

Outline

- ✓ A. Intro
- ✓ B. Basic proximity problems
- ✓ C. Basic decision and summation problems
- ➡ D. Basic 'all'-type problems
- E. Basic counting problems
- F. Further topics

'All'-type problems

- **Nearest-neighbor search** $NN(x_q) = \arg \min_r \|x_q - x_r\|$
 All-nearest neighbor (bichromatic): $\forall x_q : NN(x_q) = \arg \min_r \|x_q - x_r\|$
- **Kernel density estimation** $\hat{f}(x_q) = \frac{1}{N} \sum_{r \neq q} K_h(\|x_q - x_r\|)$
 'All' version (bichromatic): $\forall x_q : \hat{f}(x_q) = \frac{1}{N} \sum_{r \neq q} K_h(\|x_q - x_r\|)$

Almost always 'all'-type problems

- Kernel density estimation
- Nadaraya-Watson & locally-wgtd regression
- Gaussian process prediction
- Radial basis function networks

Always 'all'-type problems

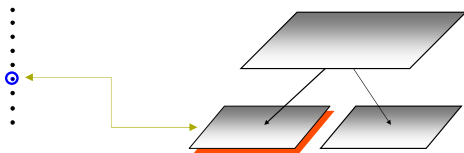
- Bichromatic all-nearest neighbor (e.g. LLE)
- Correlation dimension (pairs)
- n -point correlation (n -tuples)

Dual-tree idea

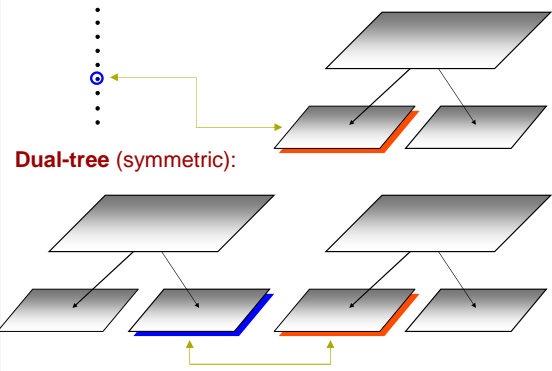
If all the queries are available simultaneously, then it is faster to:

1. Build a tree on the queries as well
2. Effectively process the queries in chunks rather than individually \rightarrow *work is shared between similar query points*

Single-tree:

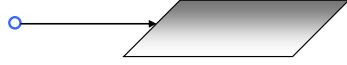


Single-tree:



**Exclusion and inclusion,
using point-node kd-tree bounds.**

O(D) bounds on distance minima/maxima:

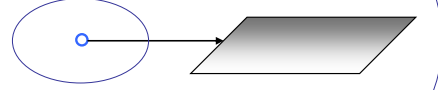


$$\min_i \|x - x_i\| \geq \sum_d^D [\max\{(l_d - x_d)^2, 0\} + \max\{(x_d - u_d)^2, 0\}]$$

$$\max_i \|x - x_i\| \leq \sum_d^D \max\{(u_d - x_d)^2, (x_d - l_d)^2\}$$

**Exclusion and inclusion,
using point-node kd-tree bounds.**

O(D) bounds on distance minima/maxima:

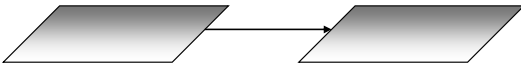


$$\min_i \|x - x_i\| \geq \sum_d^D [\max\{(l_d - x_d)^2, 0\} + \max\{(x_d - u_d)^2, 0\}]$$

$$\max_i \|x - x_i\| \leq \sum_d^D \max\{(u_d - x_d)^2, (x_d - l_d)^2\}$$

**Exclusion and inclusion,
using kd-tree node-node bounds.**

O(D) bounds on distance minima/maxima:



(Analogous to point-node bounds.)

Also needed:
Nodewise bounds.

**Exclusion and inclusion,
using kd-tree node-node bounds.**

O(D) bounds on distance minima/maxima:



(Analogous to point-node bounds.)

Also needed:
Nodewise bounds.

**Single-tree: simple recursive algorithm
(k=1 case)**

```

NN(xq, R, dlo, xsofar, dsofar)
{
  if dlo > dsofar, return.

  if leaf(R), [xsofar, dsofar] = NNBase(xq, R, dsofar).
  else,
    [R1, d1, R2, d2] = orderByDist(xq, R.l, R.r).
    NN(xq, R1, d1, xsofar, dsofar).
    NN(xq, R2, d2, xsofar, dsofar).
}

```

Single-tree → Dual-tree

- $x_q \rightarrow Q$
- $d_{lo}(x_q, R) \rightarrow d_{lo}(Q, R)$
- $x_{sofar} \rightarrow \underline{x}_{sofar}, d_{sofar} \rightarrow \underline{d}_{sofar}$
- store $Q.d_{sofar} = \max_Q \underline{d}_{sofar}$
- 2-way recursion → 4-way recursion
- $N \times O(\log N) \rightarrow O(N)$

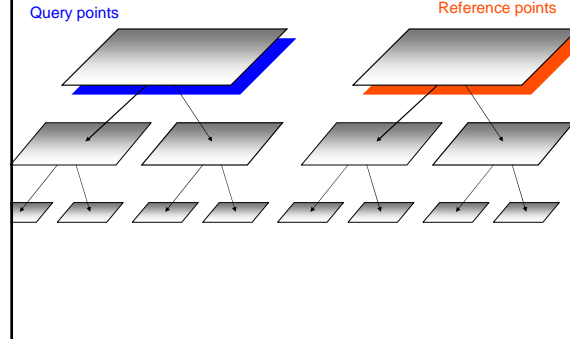
Dual-tree: simple recursive algorithm (k=1)

```

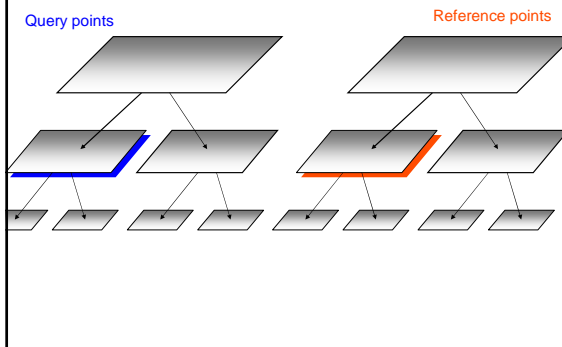
AIINN(Q,R,dlo,xsofar,dsofar)
{
  if dlo > Q.dsofar return.
  if leaf(Q) & leaf(R),
    [xsofar,dsofar]=AIINNBase(Q,R,dsofar). Q.dsofar=maxQdsofar
  else if !leaf(Q) & leaf(R), ...
  else if leaf(Q) & !leaf(R), ...
  else if !leaf(Q) & !leaf(R),
    [R1,d1,R2,d2]=orderByDist(Q.l,R.l,R.r).
    AIINN(Q.l,R1,d1,xsofar,dsofar).
    AIINN(Q.l,R2,d2,xsofar,dsofar).
    [R1,d1,R2,d2]=orderByDist(Q.r,R.l,R.r).
    AIINN(Q.r,R1,d1,xsofar,dsofar).
    AIINN(Q.r,R2,d2,xsofar,dsofar).
    Q.dsofar = max(Q.l.dsofar,Q.r.dsofar).
}

```

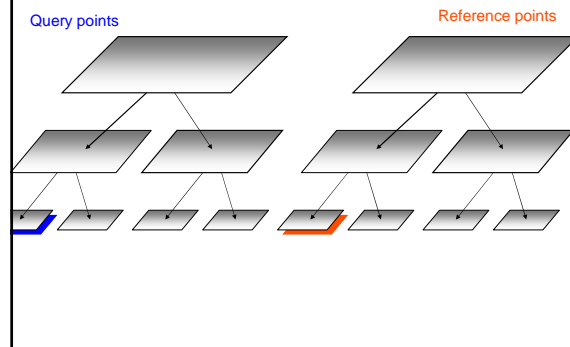
Dual-tree traversal (depth-first)



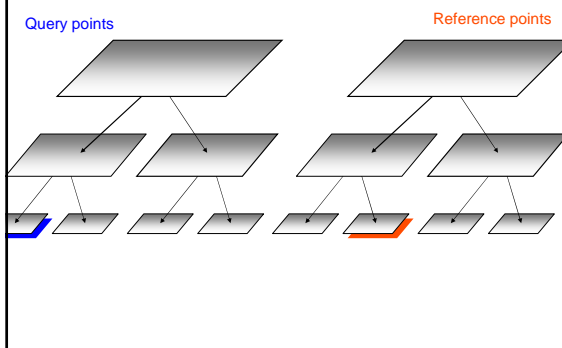
Dual-tree traversal



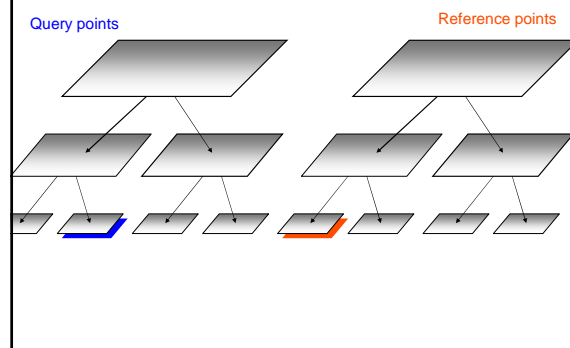
Dual-tree traversal

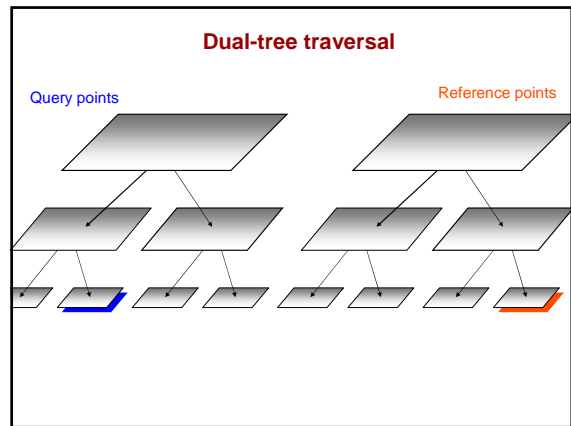
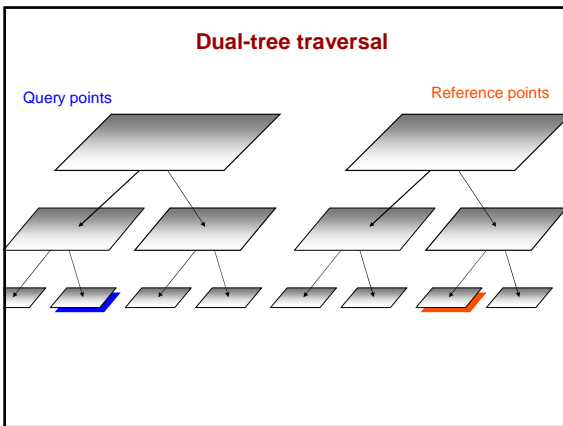
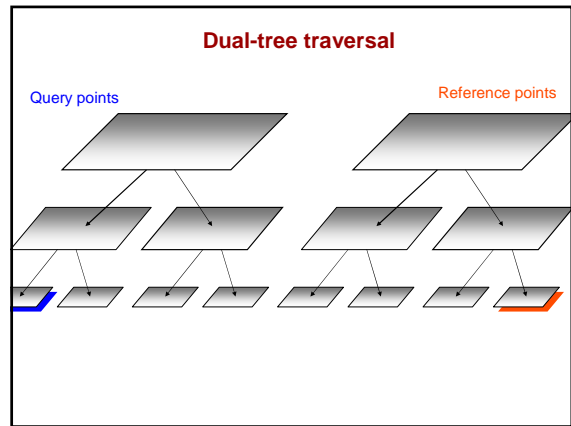
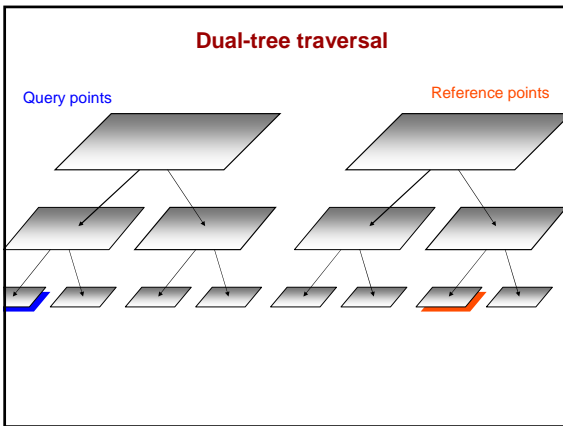
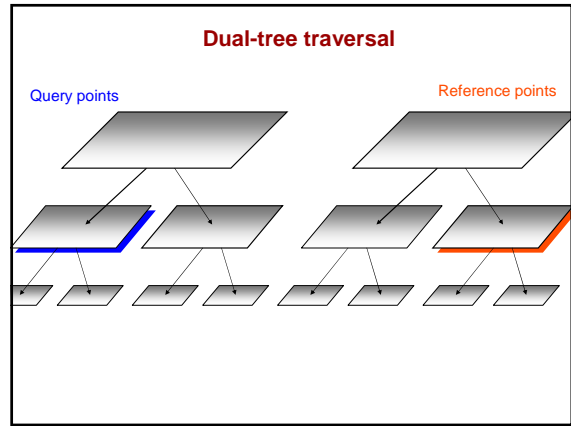
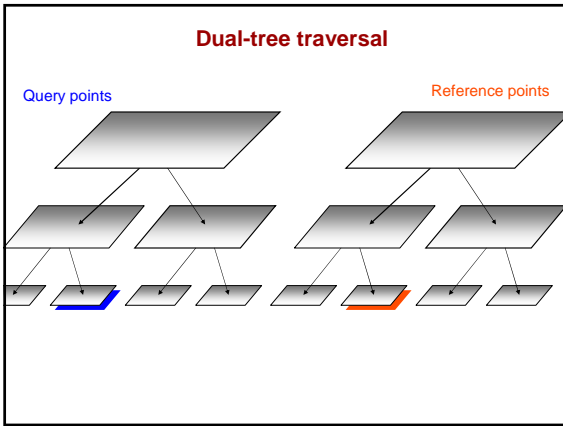


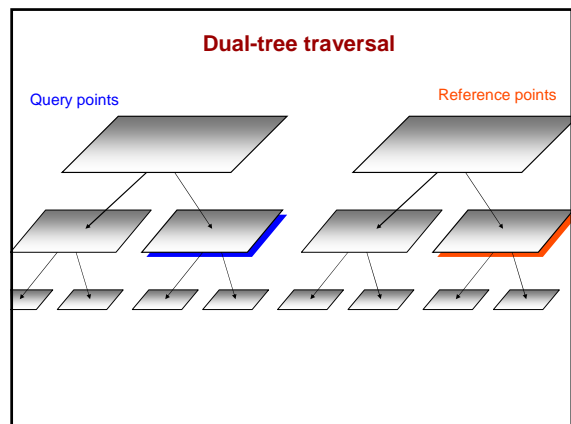
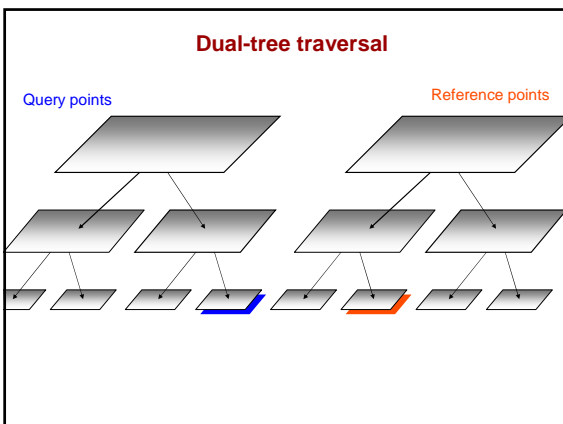
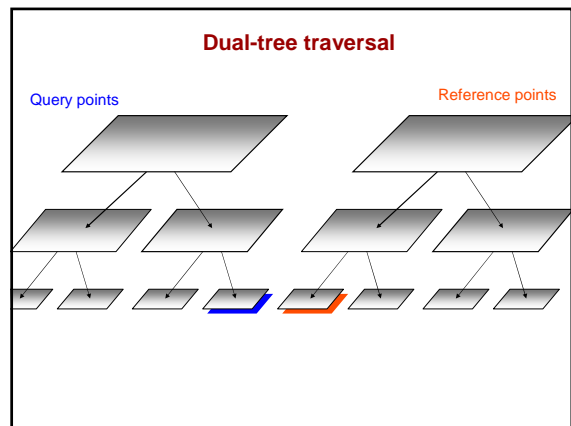
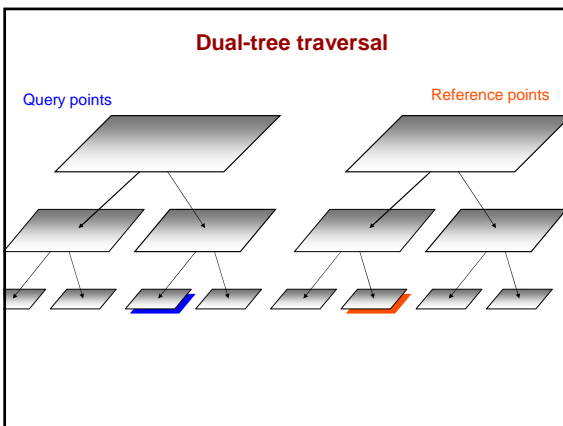
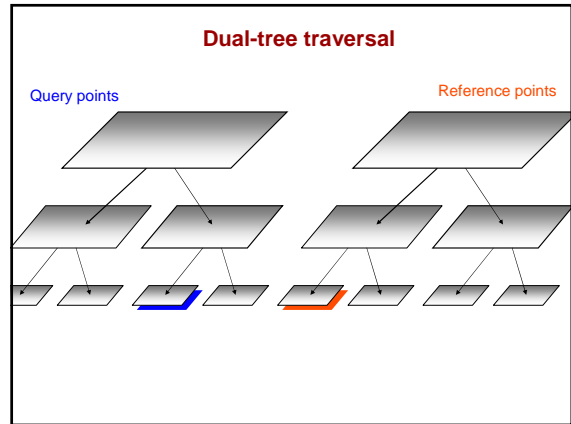
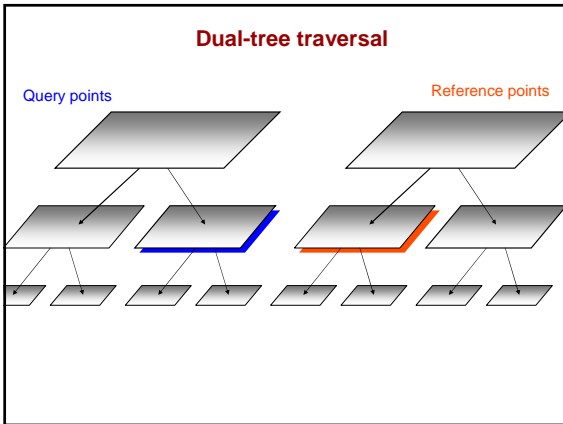
Dual-tree traversal

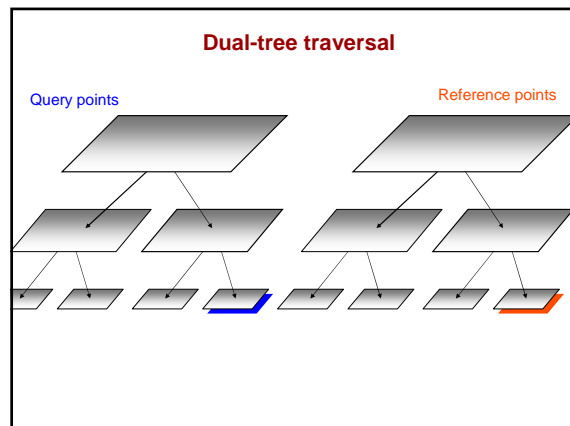
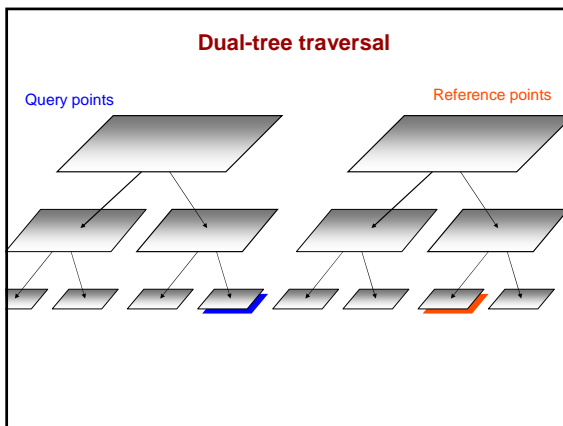
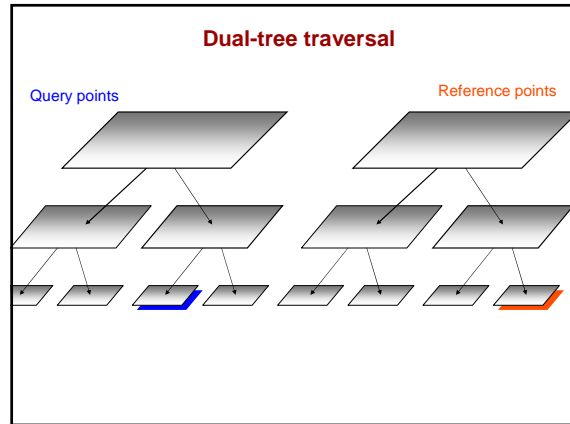
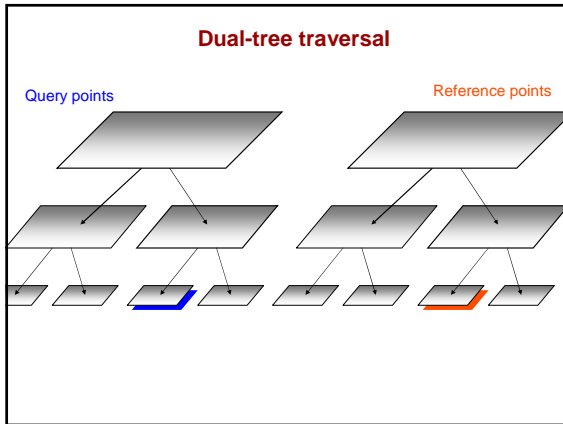


Dual-tree traversal









**Meta-idea:
Higher-order
Divide-and-conquer**

**Generalizes divide-and-conquer of a single set
to divide-and-conquer of multiple sets.**

Break each set into pieces.

Solving the sub-parts of the problem and
combining these sub-solutions appropriately
might be easier than doing this over only one set.

[Gray thesis, 2003]

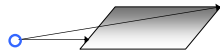
**Kernel density estimation
(All-queries version)**

$$\forall x_q : \hat{f}(x_q) = \frac{1}{N} \sum_{r \neq q}^N K_h(\|x_q - x_r\|)$$

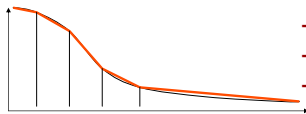
→ Now we need to do function approximation

Finite-difference function approximation.

$$f(x) = f(x_i) + \frac{1}{2} \left(\frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i} \right) (x - x_i)$$



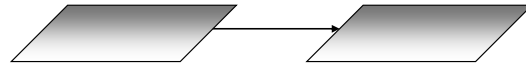
$$K_h(\delta) = K_h(\delta_{lo}) + \frac{1}{2} \left(\frac{K_h(\delta_{hi}) - K_h(\delta_{lo})}{\delta_{hi} - \delta_{lo}} \right) (\delta - \delta_{lo})$$



- stopping rule
- no tweak parameters
- hard error bounds

Approximation within dual-tree.

O(D) bounds on kernel sum contribution of R to Q:



Use distance bounds to obtain

$$K_h(\delta_{hi}), K_h(\delta_{lo})$$

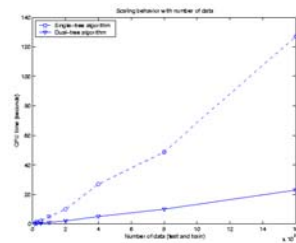
Also needed:

Nodewise bounds.

Speedup Results (KDE)

N	dual-naïve tree	
12.5K	7	.12
25K	31	.31
50K	123	.46
100K	494	1.0
200K	1976*	2
400K	7904*	5
800K	31616*	10
1.6M	35 hrs	23

5500x



One order-of-magnitude speedup over single-tree at ~2M points

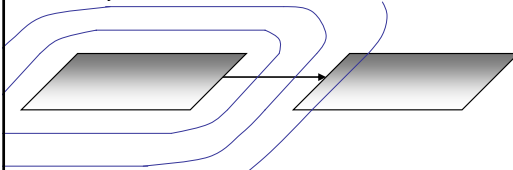
Speedup Results

N	Epan.	Gauss.
12.5K	.12	.32
25K	.31	.70
50K	.46	1.1
100K	1.0	2
200K	2	5
400K	5	11
800K	10	22
1.6M	23	51

Only about 2 - 2.5 times less efficient

Exclusion and inclusion, on multiple radii simultaneously.

Use binary search to locate critical radius:



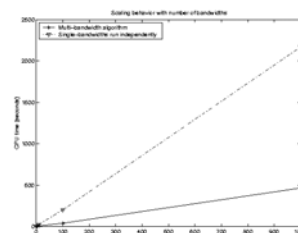
$$\min ||x - x_i|| < h_1 \Rightarrow \min ||x - x_i|| < h_2$$

Also needed:

b_lo, b_hi are arguments; store bounds for each b

Application of HODC principle

Speedup Results



One order-of-magnitude speedup over single-radius at ~10,000 radii

These are all examples of...

Generalized N-body problems

All-NN: $\{\forall, \arg \min, \delta, \cdot\}$

2-point: $\{\Sigma, \Sigma, I_h(\delta), w\}$

3-point: $\{\Sigma, \Sigma, \Sigma, I_H(\delta), w\}$

KDE: $\{\forall, \Sigma, K_h(\delta), \cdot; \{h\}\}$

SPH: $\{\forall, \Sigma, K_h(\delta), w; t\}$

General theory and toolkit for
designing algorithms for
such problems