# Fast *k*-Nearest Neighbour Search via Prioritized DCI
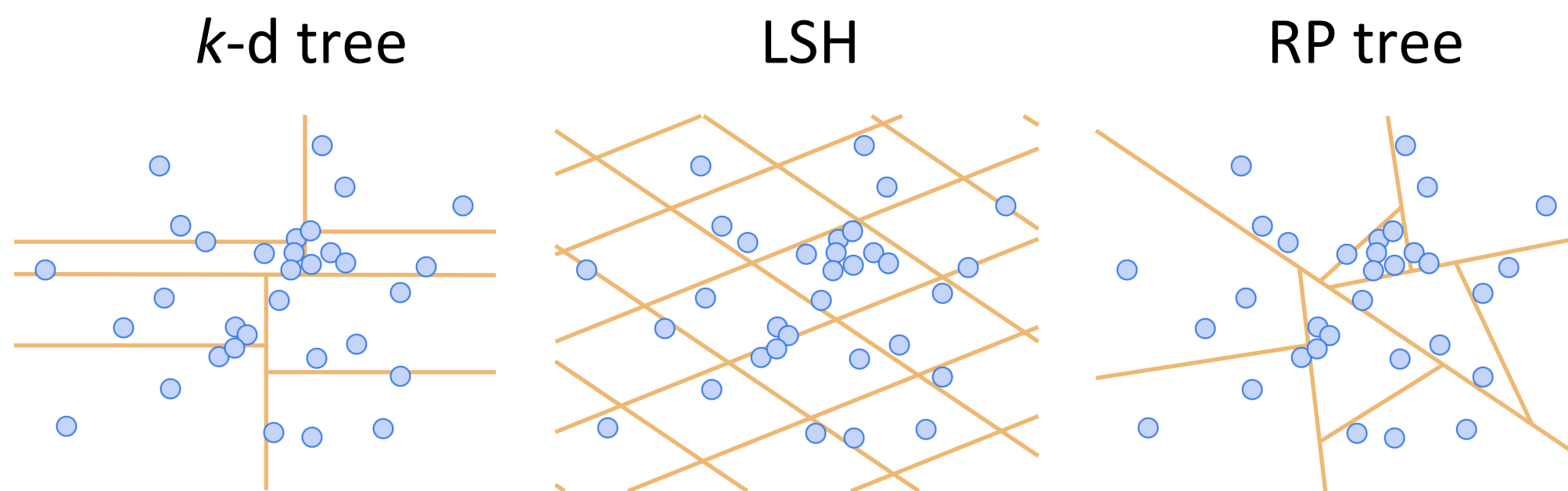
## Ke Li    Jitendra Malik

{ke.li,malik}@eecs.berkeley.edu

## Introduction

- Existing methods suffer from the curse of dimensionality.

- Most existing methods rely on a divide-and-conquer strategy known as space partitioning.

- We present a new algorithm that overcomes the curse of dimensionality, which has:

  - Time complexity: linear in ambient dimensionality, sub-linear in intrinsic dimensionality and size of the dataset.

  - Space complexity: independent of ambient dimensionality and linear in size of the dataset.

## The Case Against Space Partitioning

- Most existing methods rely on space partitioning:

  *k*-d tree          LSH          RP tree



- Problems:

  - As dimensionality increases, volume of space grows exponentially $\Rightarrow$ either the number or the size of cells must grow exponentially.

  - "Field of view" is limited to the cell containing the query; algorithm unaware of points in adjacent cells.

  - As dimensionality increases, surface area grows faster than volume $\Rightarrow$ points likely to be near cell boundaries.

  - Choosing good partitioning is non-trivial. Once chosen, cannot adapt to changes in data density.

## Algorithm



**1** Project data points and query along random directions.

**2** Add the closest point to the query along each projection direction to the frontier.

**3** Visit the point on the frontier with the shortest projected distance to the query.

**4** Add the next closest point along the most recently processed direction to the frontier.

**5** Visit the point on the frontier with the shortest projected distance to the query.

**6** Points highlighted in dark orange have been visited; visit the next point on the frontier.

**7** Current point has now been visited along all directions and is added to candidate set.

**8** Search exhaustively over all points in the candidate set and return the *k* closest ones.
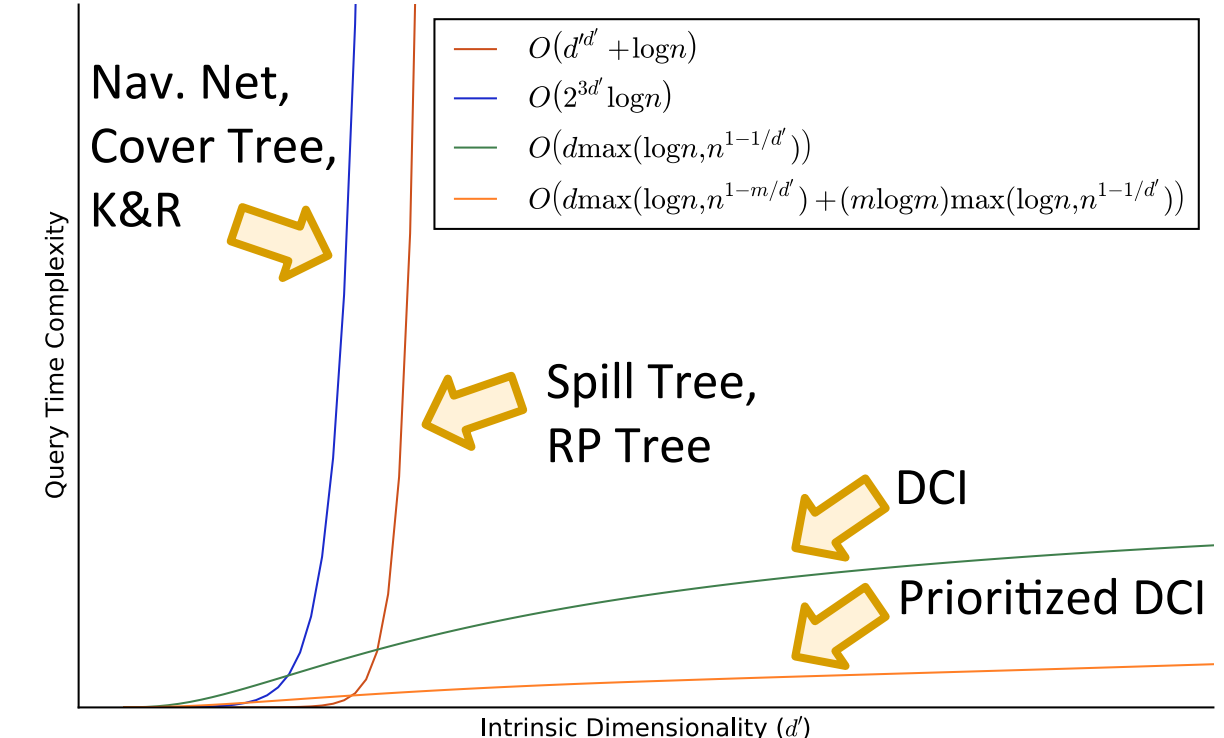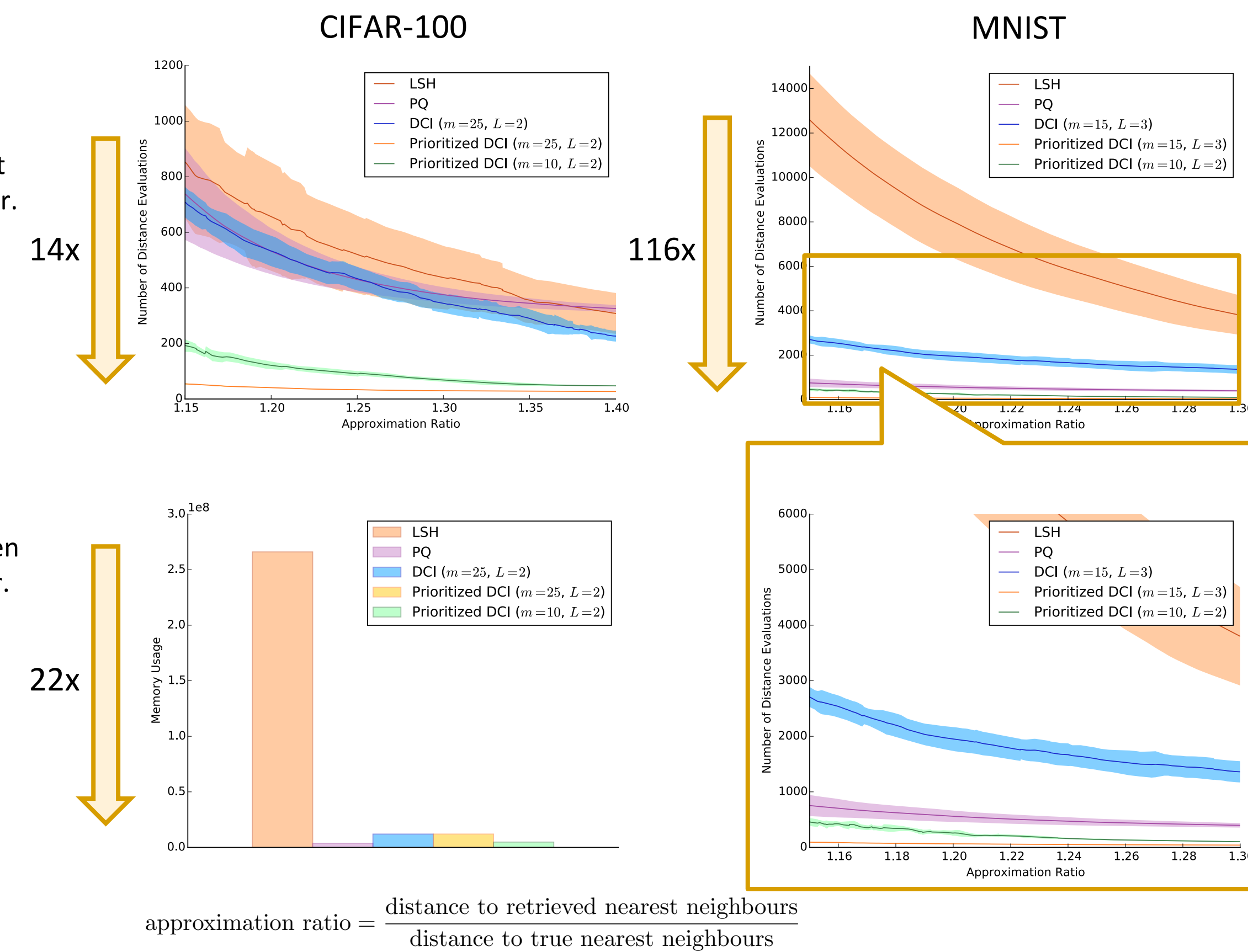
## Complexity

- Key Observation:

  - Points are added to the candidate set in the order of maximum projected distance to the query.

  - Maximum projected distance is a lower bound of the true distance.

  - As # of projection directions $\to \infty$, this $\to$ true distance.



| Construction | $O(m(dn + n \log n))$ |
|---|---|
| Query | $O(dk \max(\log(n/k), (n/k)^{1-m/d'}) + mk \log m(\max(\log(n/k), (n/k)^{1-1/d'})))$ |
| Insertion | $O(m(d + \log n))$ |
| Deletion | $O(m \log n)$ |
| Space | $O(mn)$ |

where $m \geq 1$ is # of projection directions

## Experiments



CIFAR-100     MNIST

14x     116x

22x

$$\text{approximation ratio} = \frac{\text{distance to retrieved nearest neighbours}}{\text{distance to true nearest neighbours}}$$