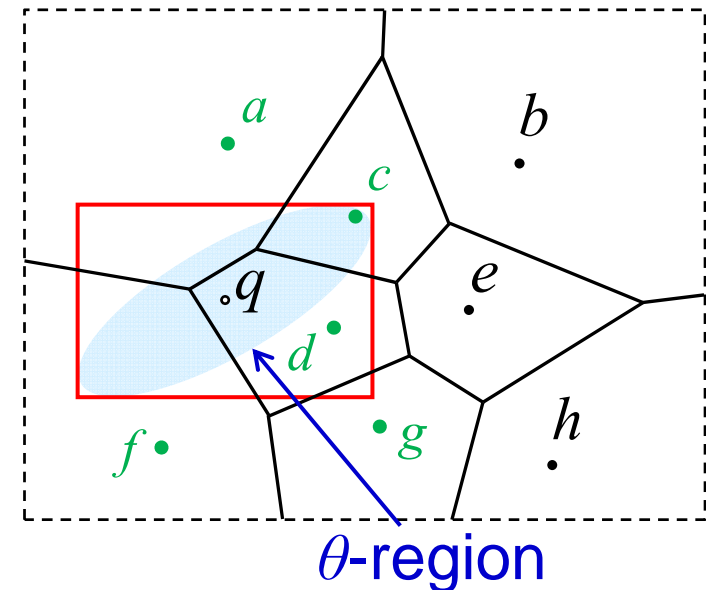       – It is constructed for the normal Gaussian ($\Sigma = \mathbf{I}, q = \mathbf{0}$): Given $\theta$, it returns appropriate $r_\theta$
     - Final $\theta$-region can be obtained by transformation
  2. Derive the bounding box of $\theta$-region
  3. Objects whose Voronoi regions overlaps with the box are the candidates
     - $\{a, c, d, f, g\}$, in this example



$\theta$-region

- Derivation details
- First, we consider the normal Gaussian
  - Using $r_\theta$-table, we can get the appropriate $r_\theta$ for given $\theta$
- Second, a spherical $\theta$-region is derived based on transformation
  - Transformation is performed by analyzing $\Sigma$

- Third, the bonding box is calculated

$$w_i = \sigma_i r_\theta$$

$$\sigma_i = \sqrt{(\Sigma)_{ii}}$$

where $(\Sigma)_{ii}$ is the $(i, i)$ entry of $\Sigma$

- **SES: Smallest Enclosing Sphere**
  - For each Voronoi region $V_o$, we compute its SES $SES_o$ beforehand



$SES_e$: SES for Voronoi region of $e$

- Integration over $SES_o$ gives the upper-bound for $\mathrm{Pr}_{NN}(q, o)$

$$\mathrm{Pr}_{NN}(q,o) = \int_{V_o} p_q(\boldsymbol{x})d\boldsymbol{x} < \int_{SES_o} p_q(\boldsymbol{x})d\boldsymbol{x}$$

  – Integration over a sphere region is more easier to compute

    - We use a table called U-catalog constructed beforehand

- What is U-catalog?
  - Given two parameters $\alpha$ and $\delta$, it returns corresponding integral
  - U-catalog is made for different $(\alpha, \delta)$ pairs by computing the integral of normal Gaussian over sphere region $R$

$$\pi(\alpha, \delta) = \int_{\boldsymbol{x} \in R} p_{\text{norm}}(\boldsymbol{x})d\boldsymbol{x}$$

| $\alpha$ | $\delta$ | $\pi(\alpha, \delta)$ |
|---|---|---|
| 0.0 | 0.1 | … |
| 0.0 | 0.2 | … |
| … | … | … |
| 1.0 | 0.1 | |
| … | … | … |

- To use U-catalog, another approximation is required since it is only useful for normal Gaussian
  - Use of upper-bounding function $p_q^\top(x)$
  - $p_q^\top(x)$ tightly bounds $p_q(x)$ and has spherical isosurfaces
  - For $p_q^\top(x)$, we can easily derive its integral over $SES_o$ using U-catalog

- In summary, we use two approximations:

$$
\begin{aligned}
\mathrm{Pr}_{NN}(q,o) \quad &= \quad \int_{V_o} p_q(x)\,dx \\
&< \quad \int_{SES_o} p_q(x)\,dx \\
&\leq \quad \int_{SES_o} p_q^\top(x)\,dx
\end{aligned}
$$

# Strategy 2: SES-Based Approach (5)

- Bounding Function
  - Original Gaussian PDF

$$p_q(\boldsymbol{x}) = \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{q})^t \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{q})\right]$$

  - Upper-Bounding Function

$$p_q^{\mathsf{T}}(\boldsymbol{x}) = \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{\lambda^{\mathsf{T}}}{2}\|\boldsymbol{x}-\boldsymbol{q}\|^2\right]$$

where

$$\boldsymbol{\Sigma}^{-1} = \sum_{i=1}^{d} \lambda_i \mathbf{v}_i \mathbf{v}_i^t$$

$$\lambda^{\mathsf{T}} = \min\{\lambda_i\}$$

$$p_q(\boldsymbol{x}) \le p_q^{\mathsf{T}}(\boldsymbol{x})$$

is satisfied

# Outline

- Background and Problem Formulation
- Related Work
- Query Processing Strategies
- <span style="color:red">Experimental Results</span>
- and Conclusions

- Target Data
  - 2D point data (50K entries)
  - Based on road line segments in Long Beach
- Computation of Voronoi regions and SESs
  - LEDA package was used
- Comparison
  - Strategies 1, 2, and their hybrid approach
  - Evaluation metric: Response time

- **Default Parameters**
  - Covariance matrix

$$\boldsymbol{\Sigma} = \gamma \begin{bmatrix} 7 & 2\sqrt{3} \\ 2\sqrt{3} & 7 \end{bmatrix}$$

  

  - For this matrix the shape of the isosurface of $p_q(\mathbf{x})$ is an ellipse titled at 30 degrees and the major-to-minor axis ratio is 3:1
  - $\gamma$ : Parameter for controlling vagueness (default: $\gamma = 10$)

  - Probability threshold value: $\theta = 0.01$

  - No. of samples for Monte Carlo method: 1,000,000

# Candidate Objects in Strategy 1



Bounding box of the $\theta$-region

Query center

Voronoi regions for candidate objects

Voronoi regions for answer objects

# Candidate Objects in Strategy 2



Query center

Voronoi regions for candidate objects

Voronoi regions for answer objects

# Candidate Objects in Hybrid Strategy



Bounding box of the $\theta$-region

Query Center

Voronoi regions for candidate objects

Voronoi regions for answer objects

Filtering   Compute Prob.   Rest

| | Strategy 1 | Strategy 2 | Hybrid |
| Response Time[s] | | | |

Strategy 1: Rest 2.45, Compute Prob. 43.70, Filtering 0.09
Strategy 2: Rest 2.48, Compute Prob. 38.70, Filtering 0.19
Hybrid: Rest 2.47, Compute Prob. 34.13, Filtering 0.13

| | Strategy 1 | Strategy 2 | Hybrid |
| --- | --- | --- | --- |
| Number of candidates | 179 | 150 | 129 |
| Number of answers | 26 | | |

- Most of the processing time is spent in numerical integration

- A strategy that can prune more objects has better performance

# Experimental Results - Different $\gamma$ Values



$\gamma = 1$
(almost exact)

$\gamma = 10$

$\gamma = 50$
(too vague)

Legend: ■ Filtering ■ Compute Prob. ■ Rest

**$\gamma = 1$ (almost exact)**

Response Time[s]

| | Strategy 1 | Strategy 2 | Hybrid |
|---|---|---|---|
| Rest | 2.45 | 2.48 | 2.48 |
| Compute Prob. | 6.28 | 6.05 | 5.53 |
| Filtering | 0.09 | 0.13 | 0.13 |

| | | | |
|---|---|---|---|
| Number of candidates | 24 | 31 | 23 |
| Number of answers | 8 | | |

**$\gamma = 10$**

| | Strategy 1 | Strategy 2 | Hybrid |
|---|---|---|---|
| Rest | 2.45 | 2.48 | 2.47 |
| Compute Prob. | 43.70 | 38.70 | 34.13 |
| Filtering | 0.09 | 0.19 | 0.13 |

| | | | |
|---|---|---|---|
| Number of candidates | 179 | 150 | 129 |
| Number of answers | 26 | | |

**$\gamma = 50$ (too vague)**

| | Strategy 1 | Strategy 2 | Hybrid |
|---|---|---|---|
| Rest | 2.45 | 2.47 | 2.46 |
| Compute Prob. | 209.22 | 70.88 | 66.41 |
| Filtering | 0.09 | 0.22 | 0.13 |

| | | | |
|---|---|---|---|
| Number of candidates | 847 | 276 | 260 |
| Number of answers | 15 | | |

Query is costly when impreciseness is high

29

# Experimental Results - Different $\theta$ Values

$\theta = 0.01$

$\theta = 0.03$

$\theta = 0.05$

Filtering ■ Compute Prob. ■ Rest

$\theta = 0.01$

| Number of candidates | 179 | 150 | 129 |
| --- | --- | --- | --- |
| Number of answers | 26 | | |

$\theta = 0.03$

| | 128 | 76 | 68 |
| --- | --- | --- | --- |
| | 7 | | |

$\theta = 0.05$

| | 107 | 44 | 40 |
| --- | --- | --- | --- |
| | 3 | | |

# Experimental Results - Different Shapes



**Circle**  **Default Ellipse**  **Narrow Ellipse**

Legend: ■ Filtering ■ Compute Prob. ■ Rest

**Circle**

| | Strategy 1 | Strategy 2 | Hybrid |
|---|---|---|---|
| Rest | 2.45 | 2.46 | 2.45 |
| Compute Prob. | 26.97 | 13.60 | 13.42 |
| Filtering | 0.09 | 0.15 | 0.13 |

| Number of candidates | 115 | 50 | 49 |
|---|---|---|---|
| Number of answers | 26 | | |

**Default Ellipse**

| | Strategy 1 | Strategy 2 | Hybrid |
|---|---|---|---|
| Rest | 2.45 | 2.48 | 2.47 |
| Compute Prob. | 43.70 | 38.70 | 34.13 |
| Filtering | 0.09 | 0.19 | 0.13 |

| Number of candidates | 179 | 150 | 129 |
|---|---|---|---|
| Number of answers | 26 | | |

**Narrow Ellipse**

| | Strategy 1 | Strategy 2 | Hybrid |
|---|---|---|---|
| Rest | 2.45 | 2.46 | 2.46 |
| Compute Prob. | 69.69 | 84.56 | 62.56 |
| Filtering | 0.09 | 0.20 | 0.13 |

| Number of candidates | 276 | 366 | 250 |
|---|---|---|---|
| Number of answers | 24 | | |

Response Time[s]

Narrow ellipse (correlated distribution) needs to consider additional candidates

31

# Conclusions of Experimental Results

- Most of the processing time is spent in numerical integration

    ⇒A strategy that can prune more objects has better performance

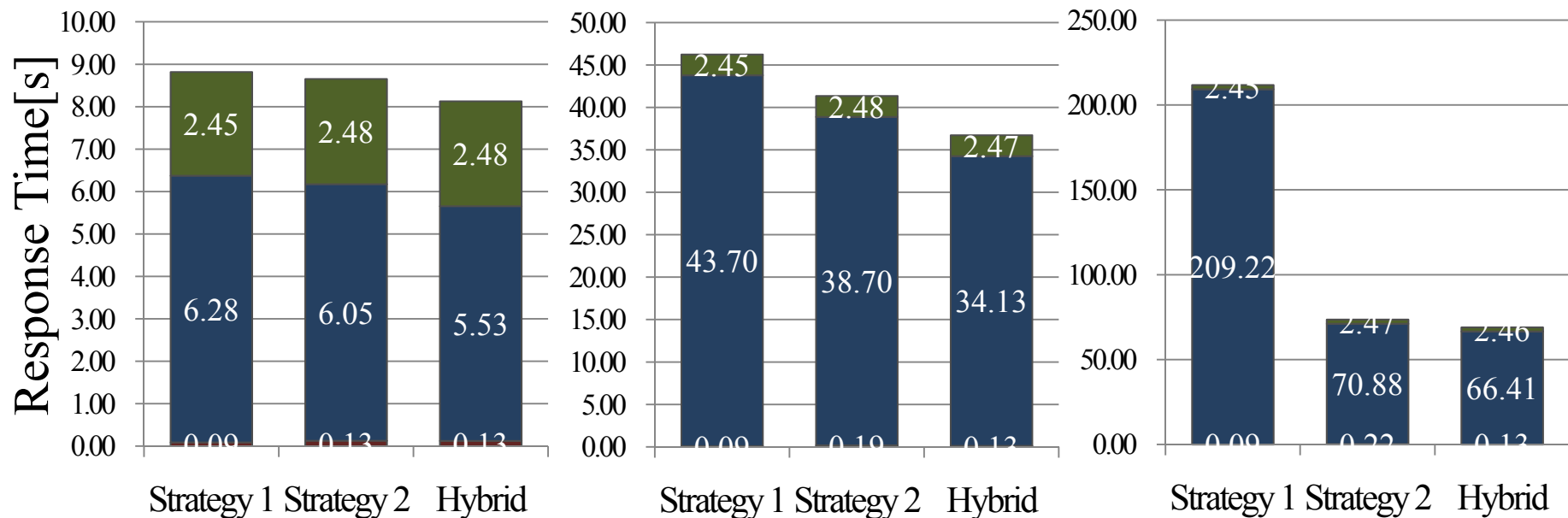- Superiority or inferiority of two strategies depends on the given query and the specified parameters

    - No apparent winner

- The hybrid strategy inherits the benefits of two strategies and shows the best performance

# Outline

- Background and Problem Formulation
- Related Work
- Query Processing Strategies
- Experimental Results
- Conclusions

# Conclusions

- **Nearest neighbor query processing methods for imprecise query objects**
  - Location of query object is represented by Gaussian distribution
  - Two strategies and their combination
  - Reduction of numerical integration is important
  - Proposal of two pruning strategies and their evaluation
- **Future work**
  - Evaluation for multi-dimensional cases ($d \geq 3$)

# Finding Probabilistic Nearest Neighbors for Query Objects with Imprecise Locations

Yuichi Iijima and Yoshiharu Ishikawa
Nagoya University, Japan