



SPEEDING UP PERMUTATION BASED INDEXING WITH INDEXING

Karina Figueroa
Universidad Michoacana
México

Kimmo Fredriksson
University of Kuopio
Finland



Content

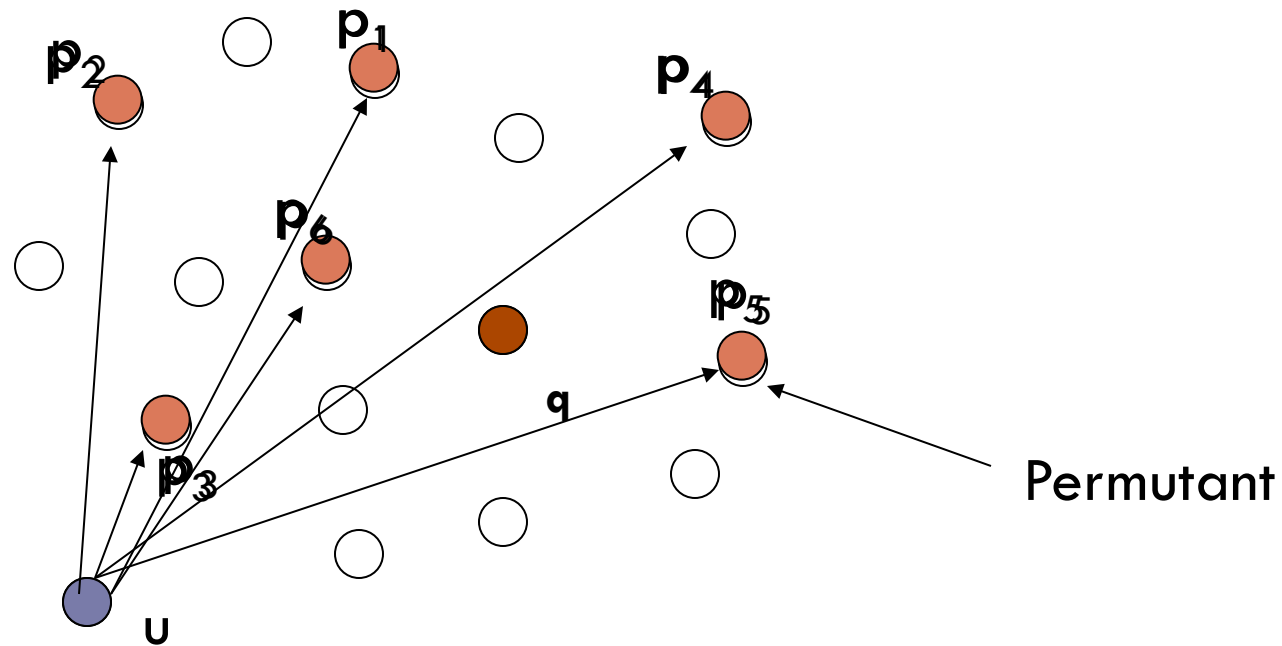
2

- Permutation-based algorithm
- Indexing Permutations
- Experimental Results
- Conclusions

Permutation based algorithms

3

IDEA



$$\begin{aligned}\Pi_u &= p_3, p_6, p_2, p_1, p_5, p_4 \\ \Pi_q &= p_6, p_5, p_4, p_1, p_3, p_2\end{aligned}$$

Spearman Footrule metric

4

$$\Pi_q = p_1, p_2, p_3, p_4, p_5, p_6$$

$$\Pi_u = p_3, p_6, p_2, p_1, p_5, p_4$$

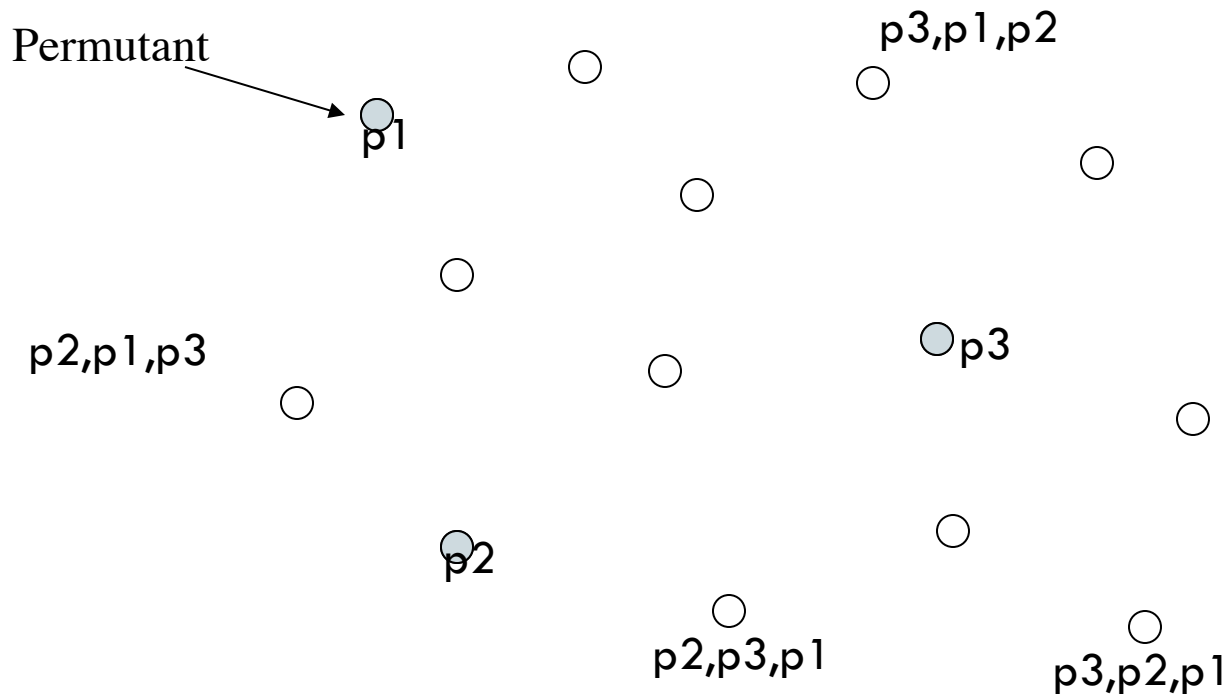
3-1, 6-2, 3-2, 4-1, 5-5, 6-4

Differences between positions

$$F(\Pi_q, \Pi_u) = |1 - 3| + |2 - 6| + |3 - 2| + |4 - 1| + |5 - 5| + |6 - 4| = 12$$

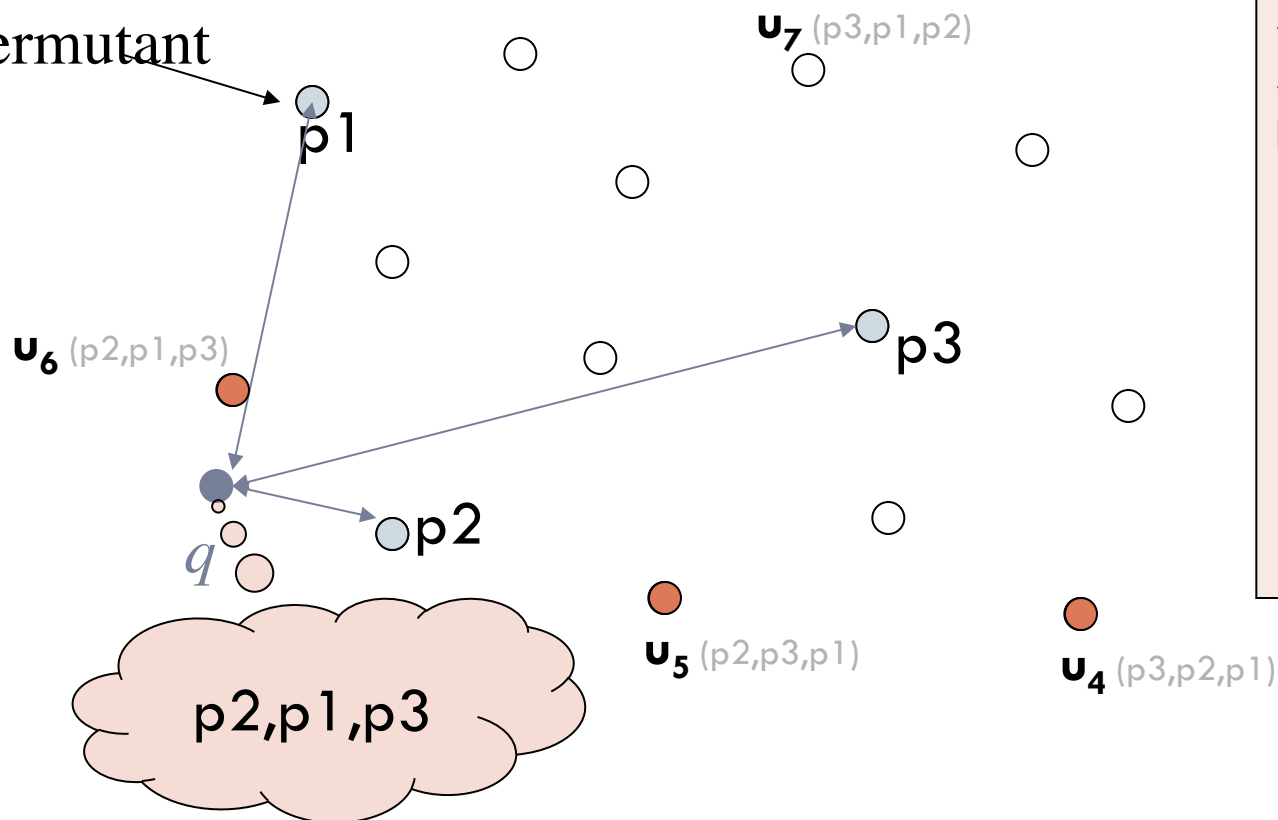
Permutation-based algorithm

Preprocessing



Search Process

Permutant



Sorting elements by
Spearman
Footrule

.....

Permutations based algorithm

Algorithm 1 Sequential-Range-Query(q, r, f)

- 7
- 1: **INPUT:** q is a query and r its radius, f is the fraction of the database to traverse.
 - 2: **OUTPUT:** Reports a subset of those $u \in \mathbb{U}$ that are at distance at most r to q .
 - 3: Let $A[1, n]$ be an array of tuples and $\mathbb{U} = \{u_1, \dots, u_n\}$
 - 4: Compute Π_q^{-1}
 - 5: **for** $i \leftarrow 1$ to n **do**
 - 6: $A[i] \leftarrow \langle u_i, S_\rho(\Pi_{u_i}, \Pi_q) \rangle$
 - 7: **end for**
 - 8: **SortIncreasing**(A) // by second component of tuples
 - 9: **for** $i \leftarrow 1$ to $f \cdot n$ **do**
 - 10: Let $\langle u, s \rangle = A[i]$
 - 11: **if** $d(q, u) \leq r$ **then**
 - 12: Report u
 - 13: **end if**
 - 14: **end for**
-

$O(kn + n \log n)$
Distance evaluations.

Our proposal

8

- Define a new metric space
 - ▣ The permutations of the database
 - ▣ The metric between permutations
- In this new metric space
 - ▣ We can use any algorithm for metric spaces
 - ▣ Finds the most similar permutations to the query
 - Avoids making comparisons through the whole database
 - Reduces the searching time

Indexing Permutations

9

Algorithm 2 Indexed-Range-Query(q, r, f)

- 1: **INPUT:** q is a query and r its radius, f is the fraction of the database to traverse.
 - 2: **OUTPUT:** Reports a subset of those $u \in \mathbb{U}$ that are at distance at most r to q .
 - 3: Let $A[1, f \cdot n]$ be an array of elements and W an index of permutations
 - 4: Compute Π_q^{-1}
 - 5: $A \leftarrow \text{NearestNeighbors}(W, \Pi_q, f \cdot n)$
 - 6: **for** $i \leftarrow 1$ to $f \cdot n$ **do**
 - 7: **if** $d(q, A[i]) \leq r$ **then**
 - 8: Report u
 - 9: **end if**
 - 10: **end for**
-

Nearest Neighbor query

10

- What Algorithm can we use?
 - ▣ Any for metric spaces
 - Because it is a metric space
 - ▣ LClusters
 - ▣ GHT
 - ▣ BKT
 - ▣ MVP
 - ▣ SAT
 - ▣ Also, permutation-based (again and again, etc)

Databases

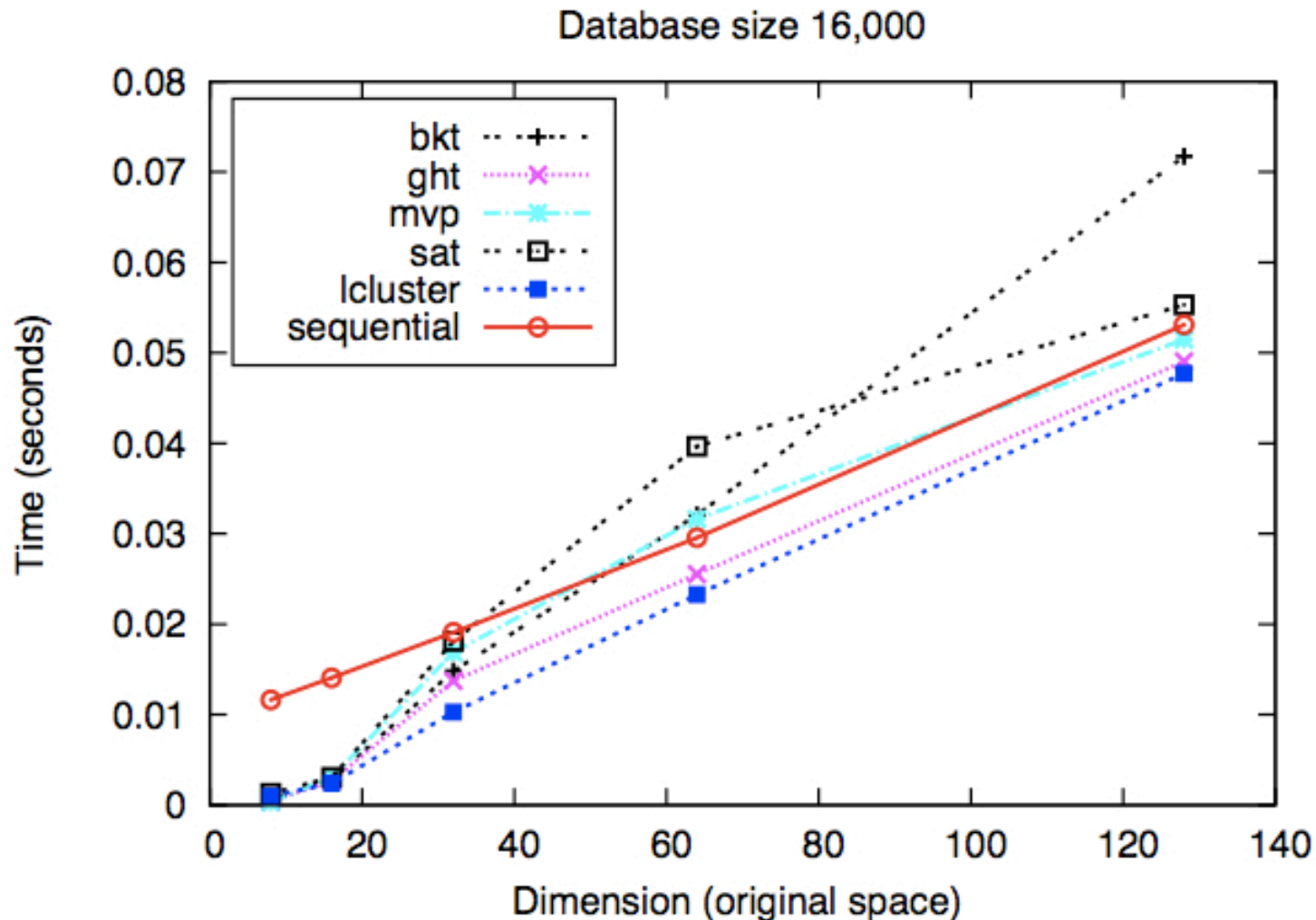
11

- Synthetic vectors in an unitary cube
 - ▣ Dimensions from 8 to 128
 - ▣ 16,000 and 40,000 objects
- Gaussians vectors
 - ▣ 10,000 object
- Spanish Language dictionary
 - ▣ 86,000 words
- Faces (762 dimensions)
 - FERET
 - 761 images

Experimental Results

Synthetic vectors in a Unitary cube

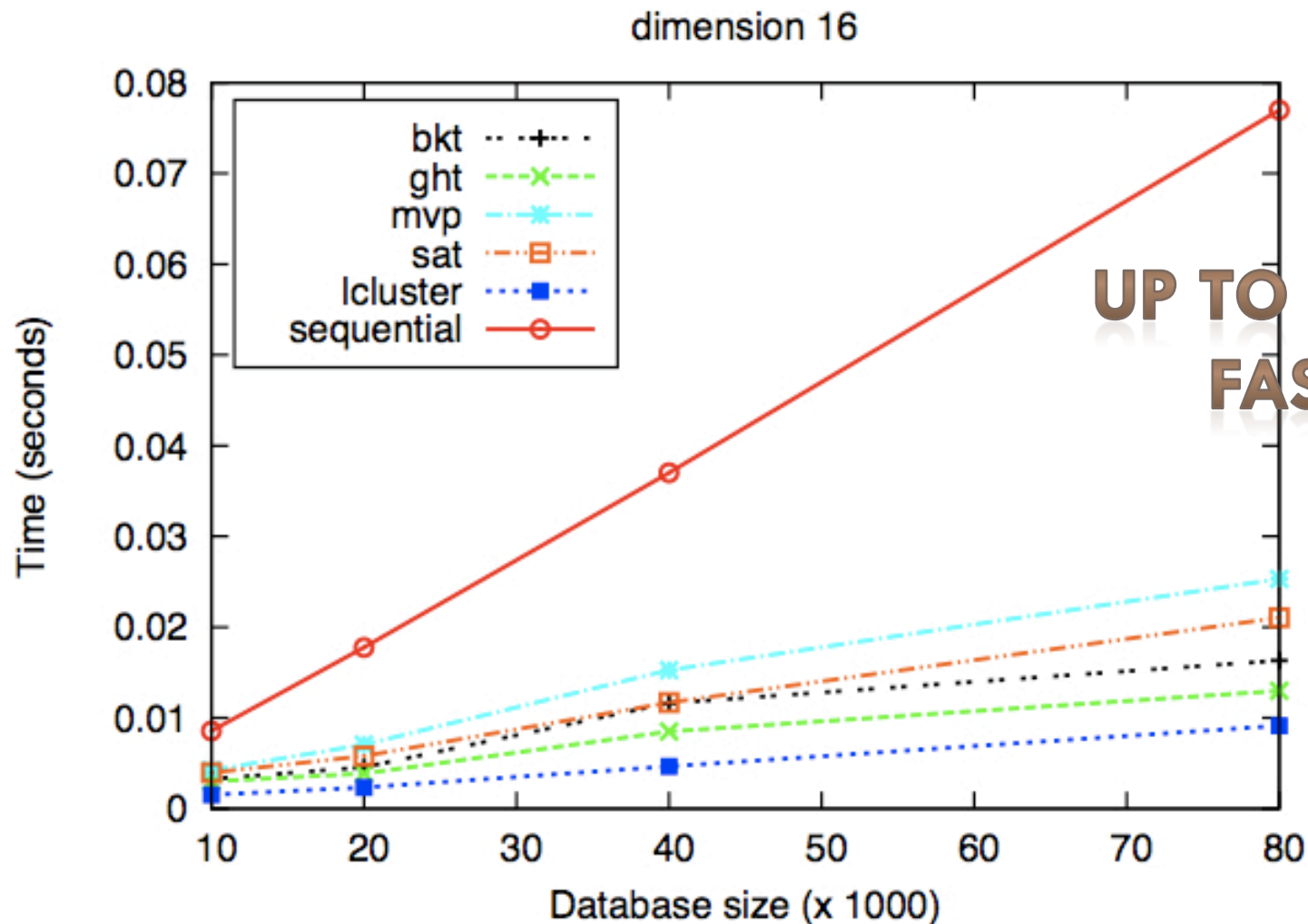
12



Experimental Results

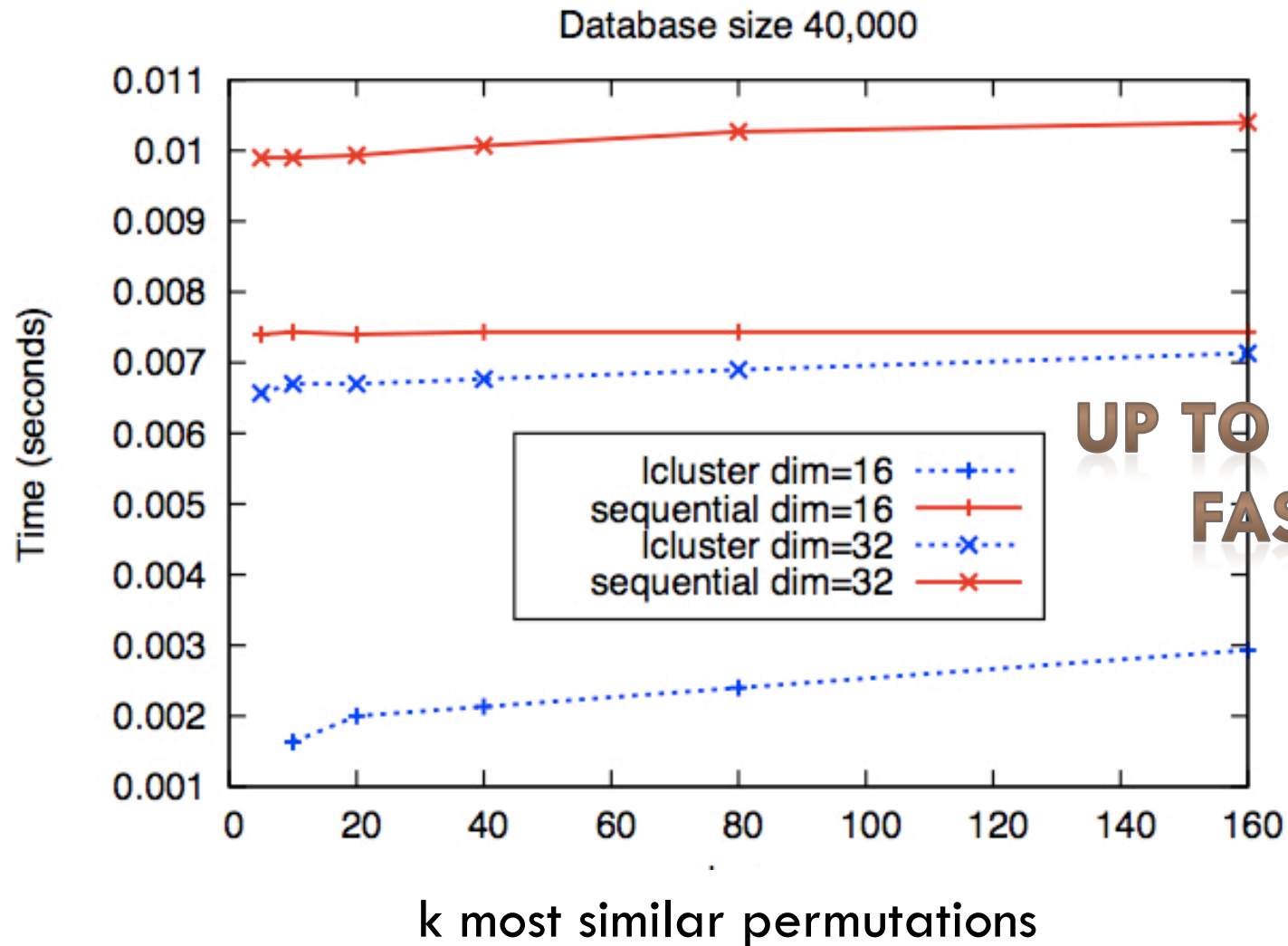
Synthetic vectors in a Unitary cube

13



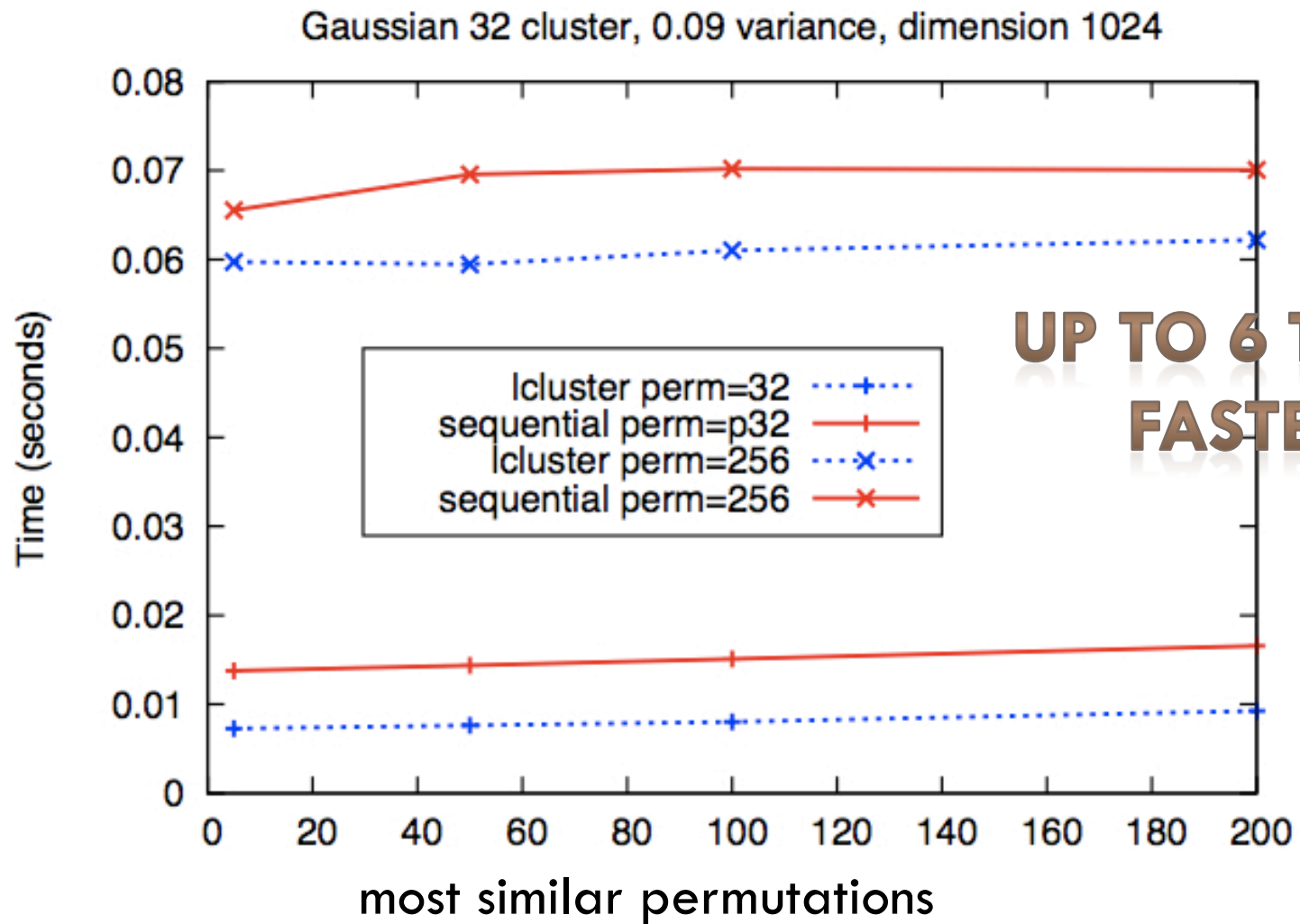
Vectors: Unitary cube

14



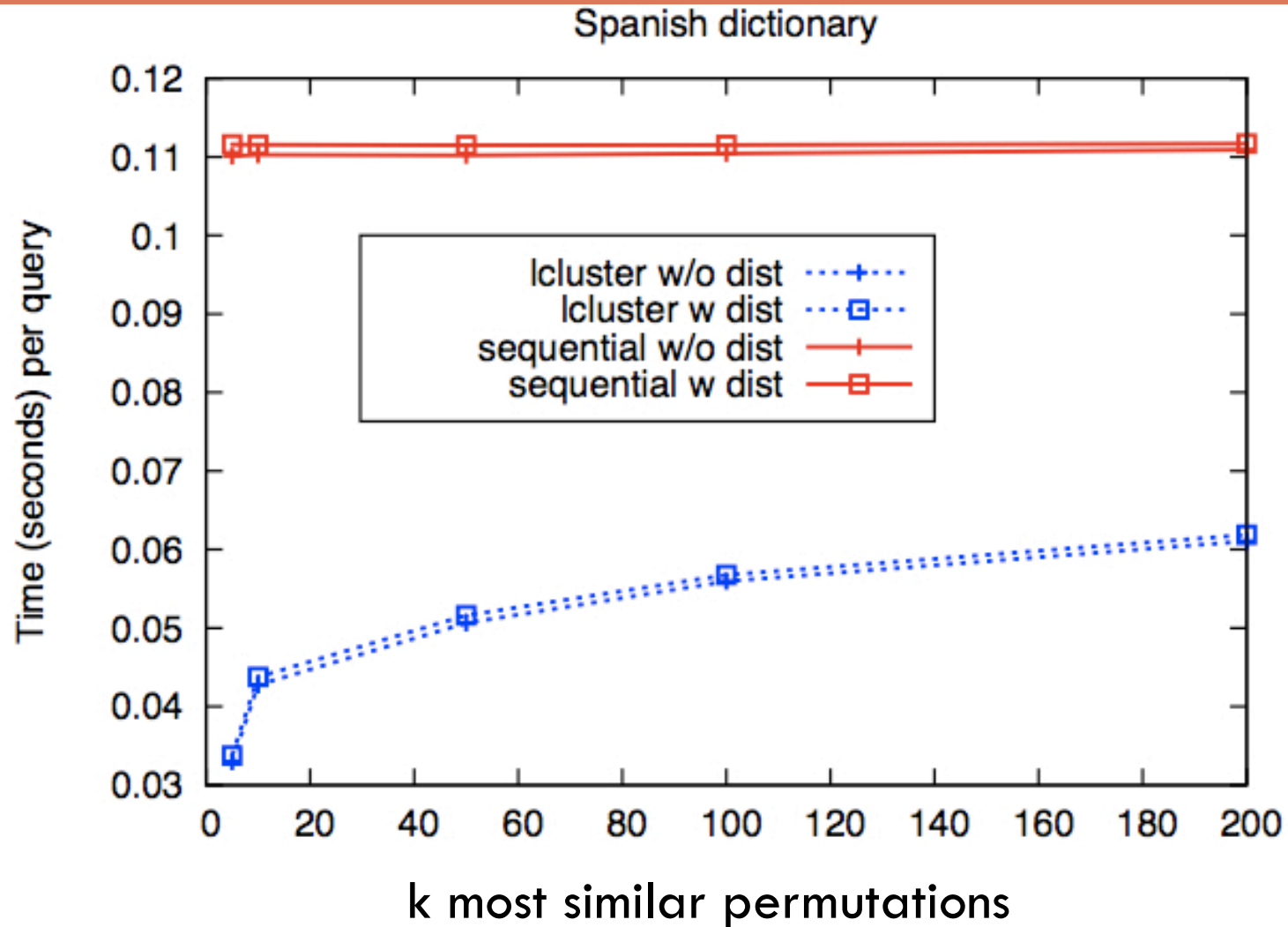
Gaussians vectors

15



Strings: Spanish Language dictionary

16



Conclusions

17

- We have significantly reduced the CPU time of one of the best algorithms for searching in very high dimensional metric spaces
- We were able to do this by identifying another metric space search problem in the algorithms internals.

18

Thanks

karina@fismat.umich.mx

Kimmo.frediksson@uku.fi

