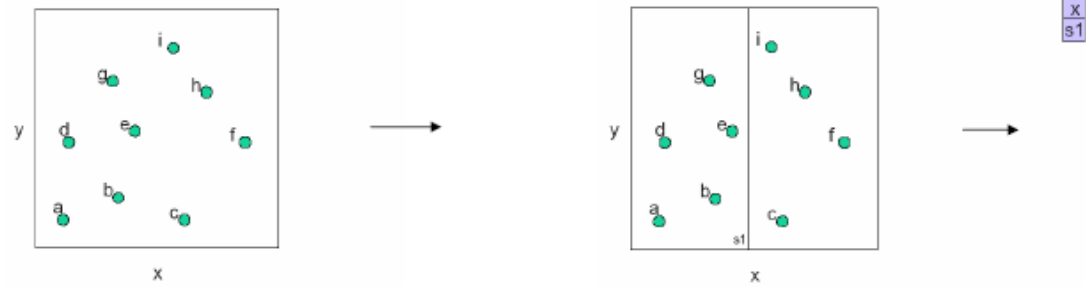


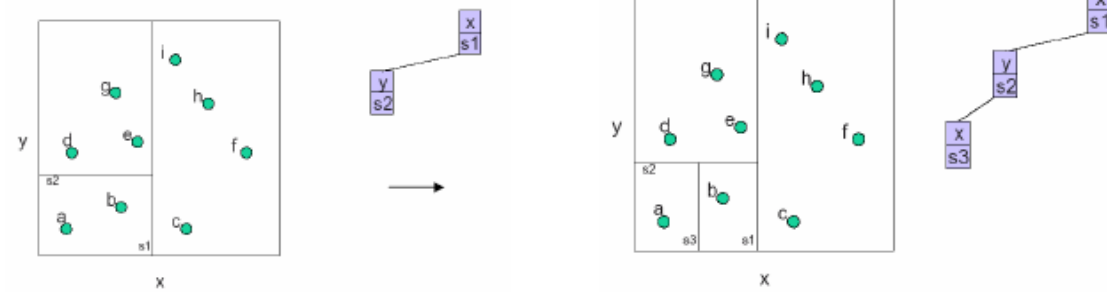
# Parallel Machine Learning

- k-NN used for classification and regression problems
- Commonly used data structure: k-d trees
- For the given multi-dimensional data, construct a k-d tree
- Similar to

# k-d tree construction example

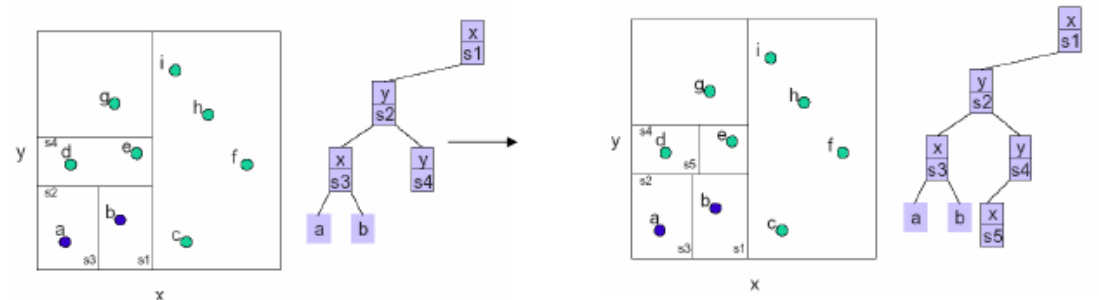
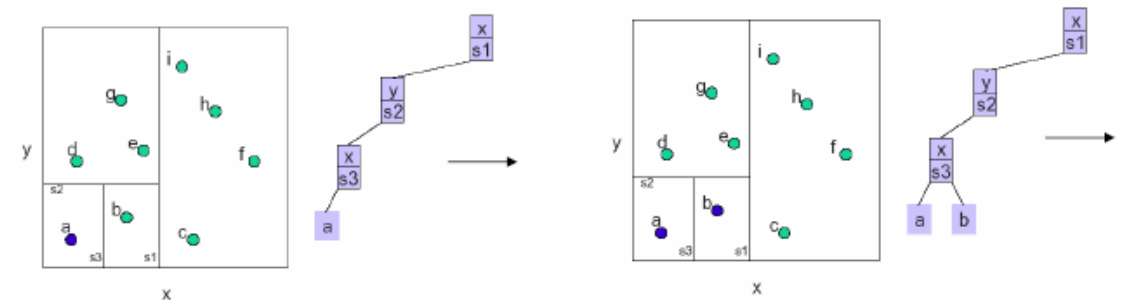


divide perpendicular to the widest spread.



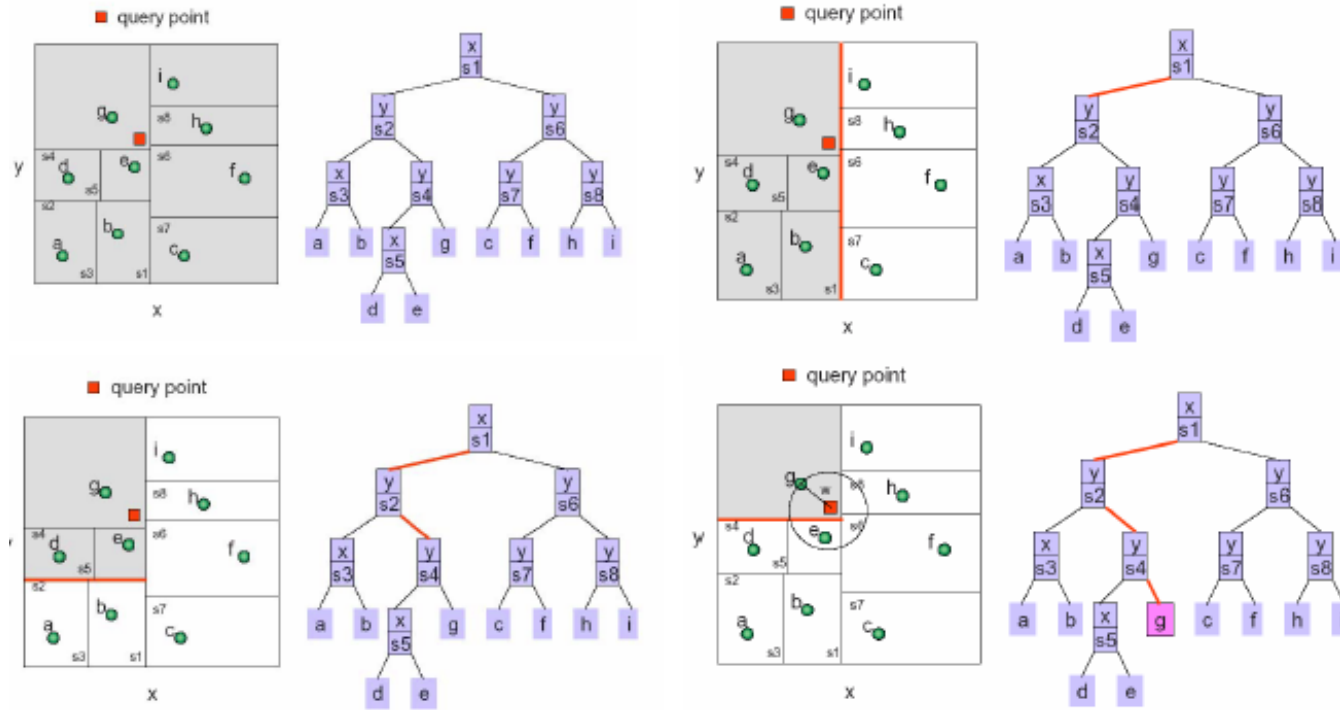
# k-d Tree Construction

## k-d tree construction example



# Nearest neighbour search using kd- tree

## k-d Tree Nearest Neighbor Search



# Two steps in k-NN

1. k-d tree construction
2. K-d tree search

## Parallelization

- Option 1:
  - Tree construction: Partition data sets among processors; Each processor constructs local k-d tree
  - Search: Query sent to all processors which perform search in local k-d trees; each processor returns the top neighbors from which the k nearest are chosen

# Parallelization

- Option 1: Poor work efficiency, i.e., wasted work
- Option 2:
  - Global k-d tree construction in all the processors
  - For each processors, one half of data is given to one half of processors, and the other half of data given to other processors
  - After this recursive division, the top part of the tree is replicated in all the processors
  - The processors then construct local k-d trees for their subdomains

# Searching the global k-d tree

- Query sent to the processors that takes care of the subdomain of the query
- The processor forms the local k nearest neighbors
- Forms a radius based on these neighbors
- Sends the query to the nearby processors consisting of subdomains that are spanned by the radius
- The processors search for points within the radius and send their results to the origin processor

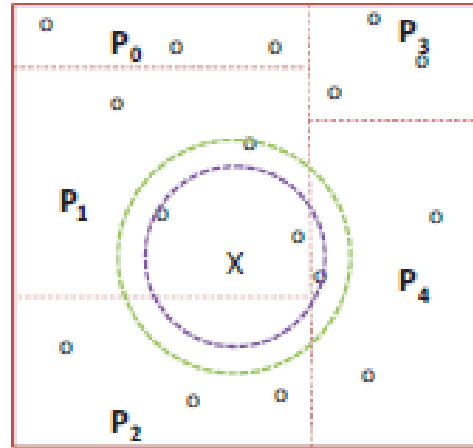


Figure 3. Figure shows the data points (denoted by  $\circ$ ) in 2D space divided among 5 nodes. Query point is shown as  $X$ . KNN with  $k = 3$  is run in node  $P_1$  (owner of  $X$ ). This returns 3 points owned by  $P_1$  and a max distance (denoted by green circle around  $X$ ). Only  $P_2$  and  $P_4$  might own points within this radius. KNN is run on  $P_2$  and  $P_4$  for  $X$  and the closest 3 points are chosen (within purple circle around  $X$ ).



# Optimization

- Queries can be batched
- Software pipelining can be performed between the stages

# Algorithm

**Algorithm 1** Finding  $k$ -nearest neighbors from the local kd-tree. Input: kd-tree  $T$ , Query  $q$ ,  $k$ , search radius,  $r$  (default  $r = \infty$ ). Output: A set,  $R$  of  $k$  nearest neighbors within  $r$ .

```
1: procedure FINDKNN( $T, q, k, r$ )
2:    $r' \leftarrow r$ ; push ( $root, 0$ ) into  $S$ 
3:   while  $S$  is not empty do
4:     ( $node, d$ )  $\leftarrow$  pop from  $S$ 
5:     if  $node$  is leaf then
6:       for each particle  $x$  in  $node$  do
7:         compute distance,  $d[x]$  of  $x$  from  $q$ 
8:         if  $d[x] < r'$  then
9:           if  $|H| < k$  then
10:            add  $x$  into  $H$ 
11:            if  $|H| = k$  then
12:               $r' \leftarrow H.max\_dis$ 
13:            else if  $d[x] < \max$  distance in  $H$  then
14:              replace the topmost point  $H$  by  $x$ 
15:               $r' \leftarrow d[x]$ 
16:           else
```

# Algorithm

```
16:     else
17:         if  $d < r'$  then
18:              $d' \leftarrow q[\text{node.dim}] - \text{node.median}$ 
19:              $d' \leftarrow \sqrt{d * d + d' * d'}$ 
20:              $C_1 \leftarrow$  closer child of  $\text{node}$  from  $q$ 
21:              $C_2 \leftarrow$  other child of  $\text{node}$ 
22:             if  $d' < r'$  then
23:                 push  $(C_2, d')$  into  $S$ 
24:             push  $(C_1, d)$  into  $S$ 
25:      $R \leftarrow H$ 
```

---

# References

- PANDA: Extreme Scale Parallel K-Nearest Neighbor on Distributed Architectures. IPDPS 2016.