# Falconn++: A locality-sensitive filtering approach for approximate nearest neighbor search

Ninh Pham, Tao Liu
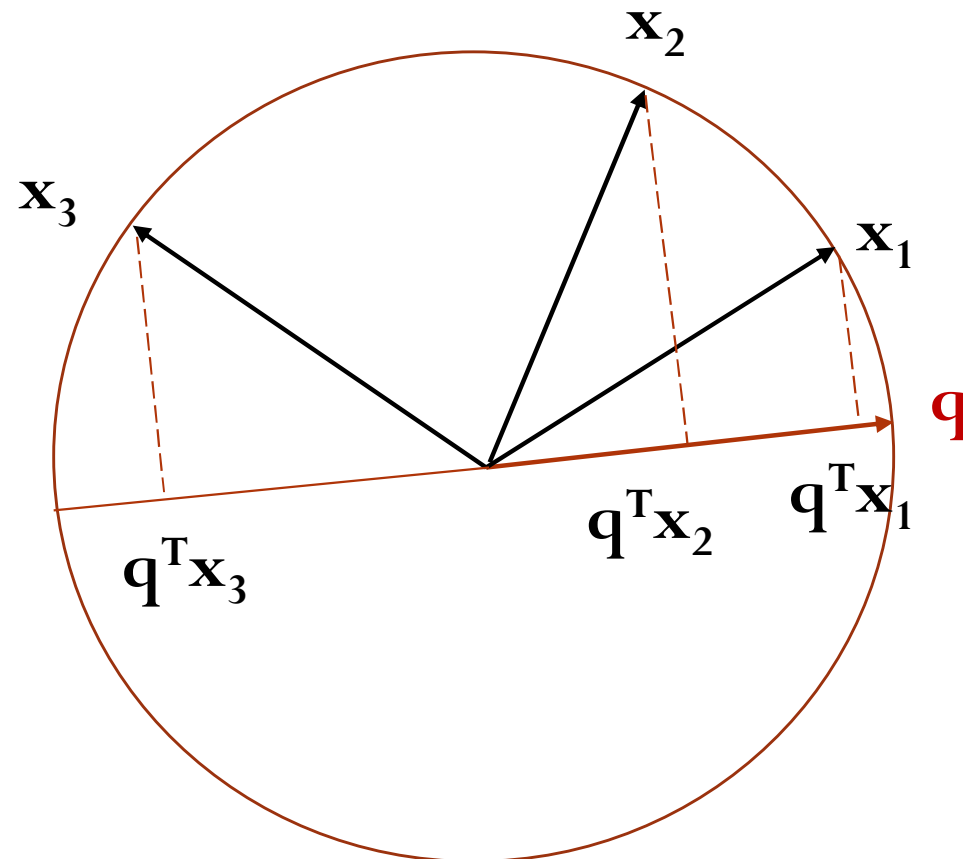
University of Auckland

MIT, Oct 15, 2022
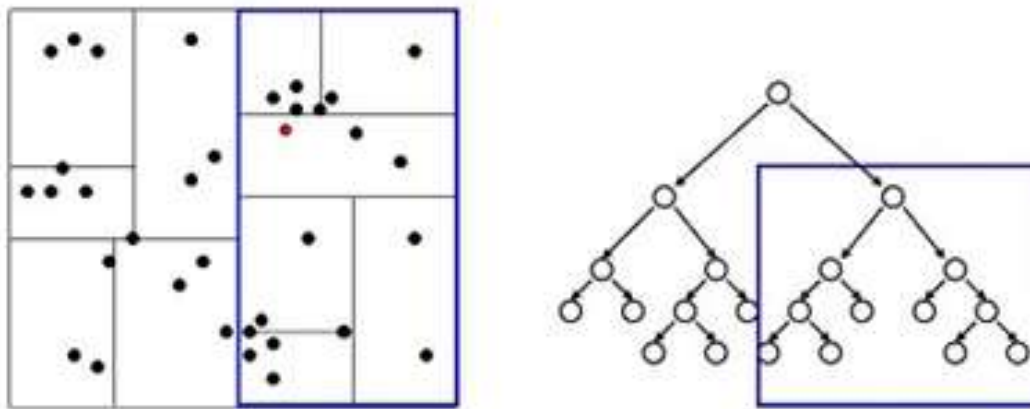
1

# Nearest neighbor search in sphere

- Nearest neighbor search (NNS) in a unit sphere:
  - Given a data set $\mathbf{X}$ of size $\mathbf{n}$ and a query $\mathbf{q}$ in $\mathbf{d}$ dimensional unit sphere, return the point $\mathbf{x} \in \mathbf{X}$ such that the $\mathbf{dist(q, x)} = \|\mathbf{x} - \mathbf{q}\|$ is minimum.



$\mathbf{x}_1$ is the top-1.

# Challenges of NNS

- Curse of dimensionality:
    - Given a polynomial indexing space $\mathbf{n^{O(1)}d^{O(1)}}$, existing a sublinear time algorithm to solve exact NNS refutes SETH [Wil18].
    - Indexing space or query time must be exponential in $\mathbf{d}$.

- Classic solutions: KD Tree



Space: $\mathbf{O(n)}$
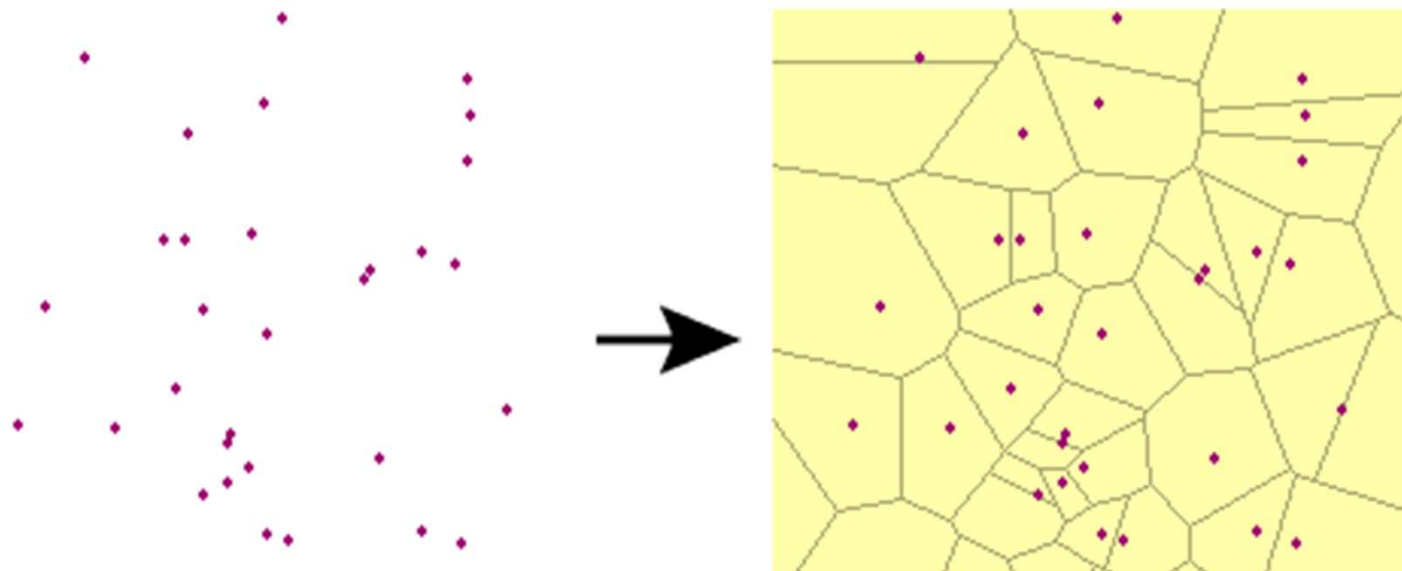Range query: $\mathbf{O(dn^{1-1/d})}$

Examine nearby points first: Explore the branch of the tree that is closest to the query point first.
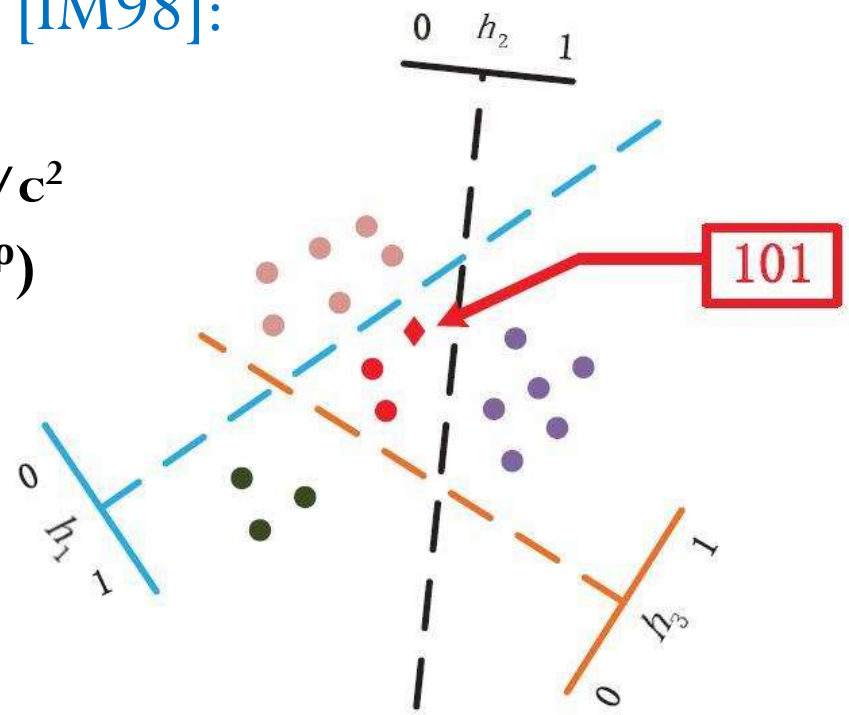
# Challenges of NNS

- Curse of dimensionality:
  - Given a polynomial indexing space $\mathbf{n^{O(1)}d^{O(1)}}$, existing a sublinear time algorithm to solve exact NNS refutes SETH [Wil18].
  - Indexing space or query time must be exponential in $\mathbf{d}$.

- Classic solutions: Voronoi decomposition
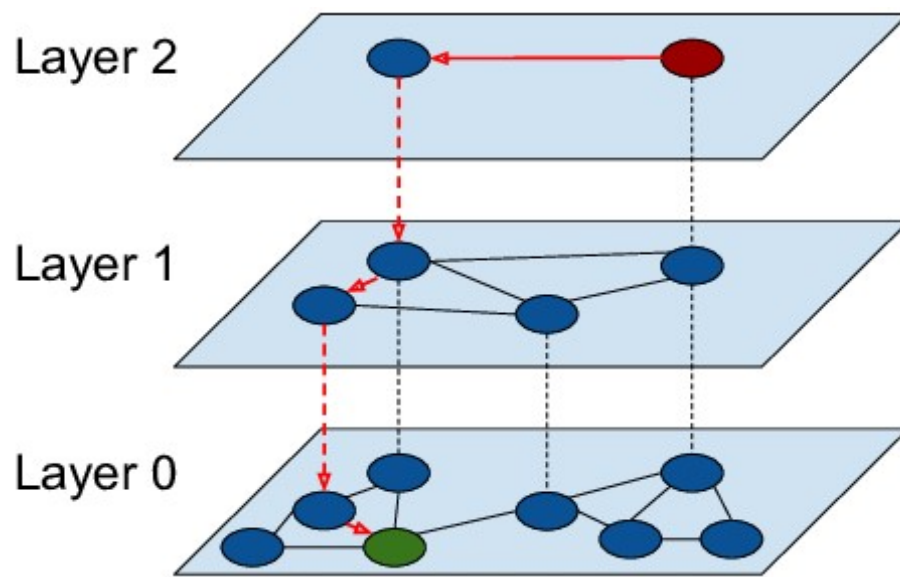


Space: $\mathbf{O(n^{O(d)})}$
Time: $\mathbf{O(d^{O(1)}\log n)}$

# Approximate NNS

- Approximate NNS (ANNS):
    - Returns **x'** such that $\mathbf{dist(q, x')} \leq \mathbf{c\ dist(q, x)}$

- Locality-sensitive hashing (LSH) [IM98]:
    - Theoretical guarantee to find **x'**
    - Sublinear query time $\mathbf{O(n^\rho)}$, $\boldsymbol{\rho} \approx \mathbf{1/c^2}$
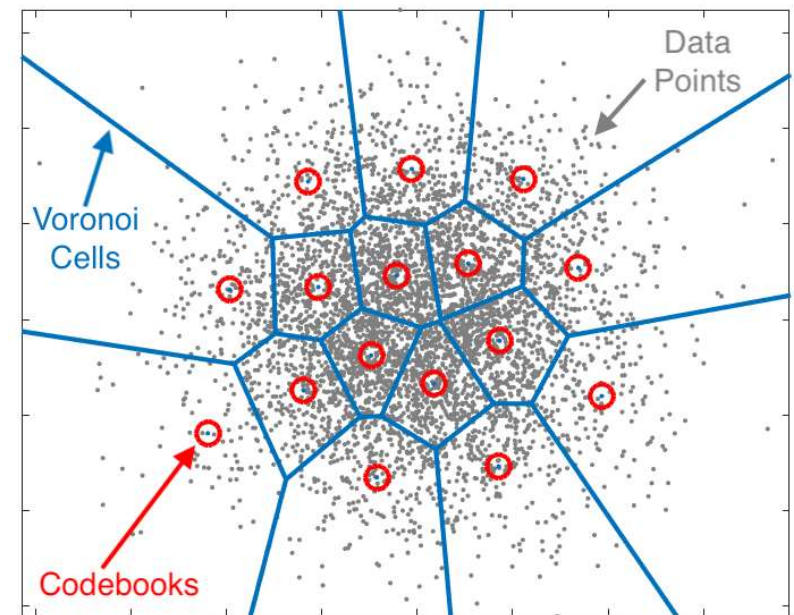    - Subquadratic indexing space $\mathbf{O(n^{1+\rho})}$

    - **But**…

# Approximate NNS

- Approximate NNS (ANNS):
  - Seek high empirical search recalls for top-$k$ NNS (e.g. tiny $c$)
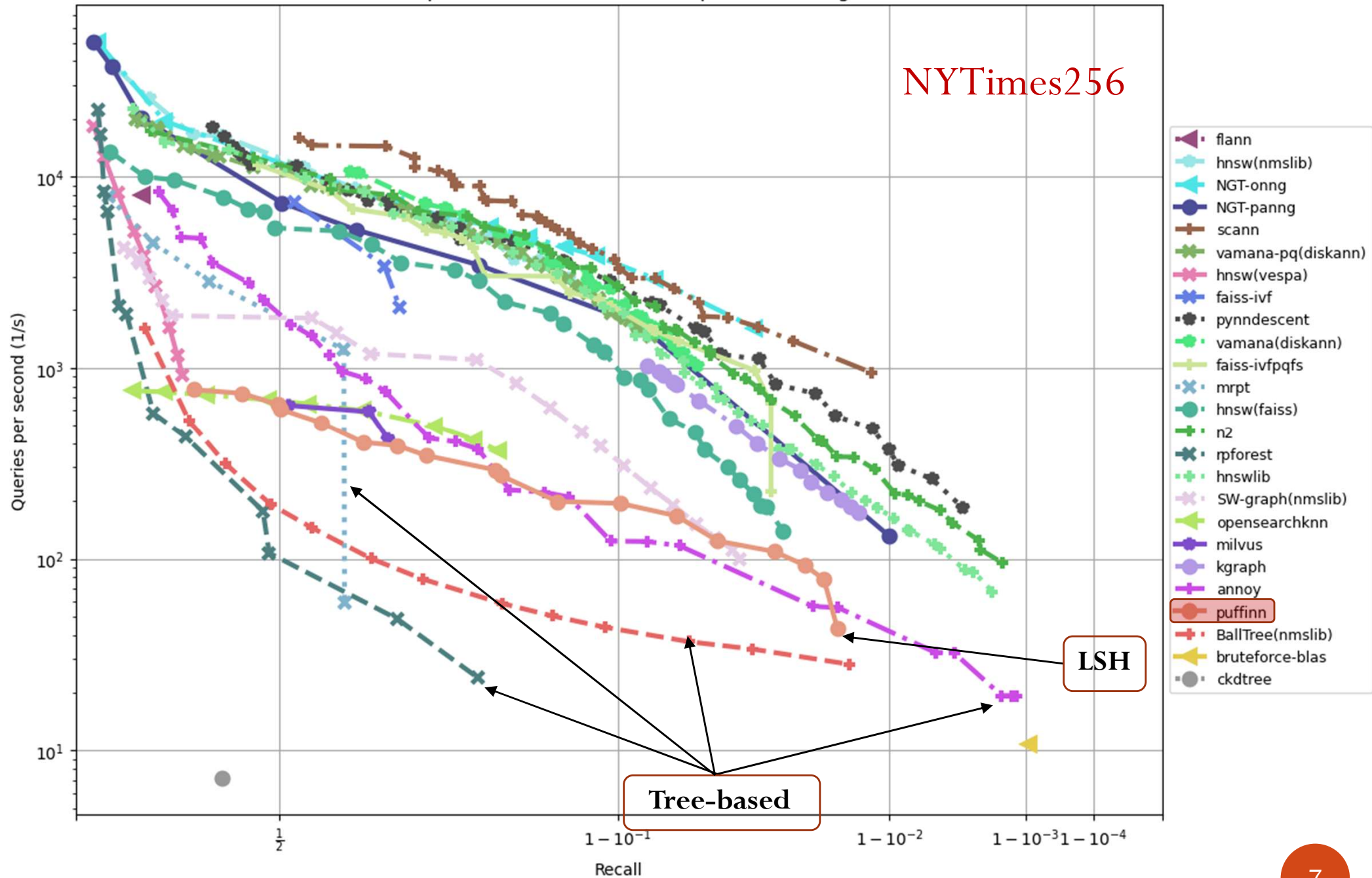
- Data-dependent approaches surpass LSH.



Graph-based index



Vector quantization

Recall-Queries per second (1/s) tradeoff - up and to the right is better

NYTimes256

https://github.com/erikbern/ann-benchmarks/

7

Recall-Queries per second (1/s) tradeoff - up and to the right is better

Glove100

LSH

Tree-based

https://github.com/erikbern/ann-benchmarks/

# Falconn++

- A practical locality-sensitive filtering approach
    - Lower query time complexity than Falconn, an optimal LSH scheme on angular distance.
    - Empirical higher recall-speed tradeoffs than Falconn
    - Competitive with HNSW on high recall regimes



1 thread

64 threads

(c) Glove200

(c) Glove200

9

# Locality-sensitive hashing (LSH)

- Definition [IM98]:
  - Given a distance function **dist(. , .)** and positive values **r, c, $p_1$, $p_2$** where $p_1 > p_2$, **c > 1**. A family of functions **H** is called **(r, cr, $p_1$, $p_2$)**-sensitive if for uniformly chosen **h ∈ H** and all **x, y ∈ $R^d$** :
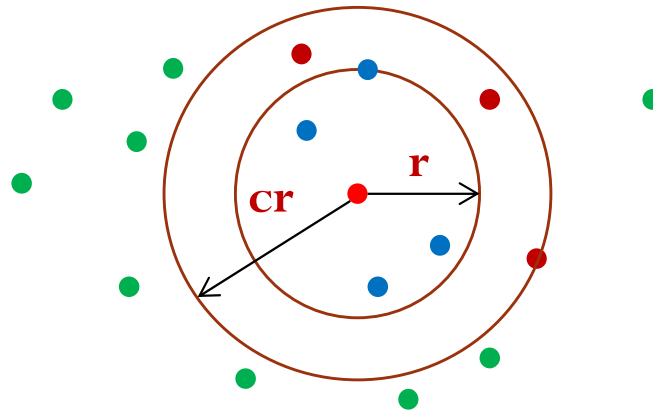    - If **dist(x, y) ≤ r** then **Pr [h(x) = h(y)] ≥ $p_1$**; (close points)
    - If **dist(x, y) ≥ cr** then **Pr [h(x) = h(y)] ≤ $p_2$**. (far away points)

# Hash tables construction



$$g_j(q) = \left( h_j^1(q), h_j^2(q), \ldots, h_j^k(q) \right), h \in H$$

Parameter settings:

$$\mathbf{k = \ln(n)/\ln(1/p_2)}$$
$$\mathbf{\rho = \ln(1/p_1)/\ln(1/p_2)}$$
$$\mathbf{L = n^\rho}$$

Space: $\mathbf{O(dn + n^{1+\rho})}$

11

# Hash tables lookup



Filtering

Refinement

$g_1(q)$

$g_2(q)$

$g_L(q)$

$q$

$T_1$

$T_2$

$T_L$

$d(q,x) \leq r$

$$g_j(q) = \left( h_j^1(q), h_j^2(q), \ldots, h_j^k(q) \right), h \in H$$

Time: Hashing time + $\mathbf{O(dn^\rho)}$

12

# Falconn [AIL+15]

- Construct a spherical Voronoi using **D** random vectors $\mathbf{r_i} \sim \mathbf{N^d(0, 1)}$

- **x** and **q** collide if sharing the closest or furthest random vector.

- Assume $\mathbf{r}_1 = \arg\max_{\mathbf{r}_i} |\mathbf{q}^\top \mathbf{r}_i|$,

$$h(\mathbf{q}) = \begin{cases} \mathbf{r}_1 \text{ if } \mathbf{sgn}(\mathbf{q}^\top \mathbf{r}_1) \geq 0, \\ -\mathbf{r}_1 \text{ otherwise}. \end{cases}$$

- Example:
  - $h(x) = h(z) = h(q) = \mathbf{r}_1$
  - $h(y) = -\mathbf{r}_2$



13

# Falconn's takeaways

- Given **D** random vectors, if **dist(x, q) = r**, then

$$\mathbf{Pr}\left[h(\mathbf{x}) = h(\mathbf{q})\right] \approx D^{-\frac{1}{4/r^2-1}}$$

- Given a small **c**, **(r, cr, p$_1$, p$_2$)**-sensitive Falconn has

$$\rho \approx \frac{4/c^2 r^2 - 1}{4/r^2 - 1} \approx 1/c^2$$

- Lower bound [OZW14]: $\rho \geq 1/c^2 - o(1)$ if $\mathbf{p_2} \geq 1/n$

# Practical multi-probe Falconn

- To reduce **L**, probe the bucket of the next closest or furthest random vectors

    - Require a large **qProbes**
    - Compute up to **0.1n** distances to achieve recall of 90% in Glove300

- Example:

    - **q** is next closest to $-r_2$
    - The blue wedge is the next probe

# Geometric intuition

- Use random projections to find $\mathbf{x}$ s.t. $\mathbf{x}^T\mathbf{q}$ is maximum.

# Geometric intuition

- Use a large number of random projections.



Among **5** random vectors $\mathbf{r_i}$, $\mathbf{r_1}$ is closest to $\mathbf{q}$ and the projections into $\mathbf{r_1}$ preserves the dot product order.

$\mathbf{x_1}$

$\mathbf{x_3}$

$\mathbf{r_2}$

$\mathbf{x_2}$

$\mathbf{r_3}$

$\mathbf{q}$

$\mathbf{q^T x_3}$    $\mathbf{q^T x_2}$    $\mathbf{q^T x_1}$

$\mathbf{r_1}^T\mathbf{x_3}$    $\mathbf{r_1}^T\mathbf{x_2}$    $\mathbf{r_1}^T\mathbf{x_1}$

$\mathbf{r_1}$

$\mathbf{r_4}$

$\mathbf{r_5}$

17

# CEOs [Pha21]

- Use a large number of random projections.



Indexing: Generate many random vectors $\mathbf{r_i}$, precompute the dot product order $\mathbf{L_i}$ for $\mathbf{r_i}$

Querying: Find $\mathbf{r_i}$ closest to $\mathbf{q}$ and return the top-1 NNS from $\mathbf{L_i}$

$\mathbf{x_1}$

$\mathbf{x_3}$

$\mathbf{x_2}$

$\mathbf{r_2}$

$\mathbf{r_3}$

$\mathbf{q^T x_3}$   $\mathbf{q^T x_2}$   $\mathbf{q^T x_1}$

$\mathbf{q}$

$\mathbf{r_1^T x_3}$   $\mathbf{r_1^T x_2}$   $\mathbf{r_1^T x_1}$

$\mathbf{r_1}$

$\mathbf{r_4}$

$\mathbf{r_5}$

# CEOs: Dimensionality reduction



Among **5** random vectors $\mathbf{r_i}$, we only use $\mathbf{r_1}$ to estimate dot products.

|       | $\mathbf{r_1}$ | $\mathbf{r_2}$ | $\mathbf{r_3}$ | $\mathbf{r_4}$ | $\mathbf{r_5}$ |
|-------|------|------|------|------|------|
| $\mathbf{x_1}$ | 4 | -2 | -4 | -3 | -3 |
| $\mathbf{x_2}$ | 3 | -1 | -3 | -4 | -4 |
| $\mathbf{x_3}$ | 1 | -3 | -2 | -3 | -3 |
| $\mathbf{q}$ | 9 | 4 | -8 | -3 | 0 |

| $\mathbf{x_1}$ | 1.2 | 0.9 |
| $\mathbf{x_2}$ | 0.5 | 0.5 |
| $\mathbf{x_3}$ | 0.1 | 0.7 |
| $\mathbf{q}$ | 1.5 | 0.1 |

19

# Concomitants of Extreme Order statistics

Random projections: Using $\mathbf{D}$ random vectors $\mathbf{r_i} \sim \mathbf{N^d(0, 1)}$, we have $\mathbf{D}$ bivariate samples $\mathbf{(Q_i, X_i)}$ from $\mathbf{N(0, 0, 1, \| x \|,} \mathbf{x^Tq)}$ where $\mathbf{Q_i = q^Tr_i}$ and $\mathbf{X_i = x^Tr_i}$

$X_i = x^Tr_i$

| | |
|---|---|
| 1.2 | 0.9 |

| | | | | |
|---|---|---|---|---|
| 5 | -2 | 4 | -3 | -3 |

$\mathbf{x}$

| | |
|---|---|
| 1.5 | 0.1 |

| | | | | |
|---|---|---|---|---|
| 8 | 4 | 9 | -3 | 0 |

$\mathbf{q}$

$Q_i = q^Tr_i$

Assume $\| \mathbf{q} \| = 1$

20

# Concomitants of Extreme Order statistics

$\mathbf{x}$

| 1.2 | 0.9 |
|-----|-----|

$X_i = x^T r_i$

| $X_{[2]}$ | $X_{[3]}$ | $X_{[1]}$ | $X_{[5]}$ | $X_{[4]}$ |
|-----|-----|-----|-----|-----|
| 5 | -2 | 4 | -3 | -3 |

$\mathbf{q}$

| 1.5 | 0.1 |
|-----|-----|

$Q_i = q^T r_i$

| 8 | 4 | 9 | -3 | 0 |
|-----|-----|-----|-----|-----|
| $Q_{(2)}$ | $Q_{(3)}$ | $Q_{(1)}$ | $Q_{(5)}$ | $Q_{(4)}$ |

Assume $\| \mathbf{q} \| = 1$

Random projections: Using $\mathbf{D}$ random vectors $\mathbf{r_i} \sim \mathbf{N^d(0, 1)}$, we have $\mathbf{D}$ bivariate samples $\mathbf{(Q_i, X_i)}$ from $\mathbf{N(0, 0, 1, \| x \|,}$ $\mathbf{x^T q)}$ where $\mathbf{Q_i = q^T r_i}$ and $\mathbf{X_i = x^T r_i}$

Order statistics: Sort $\mathbf{D}$ pairs $\mathbf{(Q_i, X_i)}$ by $Q$-value, we form the order statistics where $\mathbf{Q_{(1)}}$ is the first order statistics and $\mathbf{X_{[1]}}$ is the concomitant of the first order statistics.
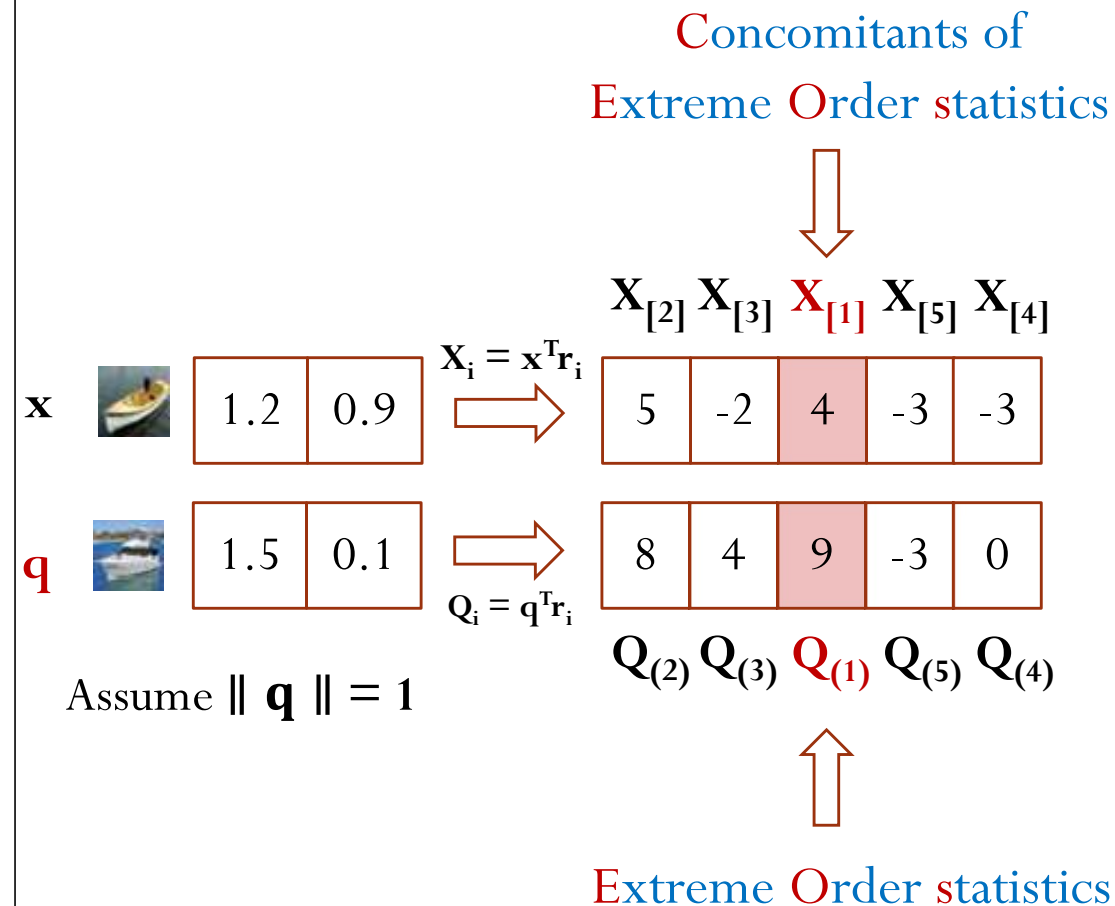
# Concomitants of Extreme Order statistics

Concomitants of
Extreme Order statistics

⬇

$X_{[2]}$ $X_{[3]}$ $X_{[1]}$ $X_{[5]}$ $X_{[4]}$

**x**

$X_i = x^T r_i$

| 1.2 | 0.9 |
|---|---|

| 5 | -2 | 4 | -3 | -3 |
|---|---|---|---|---|

**q**

| 1.5 | 0.1 |
|---|---|

$Q_i = q^T r_i$

| 8 | 4 | 9 | -3 | 0 |
|---|---|---|---|---|

$Q_{(2)}$ $Q_{(3)}$ $Q_{(1)}$ $Q_{(5)}$ $Q_{(4)}$

Assume $\| q \| = 1$

⬆

Extreme Order statistics

Random projections: Using **D** random vectors $r_i \sim N^d(0, 1)$, we have **D** bivariate samples $(Q_i, X_i)$ from $N(0, 0, 1, \| x \|, x^T q)$ where $Q_i = q^T r_i$ and $X_i = x^T r_i$

Order statistics: Sort **D** pairs $(Q_i, X_i)$ by Q-value, we form the order statistics where $Q_{(1)}$ is the first order statistics and $X_{[1]}$ is the concomitant of the first order statistics.

Extreme order statistics: When **D** is sufficiently large, $Q_{(1)}$ is the extreme order statistics and $X_{[1]}$ is the concomitant of the extreme order statistics.

22

# Theory of Concomitants of Extreme Order statistics [DG74]

Concomitants of
Extreme Order statistics

$\Downarrow$

$\mathbf{X}_{[2]}\ \mathbf{X}_{[3]}\ \mathbf{X}_{[1]}\ \mathbf{X}_{[5]}\ \mathbf{X}_{[4]}$

$\mathbf{x}$

$X_i = x^T r_i$

| 1.2 | 0.9 |
|---|---|

| 5 | -2 | 4 | -3 | -3 |
|---|---|---|---|---|

$\mathbf{q}$

| 1.5 | 0.1 |
|---|---|

$Q_i = q^T r_i$

| 8 | 4 | 9 | -3 | 0 |
|---|---|---|---|---|

$\mathbf{Q}_{(2)}\ \mathbf{Q}_{(3)}\ \mathbf{Q}_{(1)}\ \mathbf{Q}_{(5)}\ \mathbf{Q}_{(4)}$

Assume $\|\ \mathbf{q}\ \| = 1$

$\Uparrow$

Extreme Order statistics

Extreme order statistics: $\mathbf{Q}_{(1)}$ is the maximum variable among $\mathbf{D}$ random variables $\mathbf{Q_i = q^T r_i \sim N(0, 1)}$.

Extreme order statistics:

$$E[\mathbf{Q}_{(1)}] \approx \sqrt{\mathbf{2ln(D)}},\ Var[\mathbf{Q}_{(1)}] \approx 0$$

Concomitant of extreme order statistics:

$$\mathbf{X}_{[1]} \sim N(x^T q \sqrt{\mathbf{2ln(D)}},\ \|\ \mathbf{x}\ \|^2 - (x^T q)^2)$$

# Theory of Concomitants of Extreme Order statistics [DG74]

Concomitants of

Extreme Order statistics

$\mathbf{X}_{[D-1]}$    $\mathbf{X}_{[1]}$    $\mathbf{X}_{[2]}$   $\mathbf{X}_{[D]}$

$\mathbf{x}$

$\mathbf{X}_i = \mathbf{x}^T \mathbf{r}_i$

| 1.2 | 0.9 |
|---|---|

| -5 | -2 | 4 | $\cdots$ | 3 | -8 |
|---|---|---|---|---|---|

$\mathbf{q}$

| 1.5 | 0.1 |
|---|---|

| -8 | 4 | 9 | $\cdots$ | 7 | -9 |
|---|---|---|---|---|---|

$\mathbf{Q}_i = \mathbf{q}^T \mathbf{r}_i$

$\mathbf{Q}_{(D-1)}$    $\mathbf{Q}_{(1)}$    $\mathbf{Q}_{(2)}$   $\mathbf{Q}_{(D)}$

Assume $\| \mathbf{q} \| = 1$

Extreme Order statistics

Top $\mathbf{s}_0$ maximum and minimum order statistics: $\mathbf{Q}_{(i)}$ and $\mathbf{Q}_{(D-i+1)}$ where $\mathbf{i} = 1, \ldots, \mathbf{s}_0$

Extreme order statistics:

$$\mathbf{E}[\mathbf{Q}_{(i)}] \approx \sqrt{2\ln(D)}$$
$$\mathbf{E}[\mathbf{Q}_{(D-i+1)}] \approx -\sqrt{2\ln(D)}$$

Concomitant of extreme order statistics:

$$\mathbf{X}_{[i]} \sim \mathrm{N}(\mathbf{x}^T\mathbf{q}\sqrt{2\ln(D)}, \; \| \mathbf{x} \|^2 - (\mathbf{x}^T\mathbf{q})^2)$$
$$\mathbf{X}_{[D-i+1]} \sim \mathrm{N}(-\mathbf{x}^T\mathbf{q}\sqrt{2\ln(D)}, \; \| \mathbf{x} \|^2 - (\mathbf{x}^T\mathbf{q})^2)$$

$\mathbf{X}_{[i]}$ and $\mathbf{X}_{[D-i+1]}$ are independent asymptotically.

24

# Connection to multi-probe Falconn

- Falconn: If $\mathbf{r}_1 = \arg\max_{\mathbf{r}_i} |\mathbf{q}^\top \mathbf{r}_i|$,
  - Use $\mathbf{r}_1$ corresponding to $\mathbf{Q}_{(1)}$ as hash value
  - Use $\mathbf{Q}_{(i)}$ and $\mathbf{Q}_{(D-i)}$ as probing buckets where $\mathbf{i} = 1, \dots, \mathbf{s}_0$

- CEOs: If $\mathbf{r}_1 = \arg\max_{\mathbf{r}_i} |\mathbf{q}^\top \mathbf{r}_i|$,
  - Use $\mathbf{X}_{[1]} = \mathbf{x}^\mathbf{T}\mathbf{r}_1$ to estimate $\mathbf{x}^\mathbf{T}\mathbf{q}$
  - Use $\mathbf{X}_{[i]}$ and $\mathbf{X}_{[D-i]}$ as estimators of $\mathbf{x}^\mathbf{T}\mathbf{q}$ where $\mathbf{i} = 1, \dots, \mathbf{s}_0$

- Falconn++ = Falconn + CEOs

Partition **n** points into **2D** buckets

Scale each bucket by keeping **x** with largest $\mathbf{x}^\mathbf{T}\mathbf{r}_\mathbf{i}$

25

# Falconn++: A locality-sensitive filtering

- A locality-sensitive filtering (LSF) mechanism:
  - Given a distance function **dist(. , .)** and positive values **r**, **c**, $q_1$, $q_2$ where $q_1 > q_2$, **c > 1**. For an **(r, cr, $p_1$, $p_2$)**-sensitive function **h** and **x, y** in **h(q)** :
    - If **dist(x, q) ≤ r** then **Pr [x** is not filtered**] ≥ $q_1$**
    - If **dist(y, q) ≥ cr** then **Pr [y** is not filtered**] ≤ $q_2$**

- Combine LSH and LSF:
  - **Pr [h(x) = h(q), x** is not filtered**] ≥ $p_1 q_1$**
  - **Pr [h(y) = h(q), y** is not filtered**] ≤ $p_2 q_2$**

- We need **ln(1/$q_1$) / ln(1/$q_2$) ≤ ρ ≈ 1/$c^2$** to achieve a new exponent **ρ' ≤ ρ**

# Falconn++

- Asymptotic property of CEOs:

  - $X_{[1]} \sim N(x^T q \sqrt{2\ln(D)}, \; \| x \|^2 - (x^T q)^2)$
  - $Y_{[1]} \sim N(y^T q \sqrt{2\ln(D)}, \; \| y \|^2 - (y^T q)^2)$

- Filtering mechanism:

  - Define a threshold $t = (1 - r^2/2)\sqrt{2\ln(D)}$. For each bucket corresponding to $r_i$, keep any point $x$ if $x^T r_i \geq t$. Otherwise, discard it.
  - Note: $dist(x, q) = r$, then $x^T q = 1 - r^2/2$

# Falconn++'s takeaways

- For sufficiently large **D** random projections, **c > 1**,

---

- If $\|\mathbf{x} - \mathbf{q}\| \leq r$, then $\mathbf{Pr}[\mathbf{x} \text{ is not filtered}] \geq q_1 = 1/2$;
- If $\|\mathbf{y} - \mathbf{q}\| \geq cr$, then $\mathbf{Pr}[\mathbf{y} \text{ is not filtered}] \leq q_2 = \frac{1}{\gamma\sqrt{2\pi}} \exp(-\gamma^2/2) < q_1$ where

$\gamma = \frac{cr(1-1/c^2)}{\sqrt{4-c^2 r^2}} \cdot \sqrt{2 \ln D}$.

---

- New exponent $\boldsymbol{\rho'} \approx 1/(2c^2 - 2 + 1/c^2)$

$$\rho' = \frac{\ln(1/q_1 p_1)}{\ln(1/q_2 p_2)} \approx \frac{\frac{\ln 2}{\ln D} + \frac{1}{4/r^2 - 1}}{\frac{(1-1/c^2)^2}{4/c^2 r^2 - 1} + \frac{1}{4/c^2 r^2 - 1}}$$

$$\approx \frac{1}{1 + (1-1/c^2)^2} \cdot \frac{4/c^2 r^2 - 1}{4/r^2 - 1} \leq \rho.$$



c=1.5

cr= √2

28

# Connection to LSF framework [ALRW17]

- Asymmetric LSF frameworks:
  - Apply different filtering conditions on data and query to govern the space-time tradeoff
  - Let $\mathbf{t_u}$ and $\mathbf{t_q}$ be two different thresholds.
  - Collision: $\mathbf{x}$ and $\mathbf{q}$ pass the filter $\mathbf{r_i}$ with $\mathbf{Pr[x^T r_i \geq t_u , q^T r_i \geq t_q]}$

- Falconn++:
  - Use a sufficiently large $\mathbf{D}$ to ensure the asymptotic property of CEOs
  - $\mathbf{t_u = (1 - r^2/2)\sqrt{2\ln(D)}}$ , $\mathbf{t_q \approx \sqrt{2\ln(D)}}$
  - For $\mathbf{p_2 q_2 = 1/n, D = O(n^{\rho'})}$, Falconn++ yields $\mathbf{O(n^{\rho'})}$ query time where $\mathbf{\rho' \approx 1/(2c^2 - 2 + 1/c^2)}$ for $\mathbf{r \geq \sqrt{2}}$

# Practical implementations

- Data-dependent setting:
    - Select a scaling factor $0 < \alpha < 1$ to scale each bucket of size **B** to $\alpha$ **B**
        - Adapt **t** to various density
        - Easy to govern the memory footprint (i.e. # points in a table)

- Multi-probe indexing:
    - For each point **x**, hash it into **iProbes** buckets corresponding the **iProbes** closest or furthest random vectors
    - Scale each bucket of size **B** to $\alpha$ **B/iProbes**

- Other heuristics:
    - Pseudo-random rotation $\mathbf{HD_3 HD_2 HD_1}$ to simulate random projections
    - Center the data point **X**
    - Limit scaling: keep $\mathbf{max(k, \alpha\ B/iProbes)}$ points in a bucket

# Falconn++: Scaling bucket

- Experiment on Glove200 with **1.2M** points with **k = 20**:
  - **D = 256**, 2 combined LSH functions, each table has $4D^2$ buckets
  - Falconn: **qProbes = {1000, ... , 20000}**, Falconn++: **qProbes/α**



(a) iProbes=1

(b) iProbes=1

# Falconn++: Multi-probe indexing

- Observation:
  - Overfitting: Large **iProbes** degrades performance
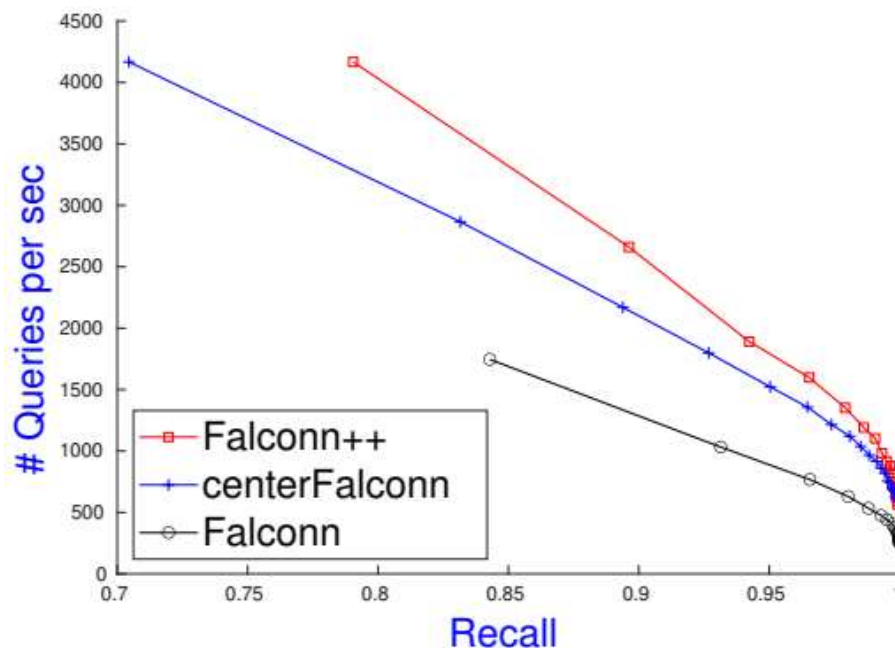  - **iProbes = 2D** is as similar as theoretical LSF framework

- Setting iProbes:
  - $k \approx n \cdot iProbes/4D^2$
  - Each bucket has roughly **k=20** points, especially sparse buckets
  - **iProbes = 3: 1.2M $\cdot$ 3/$2^{18}$=16**



(c) L=500, $\alpha$=0.1

Glove200, **k** = 20

32

# Falconn++: Centering data

- Experiment on Glove200 with **k = 20**:
  - **D = 256**, 2 combined LSH functions, **$4D^2$ buckets/table**, **L = 500**
  - Falconn++: **α = 0.1, iProbes = 3, qProbes/α**



(a) Speed-recall tradeoff

(b) # Dot products-recall tradeoff

# Falconn++: Limit scaling & centering

- Observations:
  - After centering, buckets are more balanced.
  - With **α = 0.1, iProbes = 3,** and keep **max(k, αB/iProbes)** points, # points/table ≈ **2.42 n**
  - Less **qProbes** and hash evaluation time than Falconn
- Future improvement:
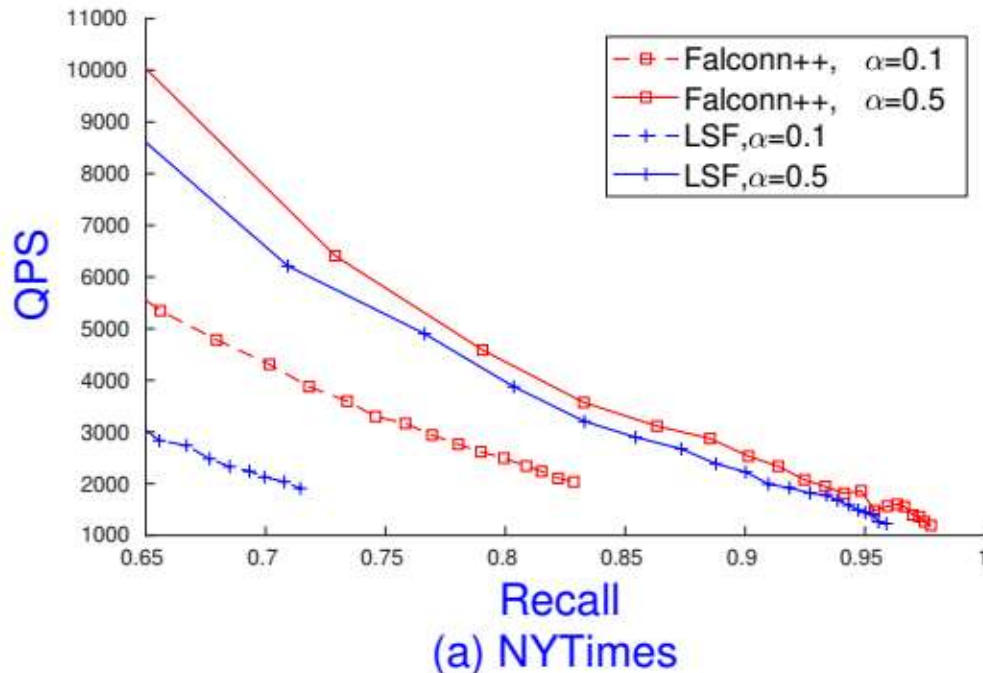  - SimHash signatures to reduce distance computation time



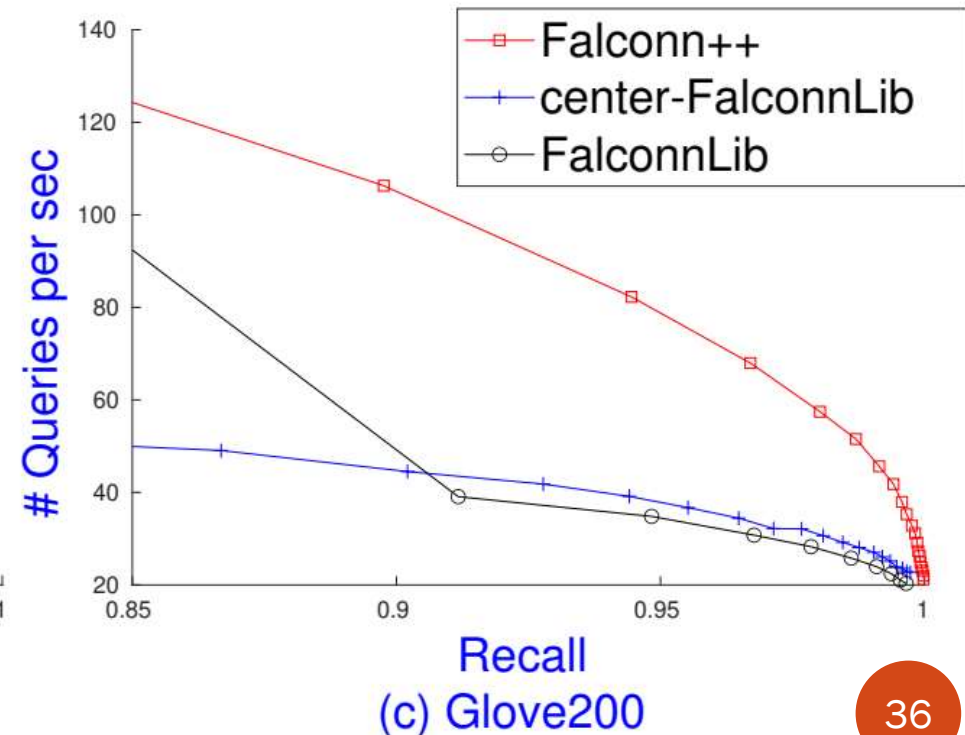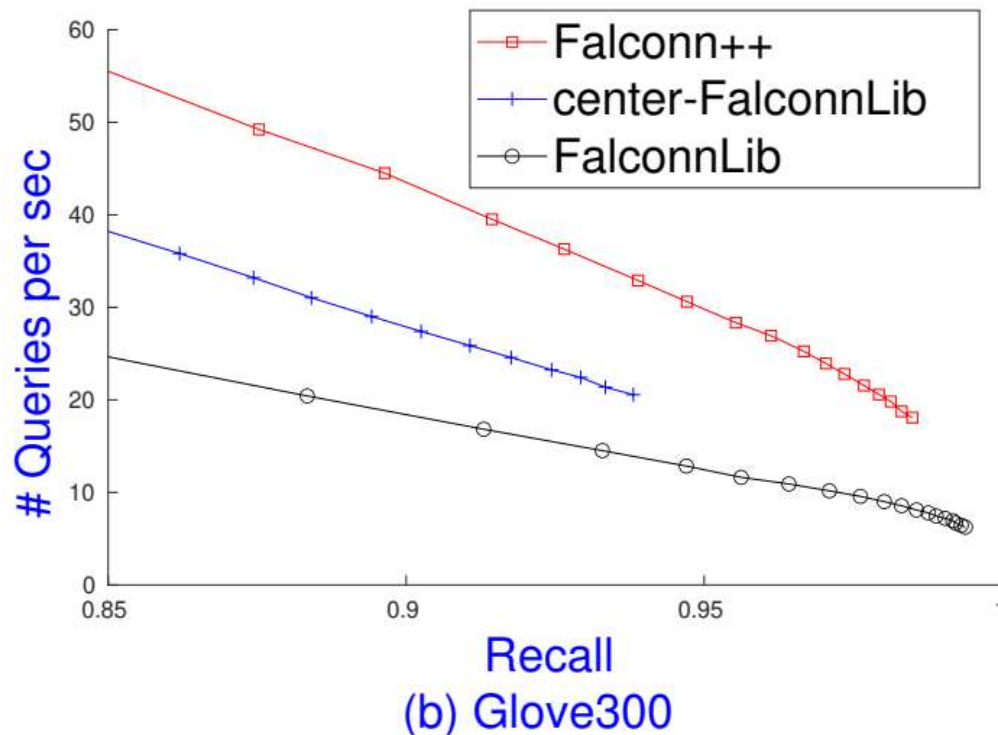(c) Query time components

Glove200, **k** = 20

34

# Falconn++ vs. Theoretical LSF

- LSF: Select $\mathbf{t_u}$ s.t. $\mathbf{Pr}\left[\mathbf{x}^\top \mathbf{r}_i \geq t_u\right]^2 = \alpha/4D^2$
  - Falconn++: No limit scaling, only centering
  - **iProbes = 1, D = {128, 256}, L = 100**, 2 combined LSH/LSF functions



(a) NYTimes  (c) Glove200

# Falconn++ vs. FalconnLib

- Glove300 and Glove200 with **k = 20** and **1 thread**:
  - **L = 500, D = 256, α = 0.1, iProbes = {1, 3}, 4D² buckets/table**
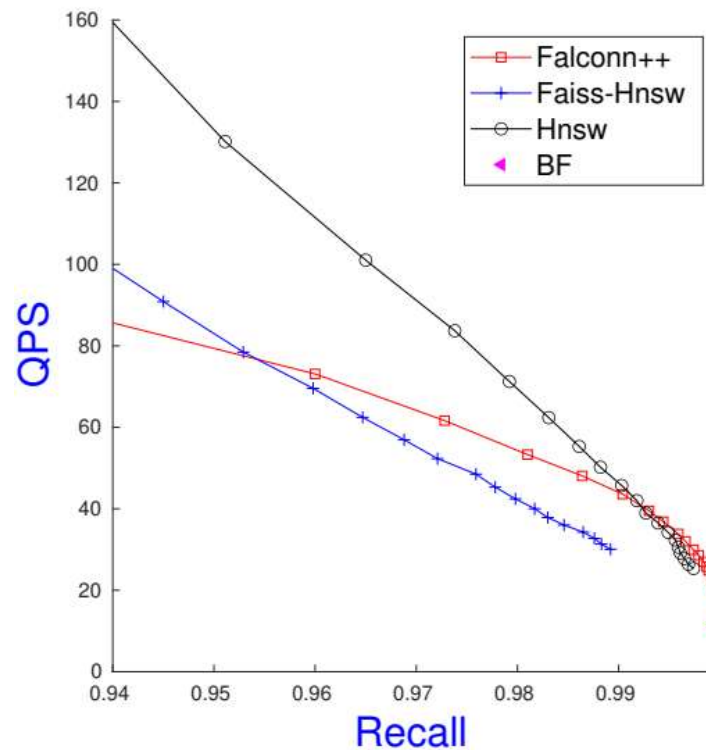  - Falconn: **L = 50** (Glove300), **L = 1210** (Glove200)



(b) Glove300

(c) Glove200

# Falconn++ vs Hnsw

- Parameter settings:
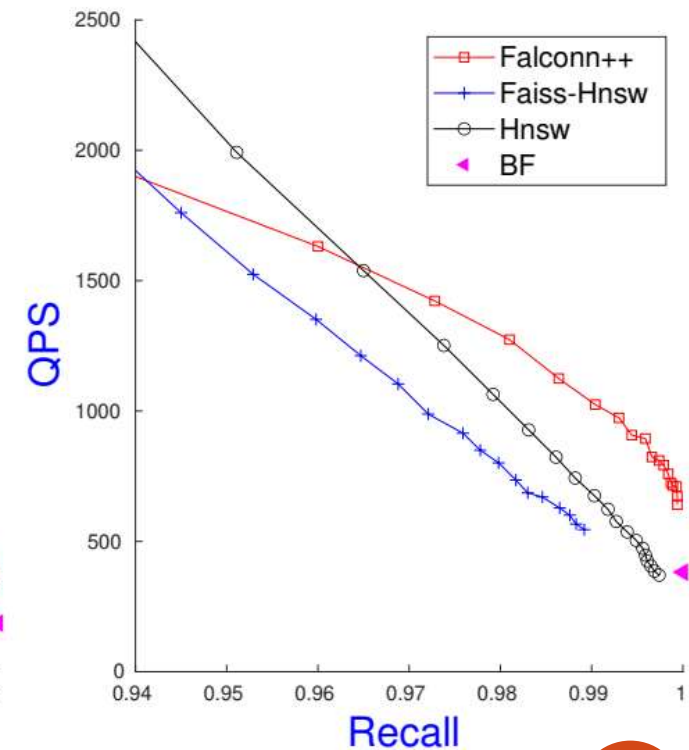  - Hnsw: **ef-index = 200, M = 512,** vary **ef-query**
  - Falconn++: **D = 256, L = 350, α = 0.01, iProbes = 3**, vary **qProbes**

Indexing

Space: **5.4GB**
Hnsw: **13.7 mins**
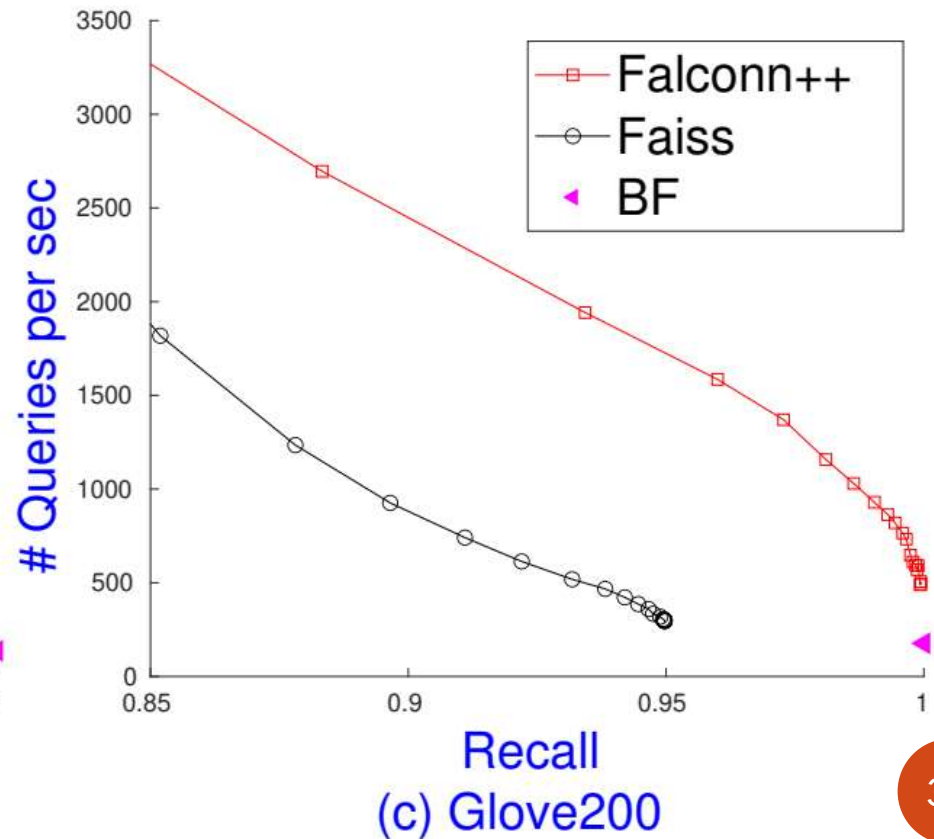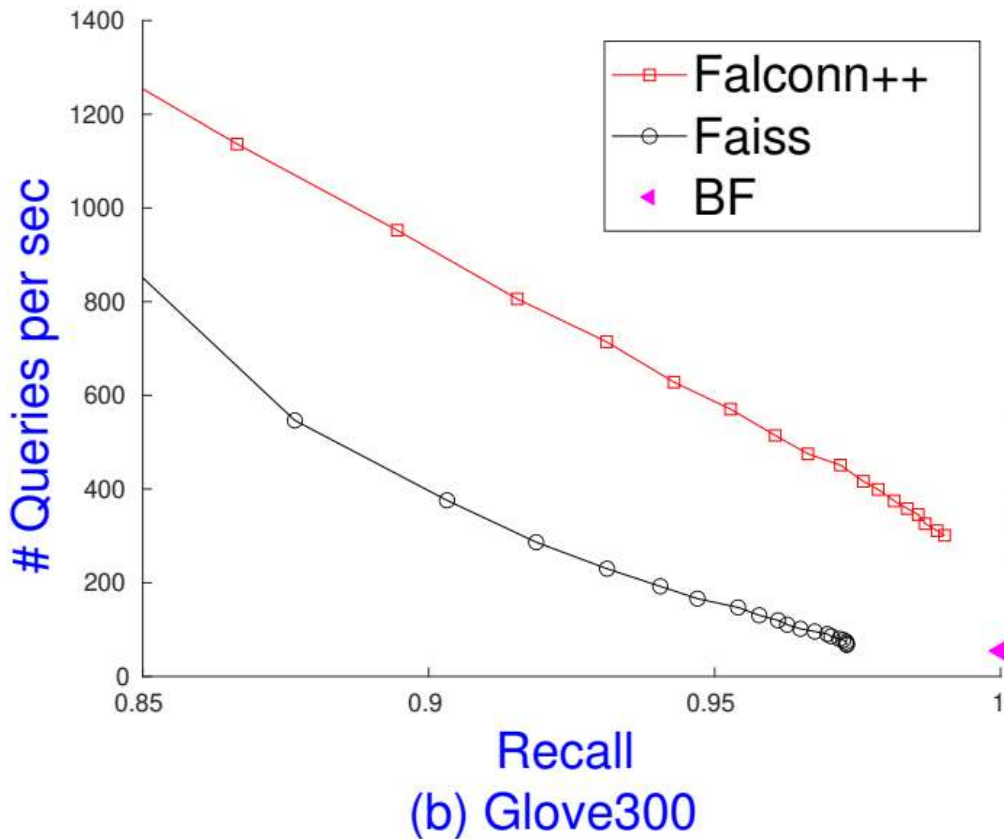Falconn++: **1.1 mins**



(c) Glove200, 1 thread

(d) Glove200, 64 thread

# Falconn++ vs Faiss

- Parameter settings:
  - Faiss: **m = 256, nlist=1000, 8 bits/centroid**, vary **probe**
  - Falconn++: **D = 256, L = 350, α = 0.01, iProbes = 3**, vary **qProbes**



(b) Glove300

(c) Glove200

# Open problem

- Practical LSH & LSF pattern:
  - Existing for Euclidean distance with $\rho = 1/c^2$

- Characterize # random projections **D**

- Characterize the scaling factor $\boldsymbol{\alpha}$