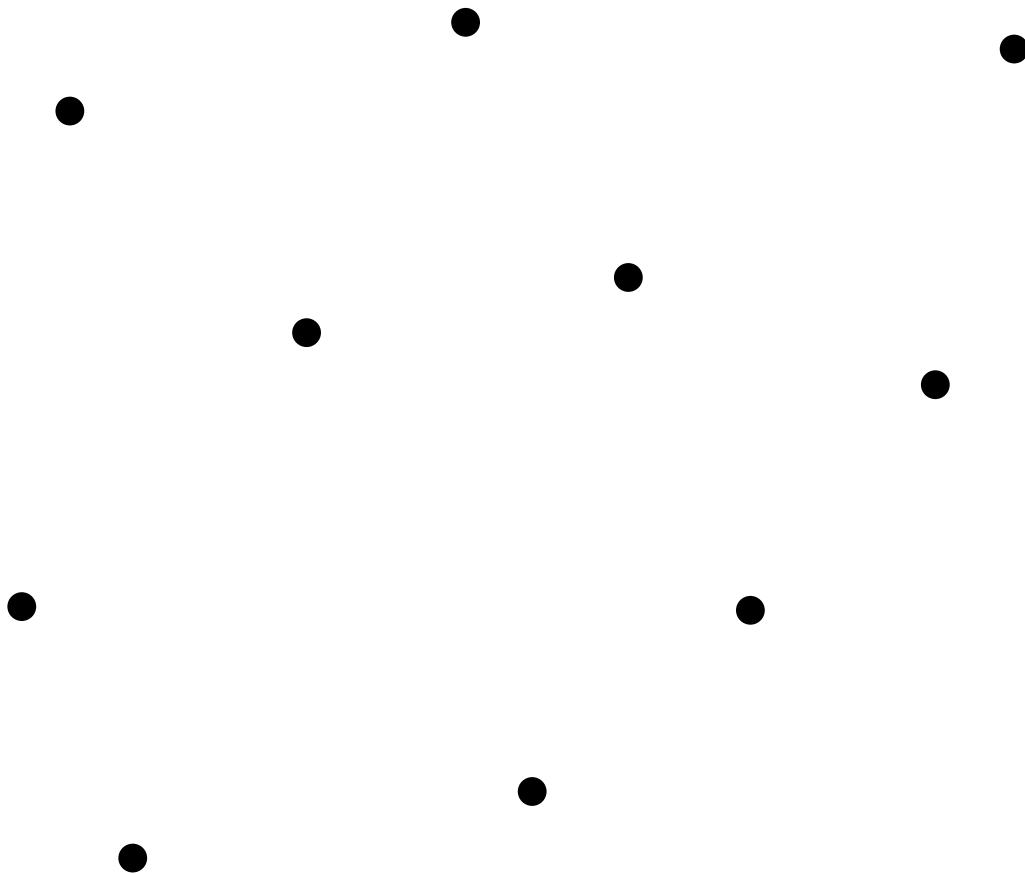


# Nearest Neighbor Search

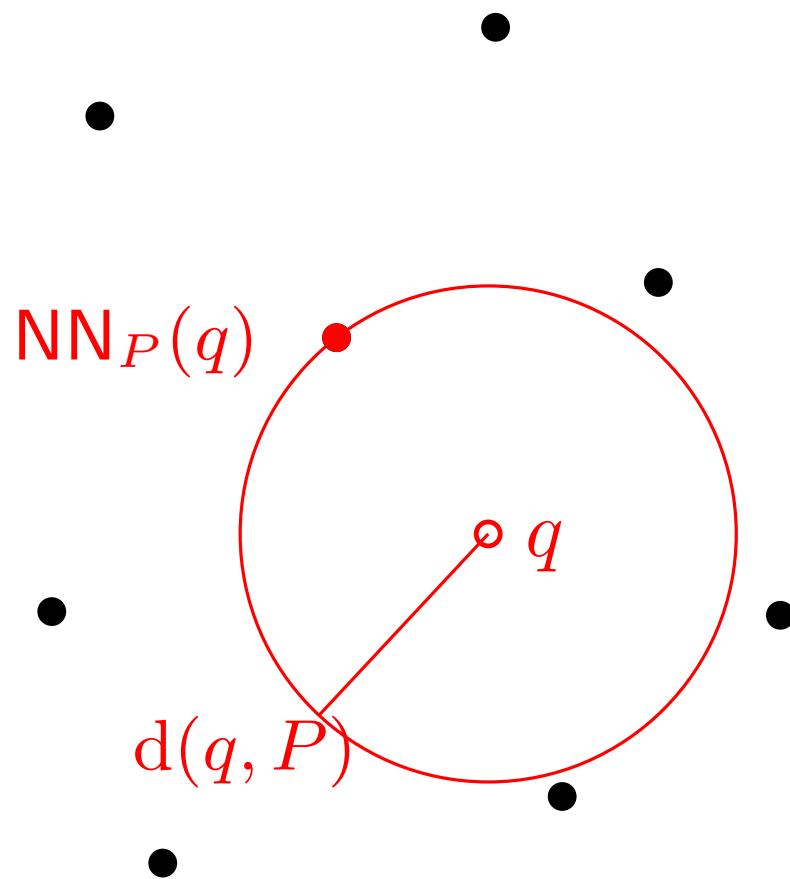
Steve Oudot

# Nearest-Neighbor problem

pre-processing input:  $P$



# Nearest-Neighbor problem



pre-processing input:  $P$

query input:  $q$

goal: find  $p \in NN_P(q)$

$$d(q, p) = \min_{p' \in P} d(q, p')$$

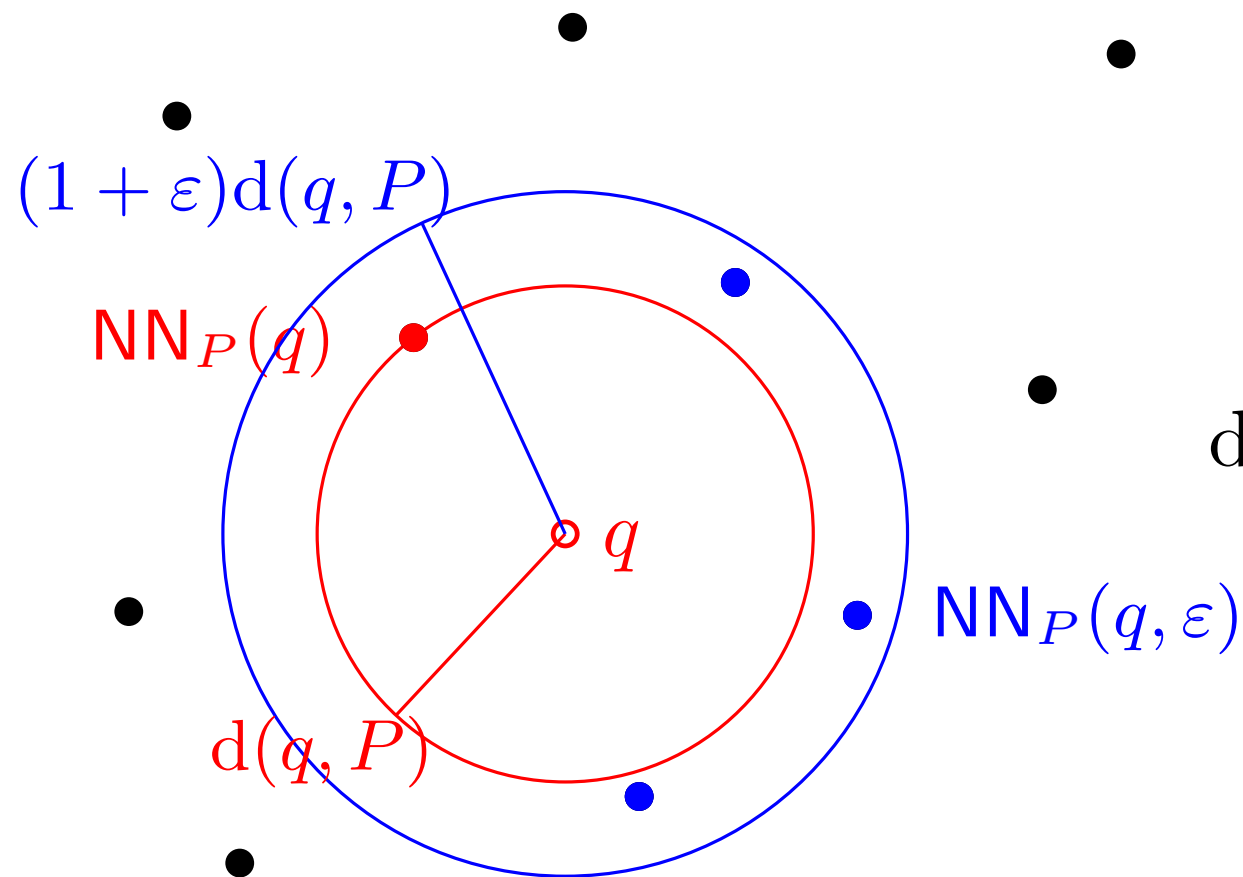
# $\varepsilon$ -Nearest-Neighbor problem

pre-processing input:  $P, \varepsilon$

query input:  $q$

goal: find  $p \in \text{NN}_P(q, \varepsilon)$

$$d(q, p) \leq (1 + \varepsilon) \min_{p' \in P} d(q, p')$$



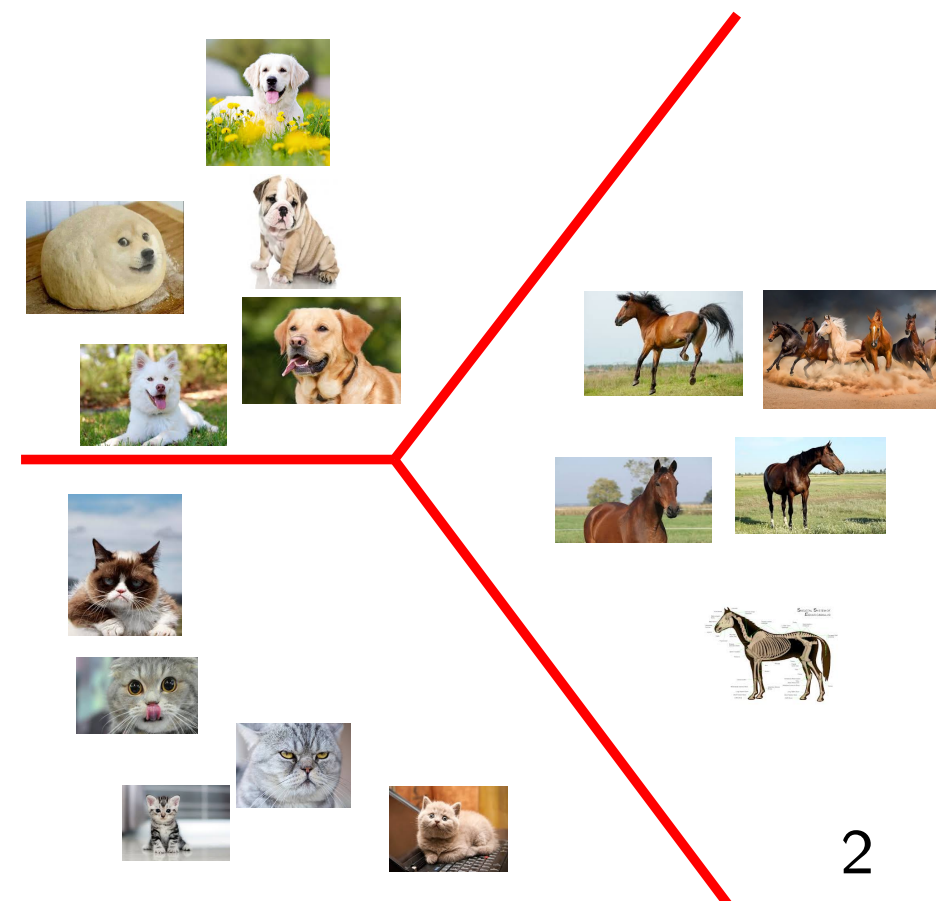
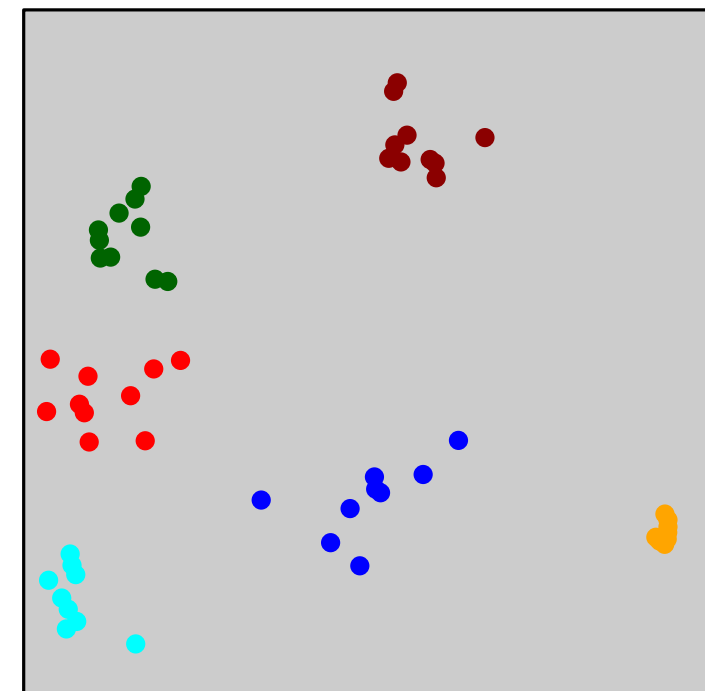
# Nearest-Neighbor problem

## Variants:

- $k$ -nearest neighbors: find the  $k$  points closest to  $q$  in  $P$
- $r$ -nearest neighbor: find a point  $p \in P$  such that  $d(q, p) \leq r$
- metrics:
  - ▶  $\ell_2, \ell_p, \ell_\infty$
  - ▶ strings: Hamming distance
  - ▶ images: optimal transport distances
  - ▶ point clouds: (Gromov-)Hausdorff distances
  - ▶ proteins: RMSD distances
  - ▶ ...

# Applications

- clustering, e.g. k-means, mean-shift
- information retrieval in databases
- information theory, e.g. vector quantization
- supervised learning, e.g. NN-classifiers
- . . .



# Strategy and Challenges

## Strategy:

- ▶ preprocess the  $n$  point of  $P$  in  $\mathbb{R}^d$  into some data structure DS for fast nearest-neighbor queries answers

## Ideal wish list:

- ▶ DS should have **linear size** in  $n$  and polynomial size in  $d$
- ▶ a query should take **sublinear time** in  $n$  and polynomial time in  $d$   
e.g. binary search trees in  $d = 1$ : linear size,  $O(\log n)$  time

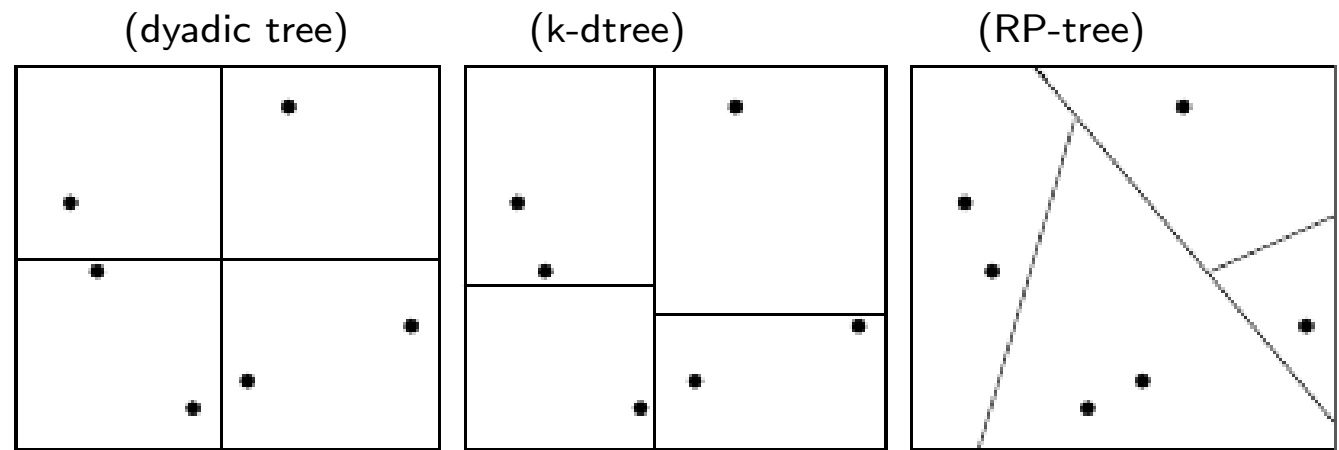
## Core difficulties:

- ▶ ***Curse of dimensionality***: hard to outperform linear scan in high  $d$
- ▶ Interpretation: meaningfulness of distances in high  $d$  (**concentration**)

$O(dn)$  space and time

# Approaches

- Linear scan
- Voronoi diagrams
- Tree-like data structures



- ▶ quadtrees (split at midpoint in all coordinates)
- ▶ tries / dyadic trees (split at mean, cycle around coordinates)
- ▶ kd-trees (split at median, cycle around coordinates)
- ▶ Random Projection trees (split at median along random coordinates)
- ▶ PCA trees (split at median along 1st eigenvector of covariance matrix)
- ▶ ...

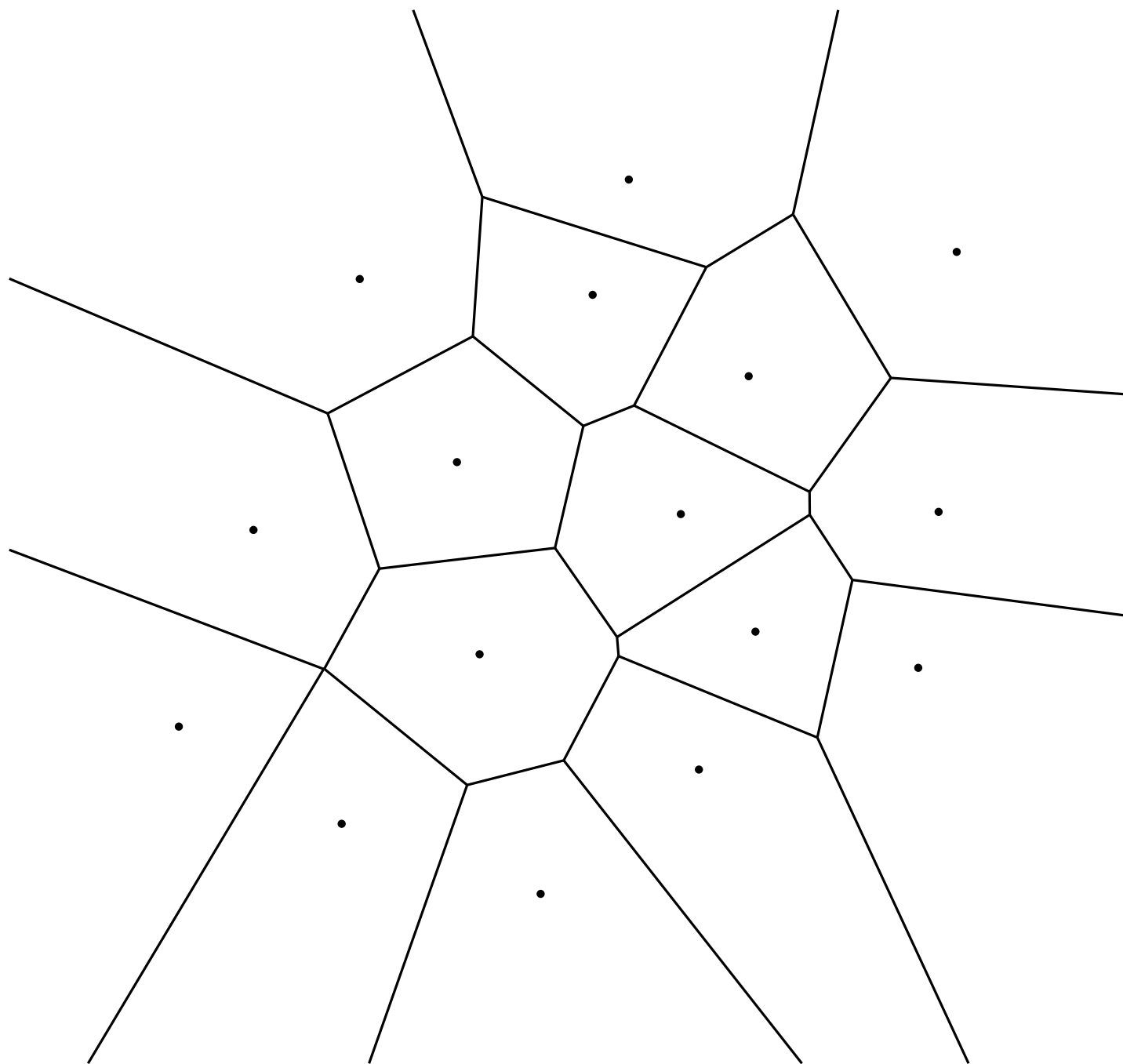
Binary  
Space  
Partitions

- Locality Sensitive Hashing



# Voronoi diagrams

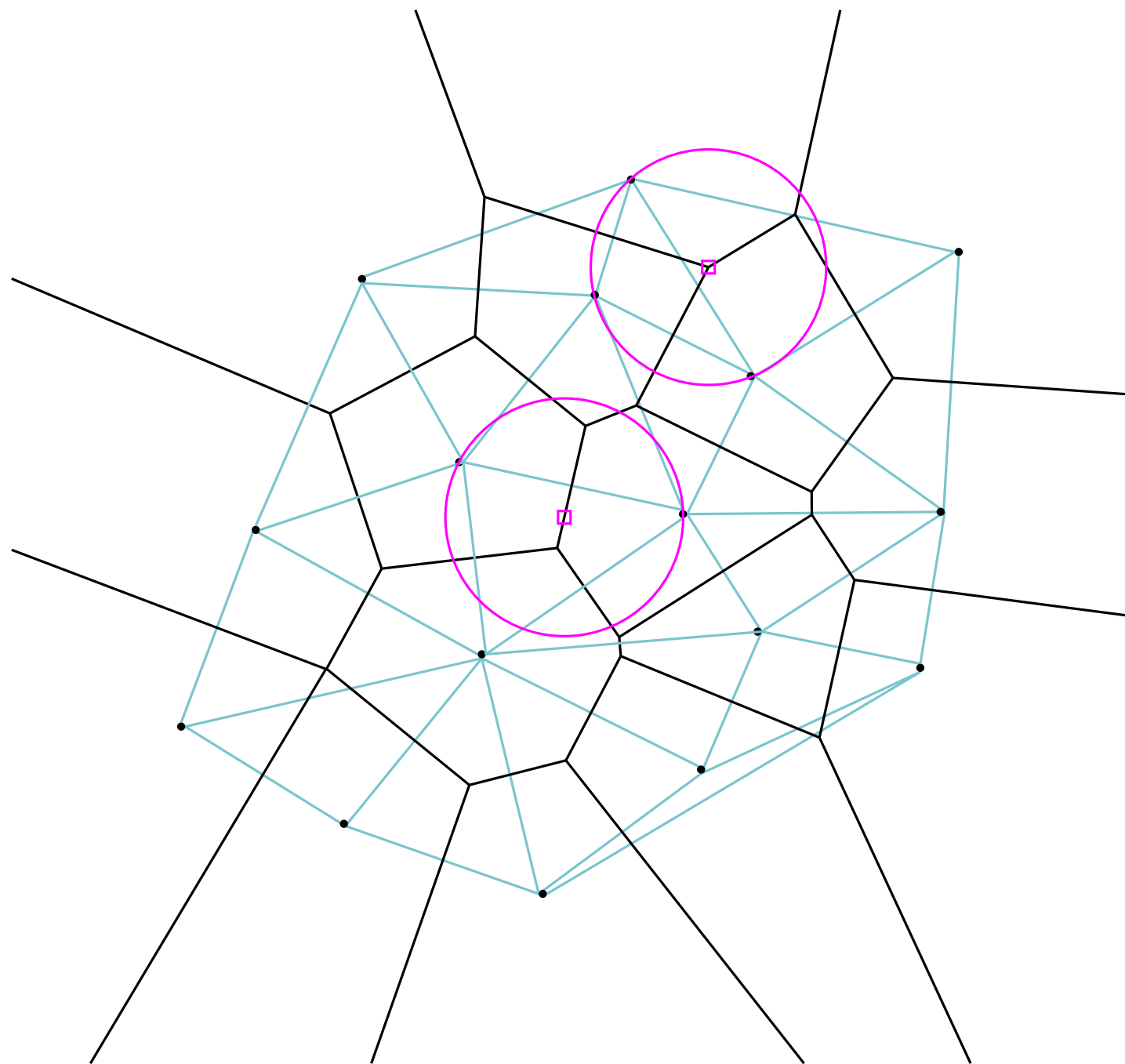
# Definition



$$V(p) := \{q \in \mathbb{R}^d \mid p \in \text{NN}_P(q)\}$$

affine diagram

# Definition



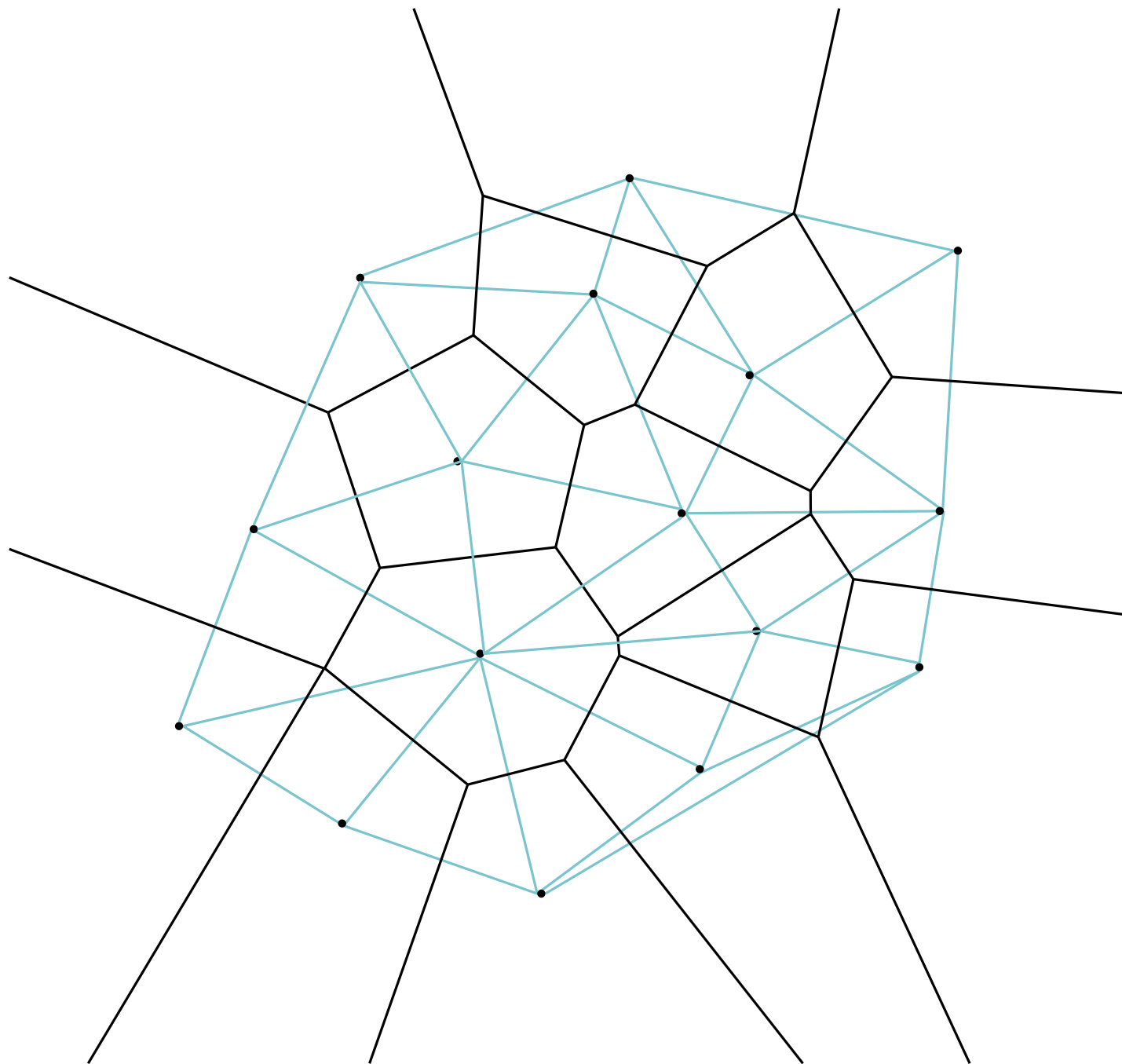
$$V(p) := \{q \in \mathbb{R}^d \mid p \in \text{NN}_P(q)\}$$

affine diagram

computed/stored via dual

(Delaunay triangulation)

# Definition



$$V(p) := \{q \in \mathbb{R}^d \mid p \in \text{NN}_P(q)\}$$

affine diagram

computed/stored via dual  
(Delaunay triangulation)

size:

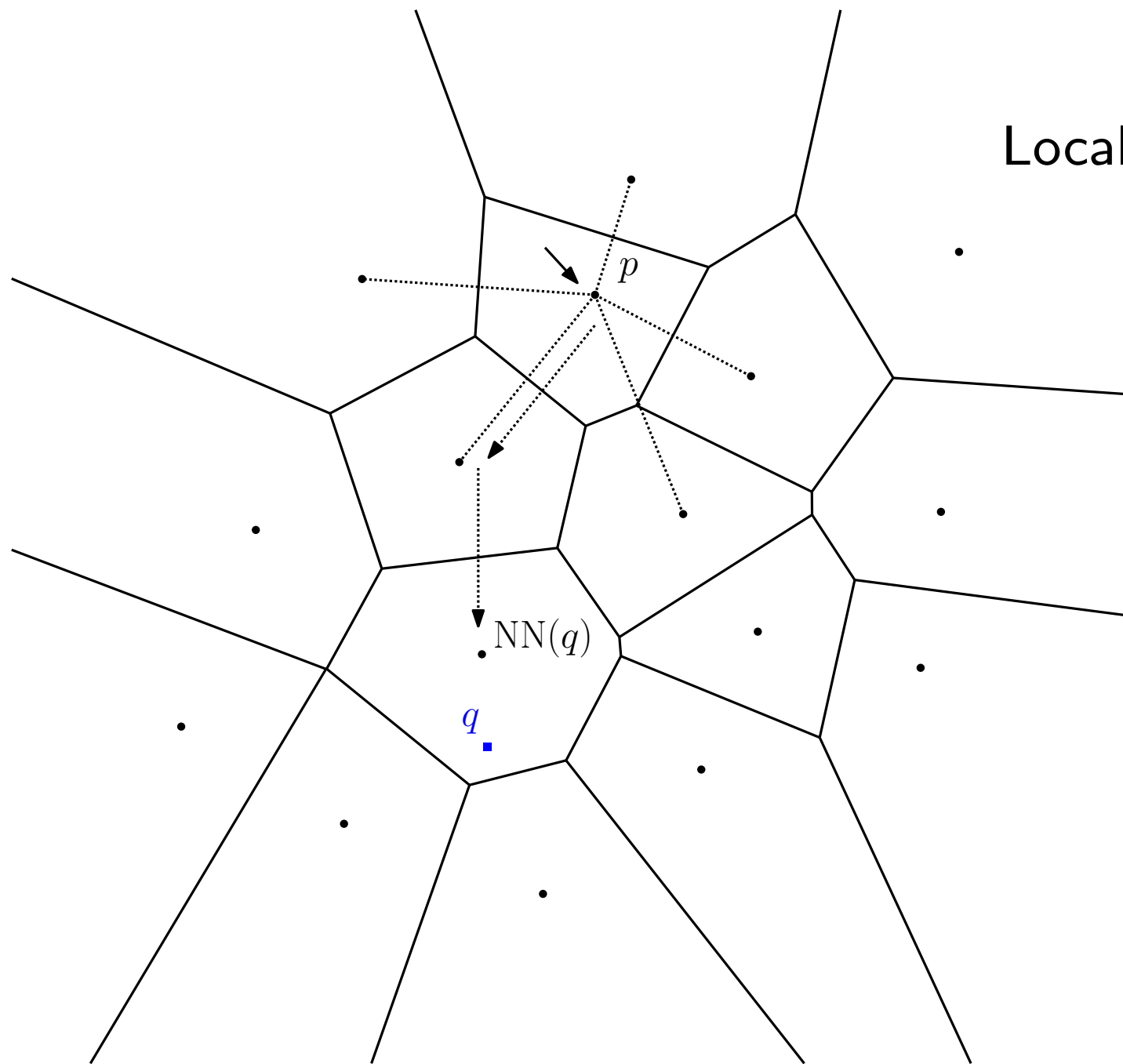
- worst case:  $\Theta\left(n^{\lceil d/2 \rceil}\right)$

Upper Bound Thm [McMullen'70]

- average case (unif. distrib.):

$$2^{O(d \log d)} n$$

# Usage for NN-search



Localizing by **walk**:

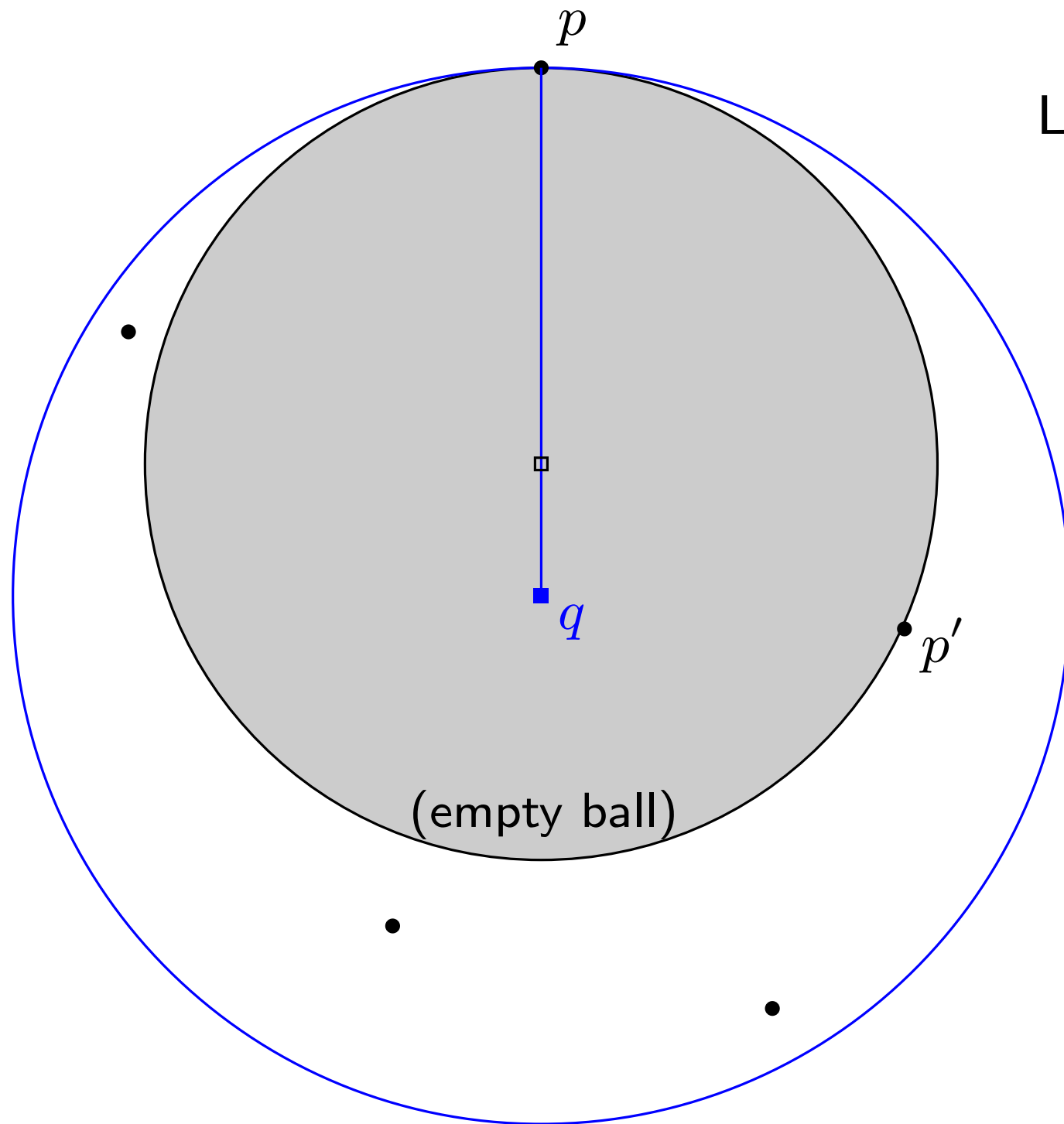
start from  $p \in P$  random

**while**  $\exists p'$  neighb. of  $p$  in Del.

s.t.  $d(q, p') < d(q, p)$ :

update  $p := p'$

# Usage for NN-search



Localizing by **walk**:

start from  $p \in P$  random

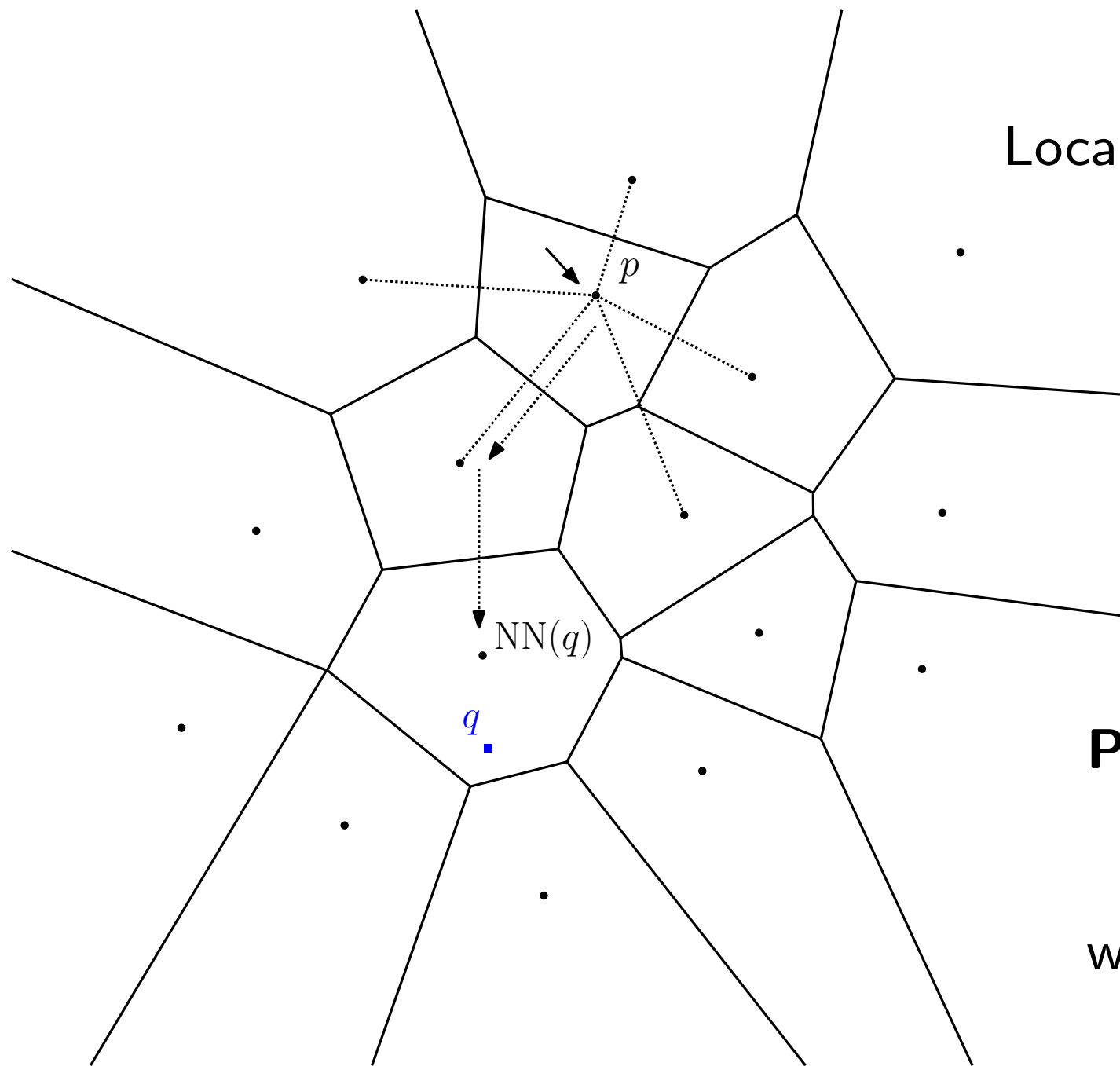
**while**  $\exists p'$  neighb. of  $p$  in Del.

s.t.  $d(q, p') < d(q, p)$ :

update  $p := p'$

**Prop:** Del. neighborhood is complete

# Usage for NN-search



Localizing by **walk**:

start from  $p \in P$  random

**while**  $\exists p'$  neighb. of  $p$  in Del.

s.t.  $d(q, p') < d(q, p)$ :

update  $p := p'$

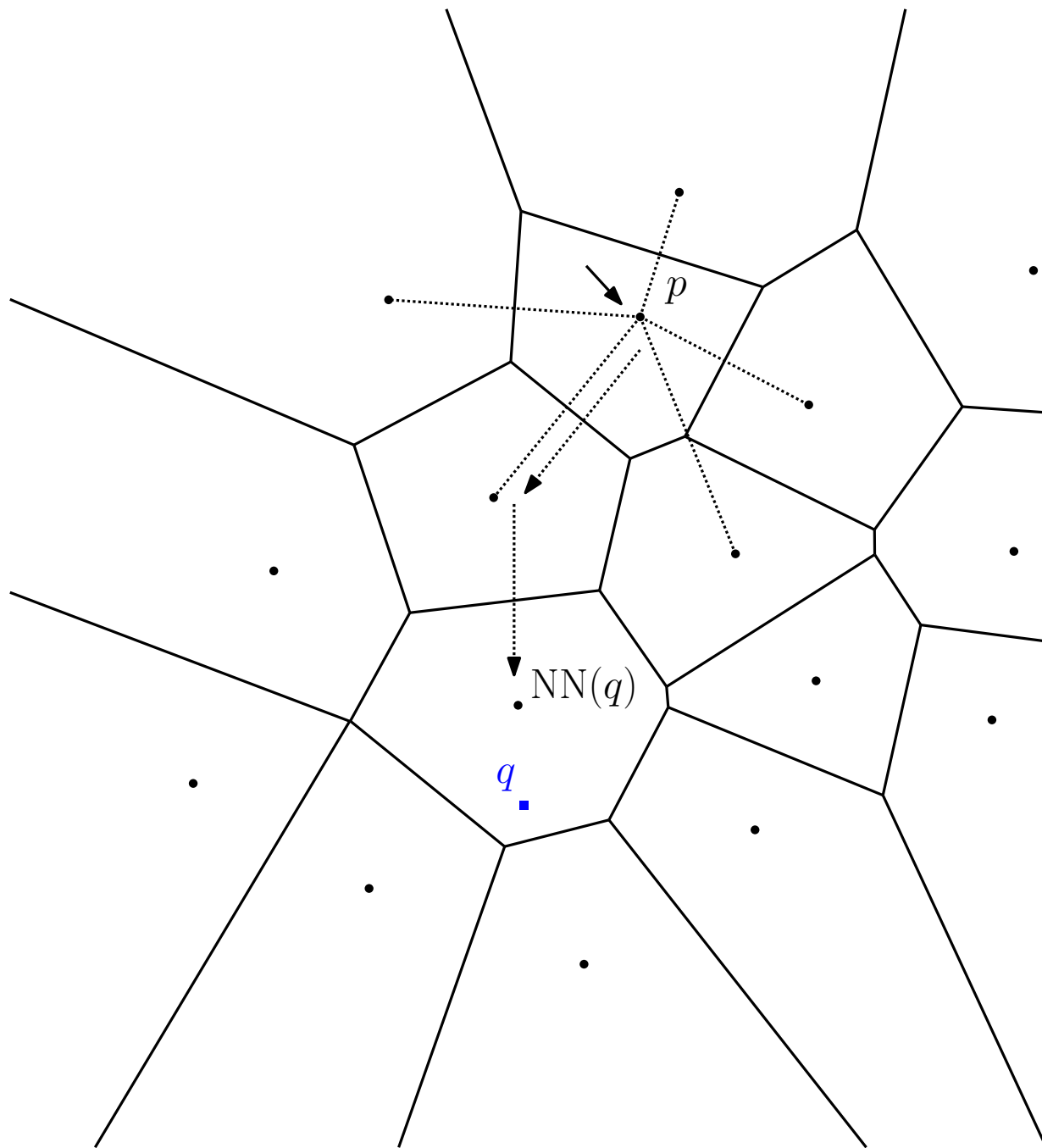
**Prop:** Del. neighborhood is complete

walk time:

worst case:  $O(|\text{Del}(P)|)$

avg. case (2d):  $O(\sqrt{n})$

# Usage for NN-search



Localizing by **hierarchy**:

- Voronoi subdivision [Kirk.'83, Meiser'93]:

(2D)  $O(n)$  space,  $O(\log n)$  time

(dD)  $\Theta(n^d)$  space,  $O(d^5 \log n)$  time

• Delaunay tree [Mulmuley'91]:

(2D)  $O(n \log n)$  space,  $O(\log n)$  time

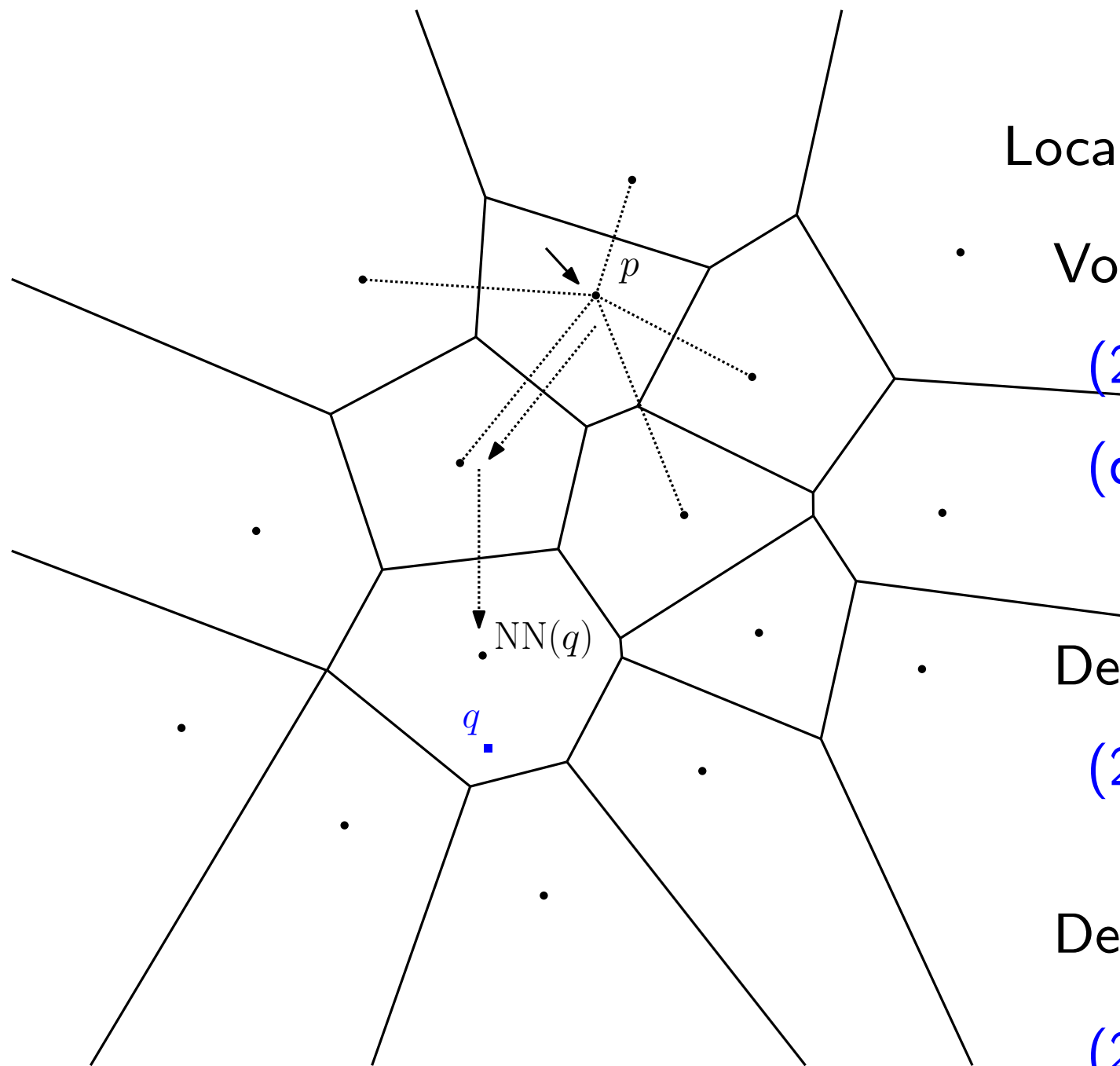
Delaunay tree + walk [Devillers'02]:

(2D)  $O(n \log n)$  space,  $O(\log n)$  time

(dD)  $O(n^{\lceil \frac{d}{2} \rceil})$  space,  $O(n^{\lceil \frac{d-2}{2} \rceil})$  time



# Usage for NN-search



Localizing by **hierarchy**:

- Voronoi subdivision [Kirk.'83, Meiser'93]:

(2D)  $O(n)$  space,  $O(\log n)$  time

(dD)  $\Theta(n^d)$  space,  $O(d^5 \log n)$  time

• Delaunay tree [Mulmuley'91]:

(2D)  $O(n \log n)$  space,  $O(\log n)$  time

Delaunay tree + walk [Devillers'02]:

(2D)  $O(n \log n)$  space,  $O(\log n)$  time

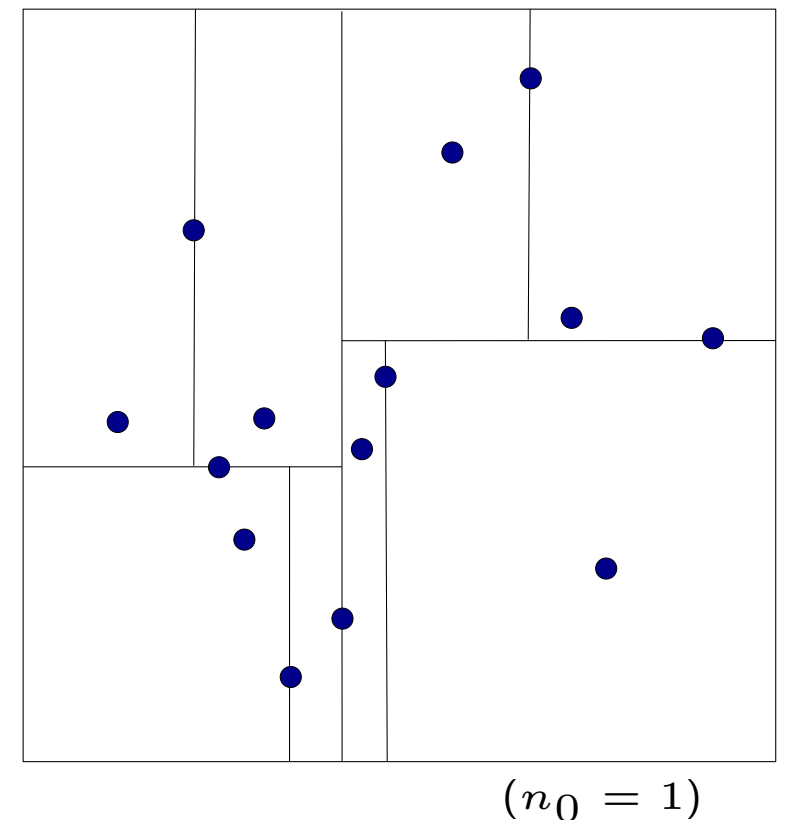
(dD)  $O(n^{\lceil \frac{d}{2} \rceil})$  space,  $O(n^{\lceil \frac{d-2}{2} \rceil})$  time

For small dimensions (2 or 3) only!

k-d trees

# Definition

- a binary tree
- each internal node implements a spatial partition induced by a hyperplane  $H$ , splitting the point cloud into two equal subsets
  - ▶ right subtree: all points lying on one side of  $H$
  - ▶ left subtree: remaining points
- subdivision stops whenever fewer than  $n_0$  remain  
     $\rightsquigarrow$  size:  $O(dn)$



# Definition

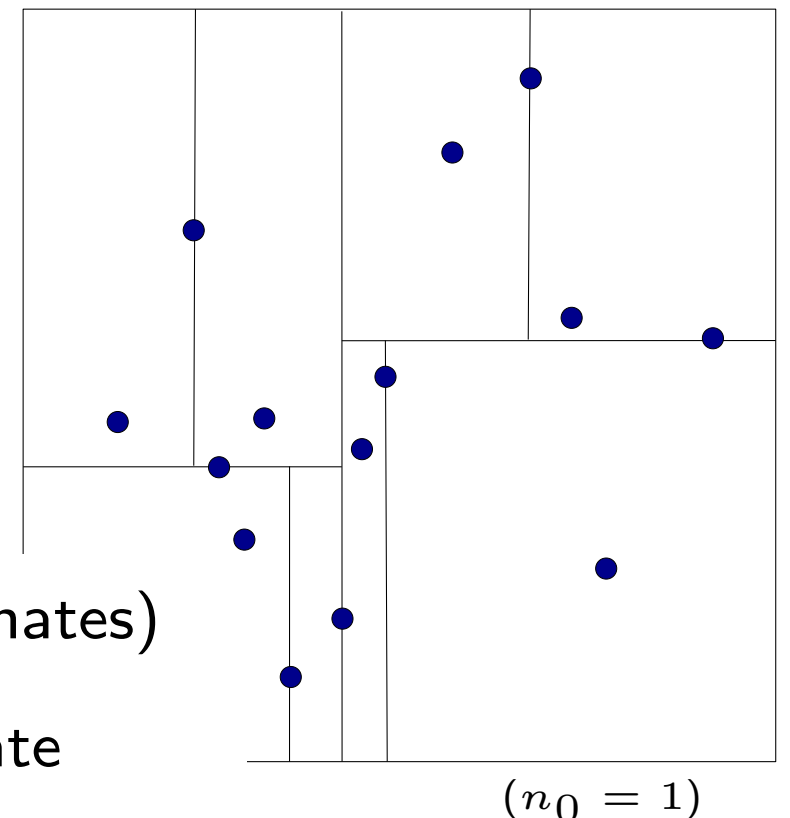
- a binary tree
- each internal node implements a spatial partition induced by a hyperplane  $H$ , splitting the point cloud into two equal subsets
  - ▶ right subtree: all points lying on one side of  $H$
  - ▶ left subtree: remaining points
- subdivision stops whenever fewer than  $n_0$  remain

↪ size:  $O(dn)$

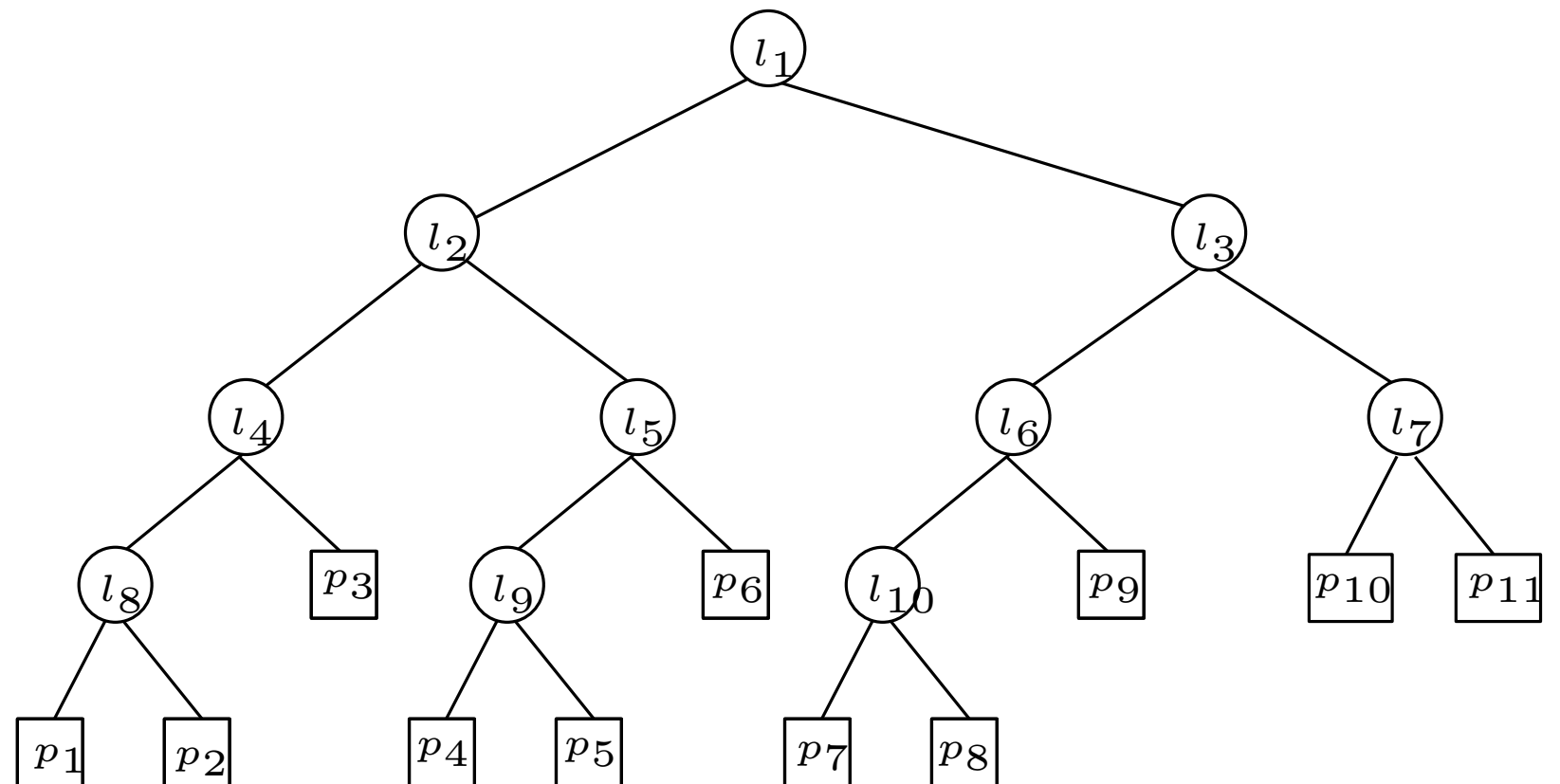
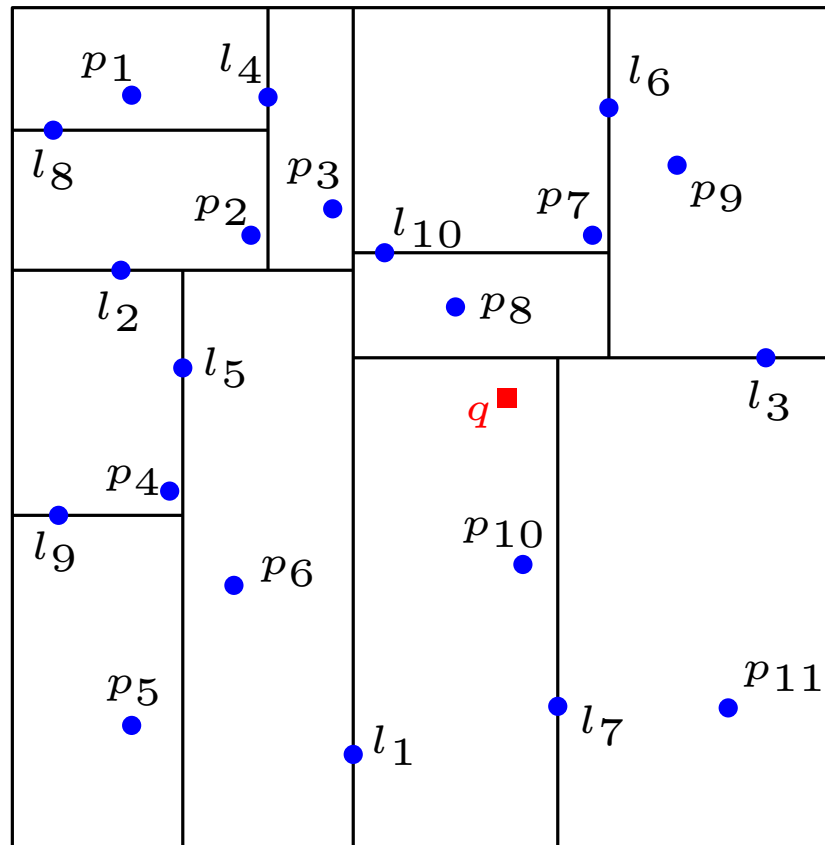
kd-tree specifics:

$H$  orthogonal to coordinate axis (cycle through coordinates)

$H$  goes through the median in the considered coordinate



# Example



$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Usage for NN search

## Strategy 1: defeatist search

$d_{\min} := \infty$  (dist. to pts viewed so far)

search (*node*): (*node* = *root* initially)

**if** *node* = *leaf*:

$d_{\min} := \min\{d_{\min}, \min_{p \in \text{node.batch}} d(q, p)\}$

**else:**

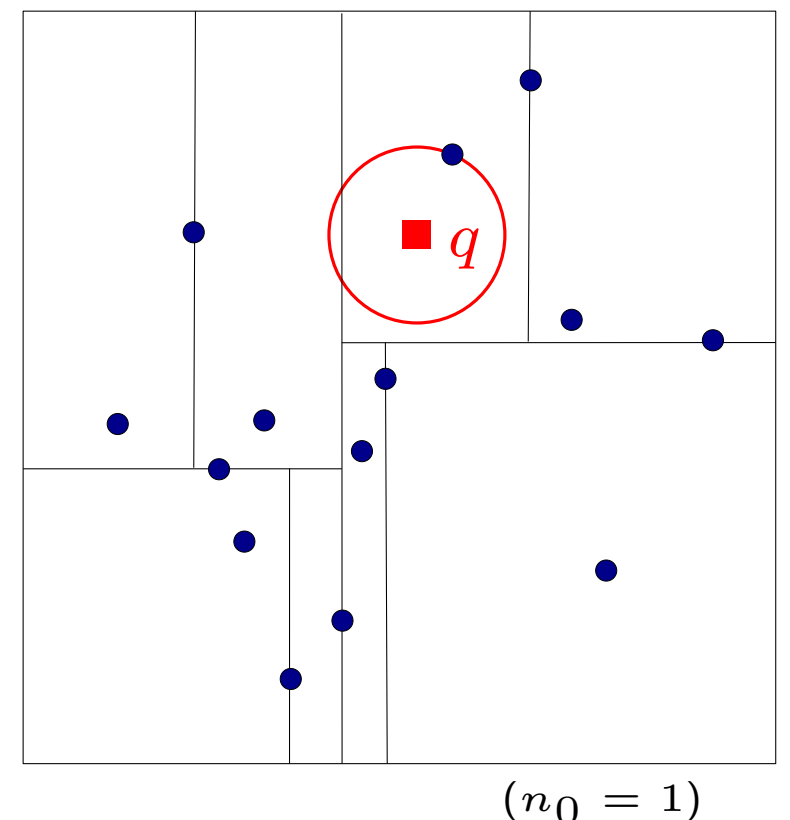
$d_{\min} := \min\{d_{\min}, d(q, \text{node.point})\}$

**if** *q* on "left" side of *node.H*

**recurse** on *node.left*

**else** (*q* on "right" side of *node.H*)

**recurse** on *node.right*



# Usage for NN search

## Strategy 1: defeatist search

$d_{\min} := \infty$  (dist. to pts viewed so far)

Query time:  $O(d(n_0 + \log \frac{n}{n_0}))$

search (*node*): (*node* = *root* initially)

**if** *node* = *leaf*:

$d_{\min} := \min\{d_{\min}, \min_{p \in \text{node.batch}} d(q, p)\}$

**else:**

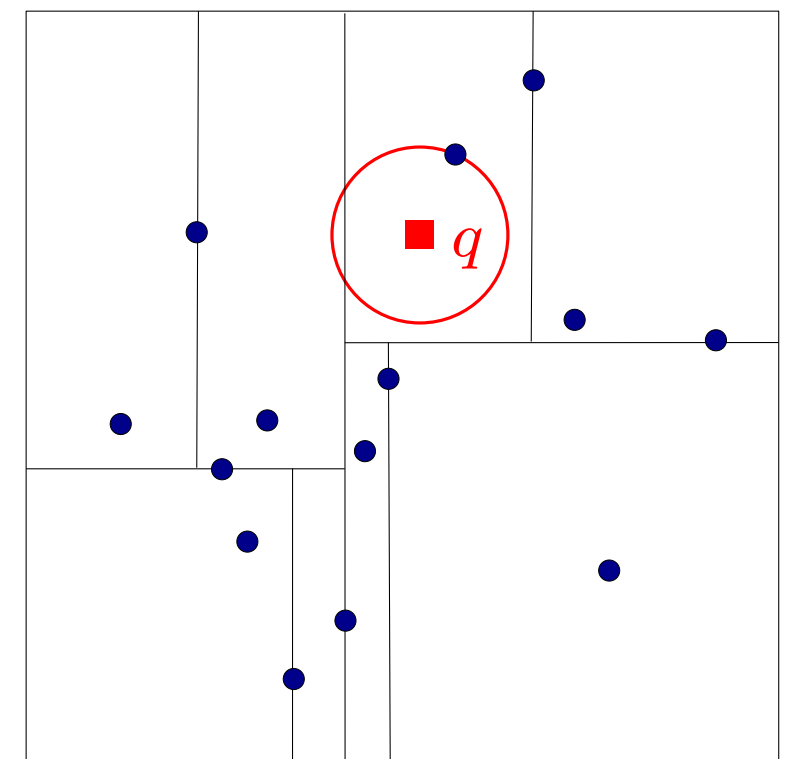
$d_{\min} := \min\{d_{\min}, d(q, \text{node.point})\}$

**if** *q* on "left" side of *node.H*

**recurse** on *node.left*

**else** (*q* on "right" side of *node.H*)

**recurse** on *node.right*



( $n_0 = 1$ )

# Usage for NN search

## Strategy 1: defeatist search

$d_{\min} := \infty$  (dist. to pts viewed so far)

Query time:  $O(d(n_0 + \log \frac{n}{n_0}))$

search ( $node$ ): ( $node = root$  initially)

May fail!

**if**  $node = leaf$ :

$d_{\min} := \min\{d_{\min}, \min_{p \in node.batch} d(q, p)\}$

**else:**

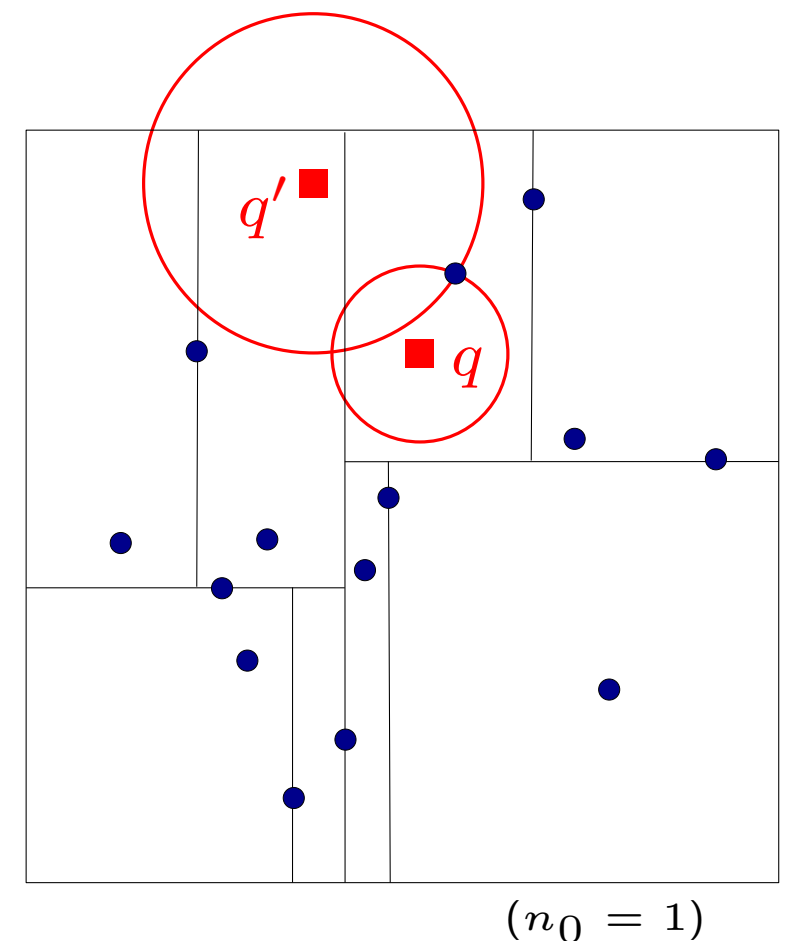
$d_{\min} := \min\{d_{\min}, d(q, node.point)\}$

**if**  $q$  on "left" side of  $node.H$

**recurse** on  $node.left$

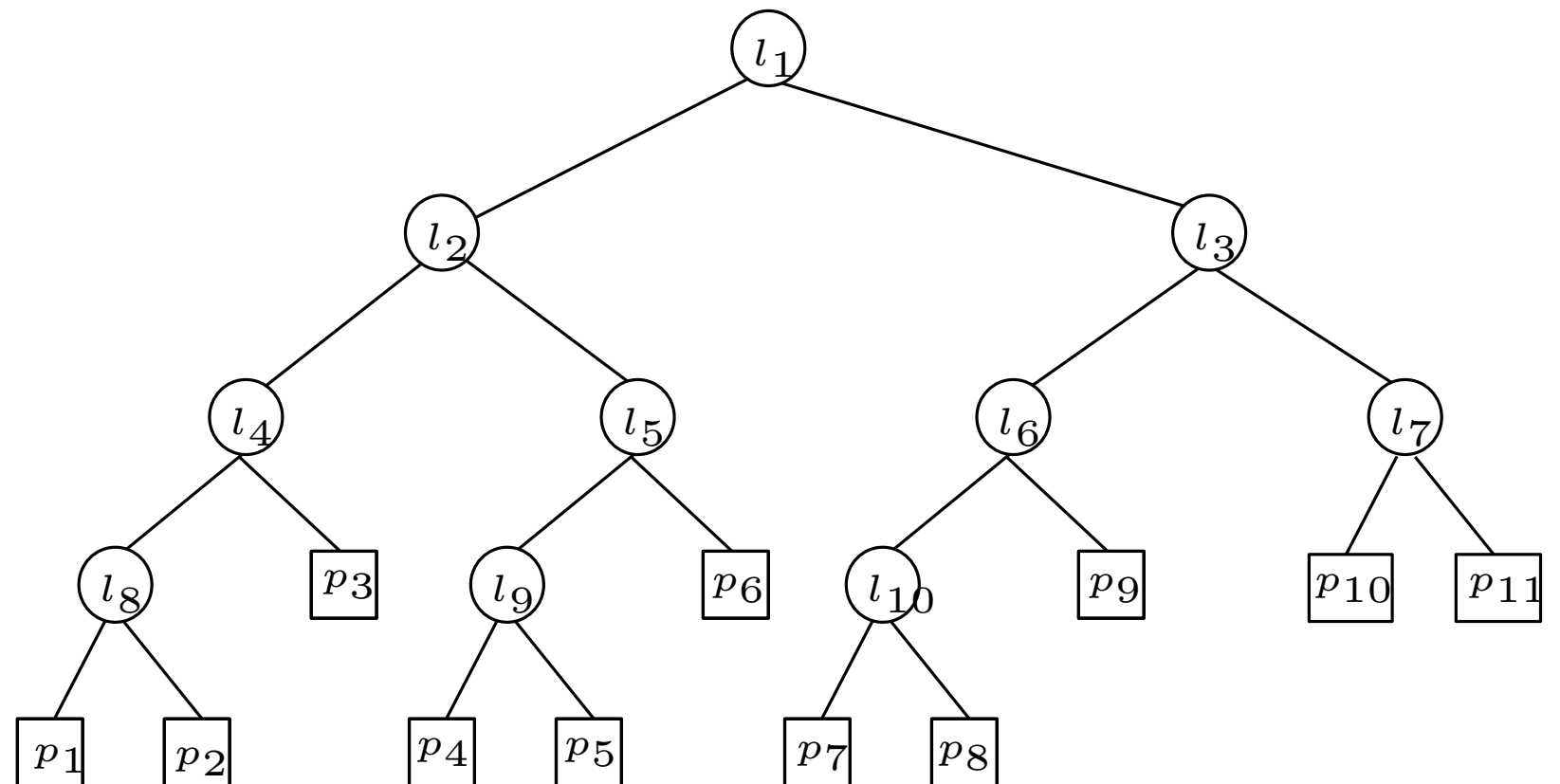
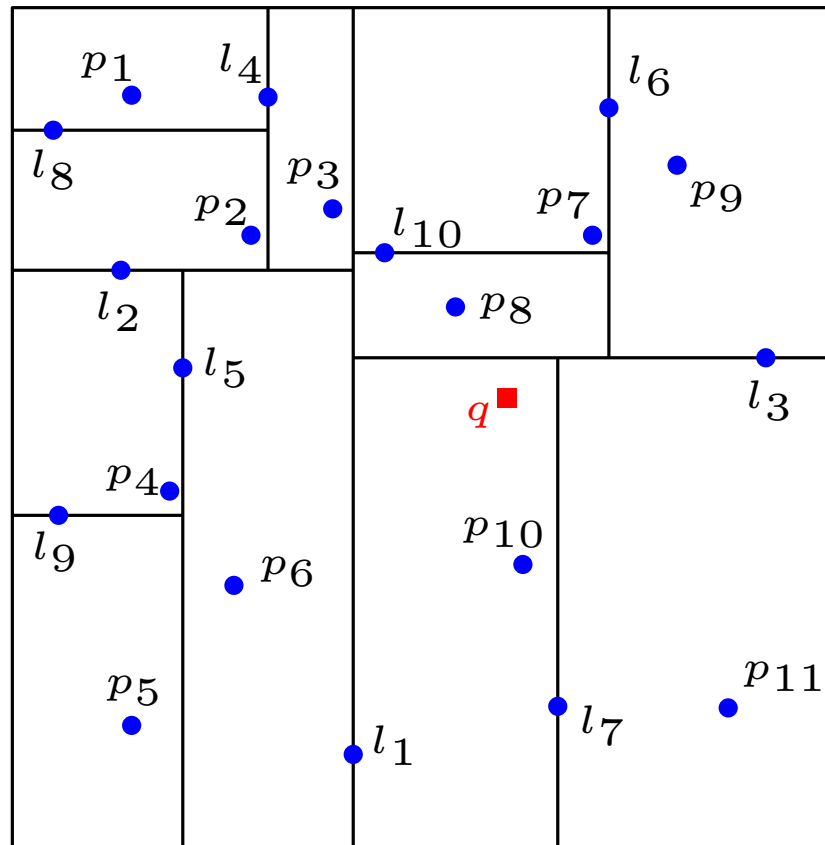
**else** ( $q$  on "right" side of  $node.H$ )

**recurse** on  $node.right$





# Example

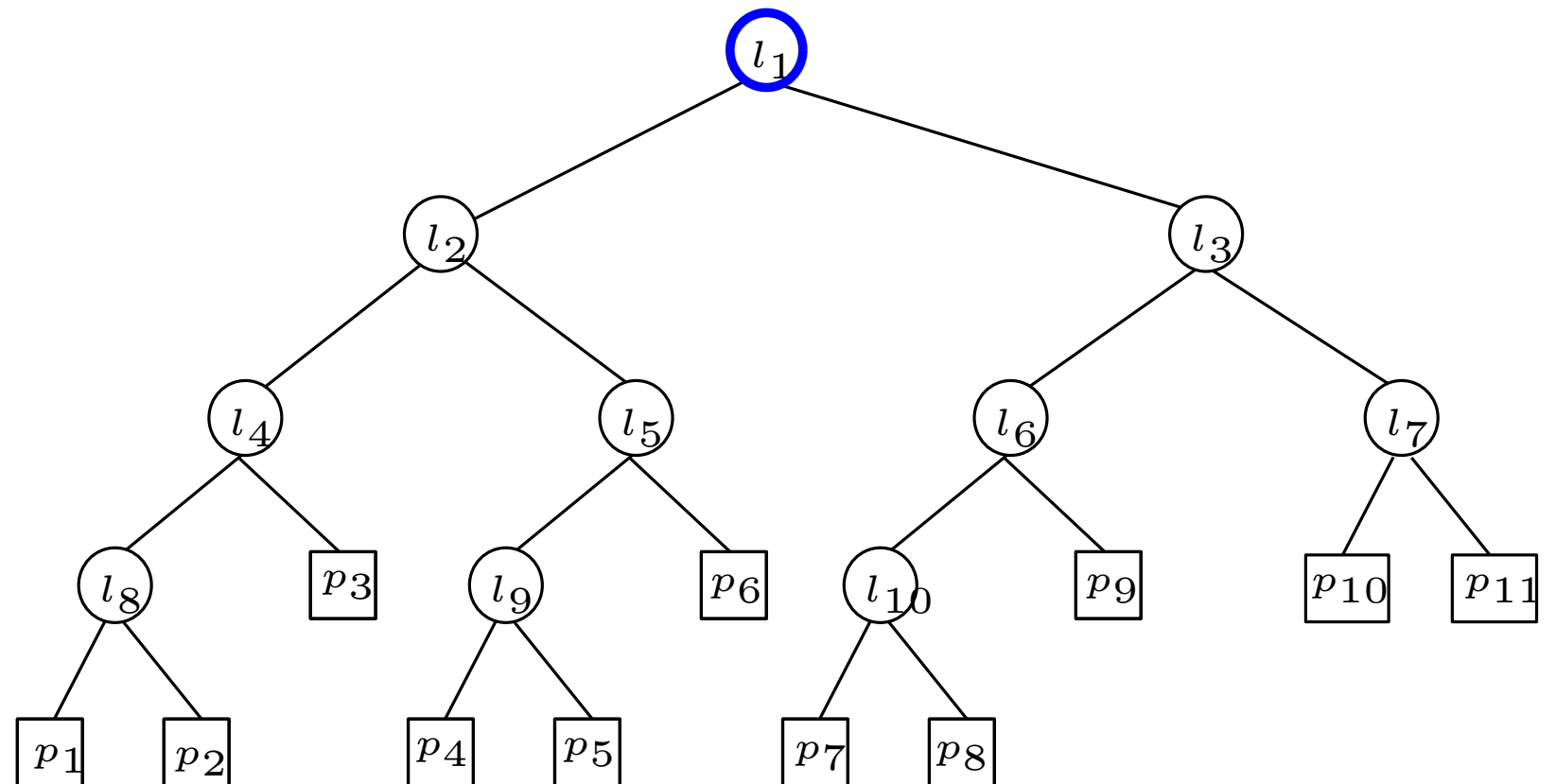
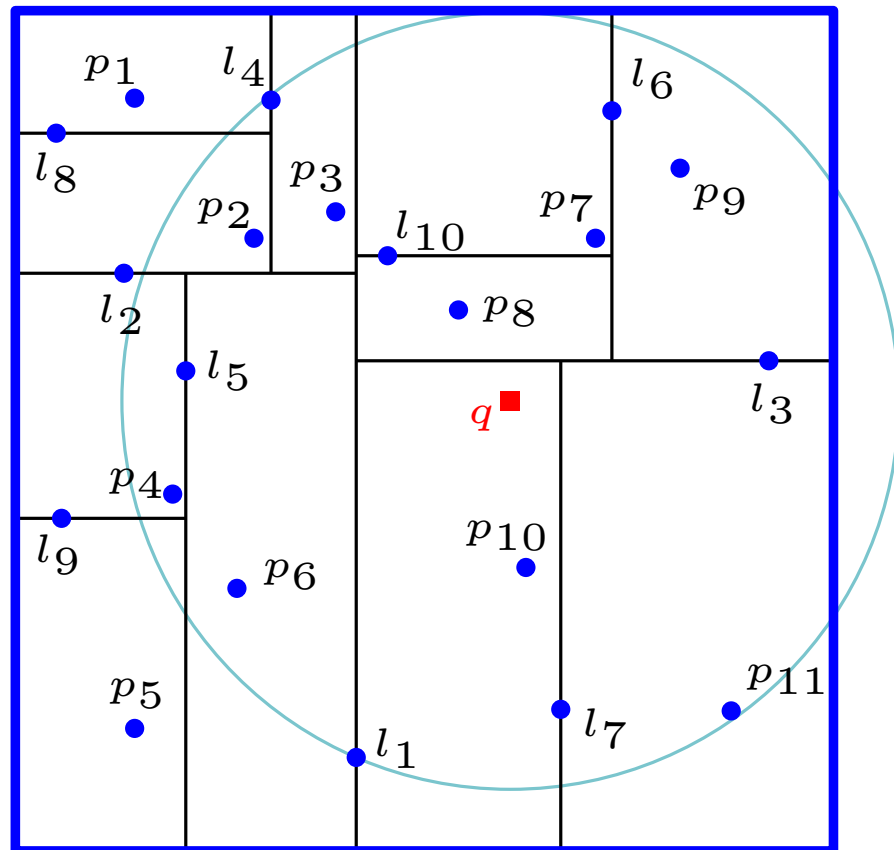


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

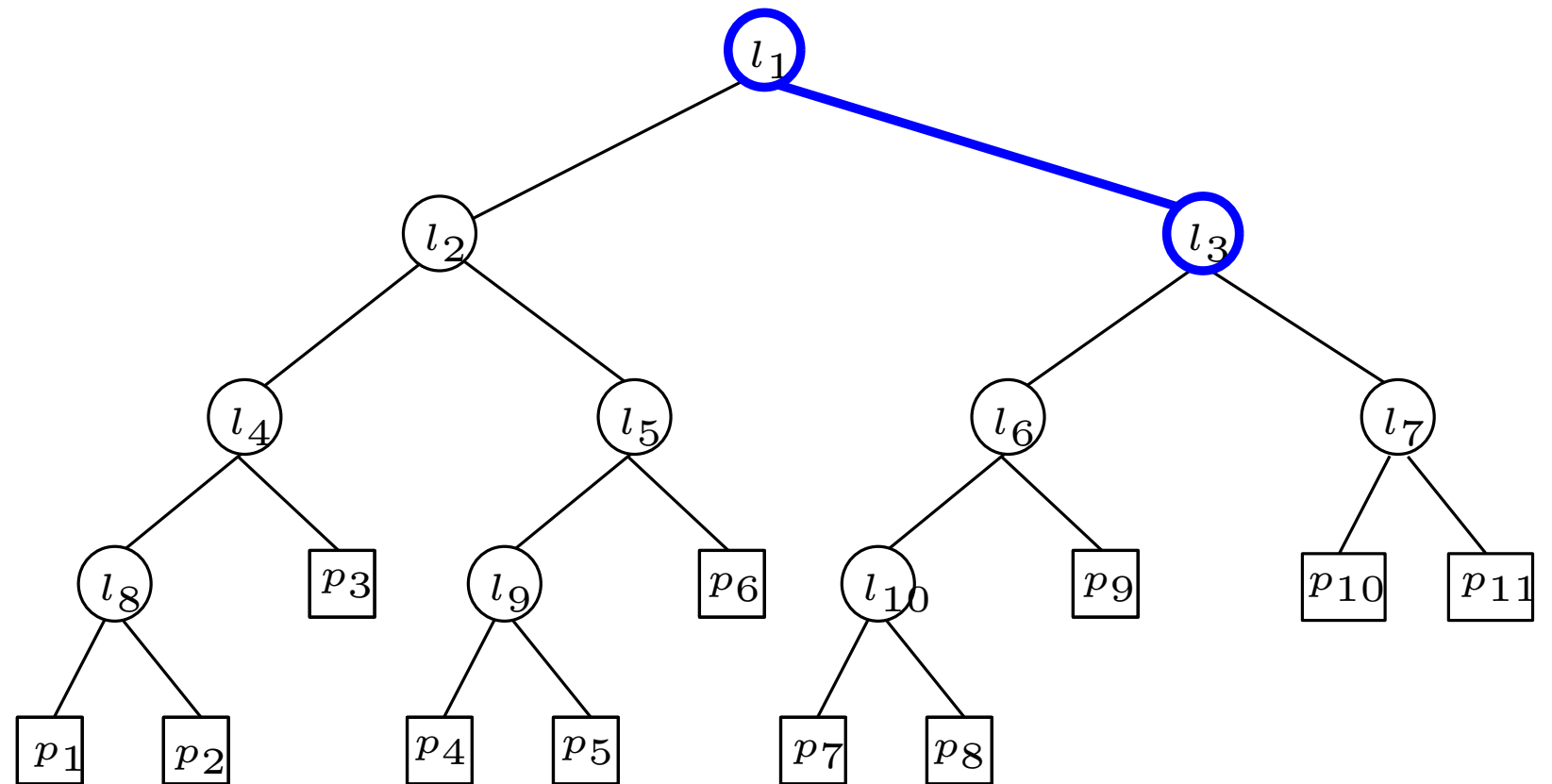
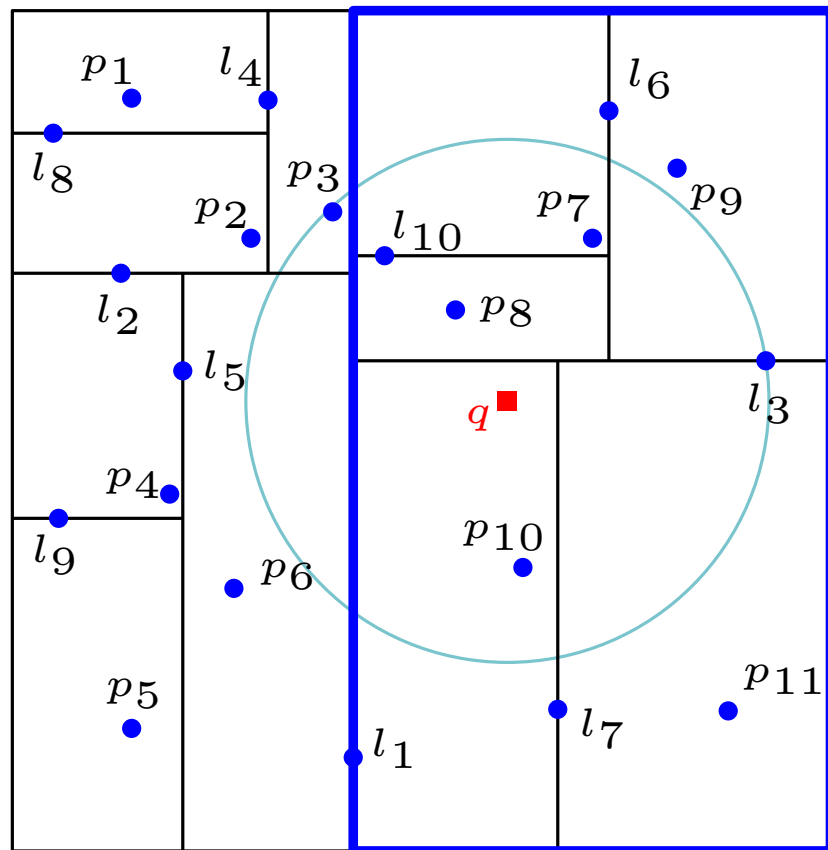


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

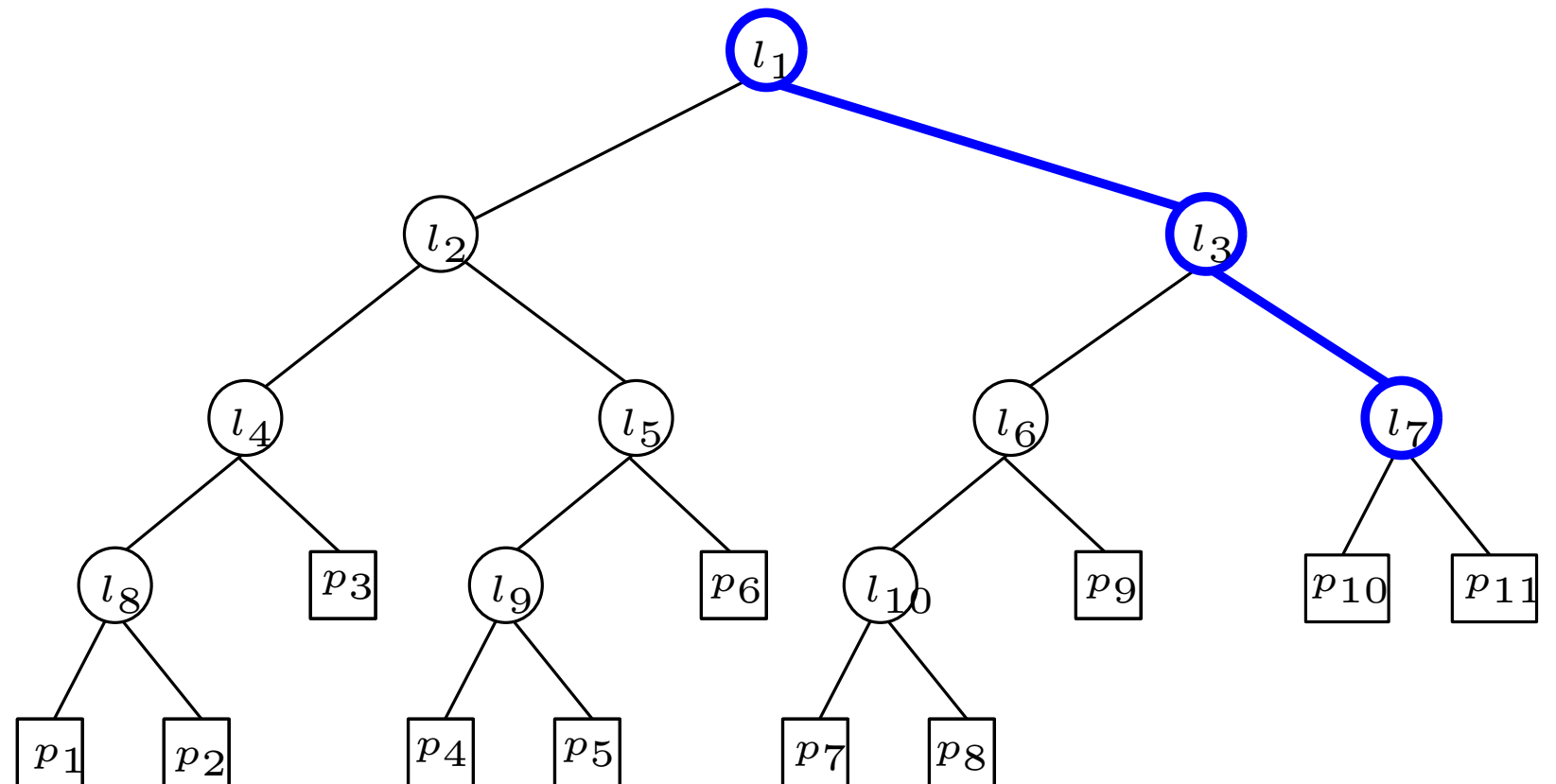
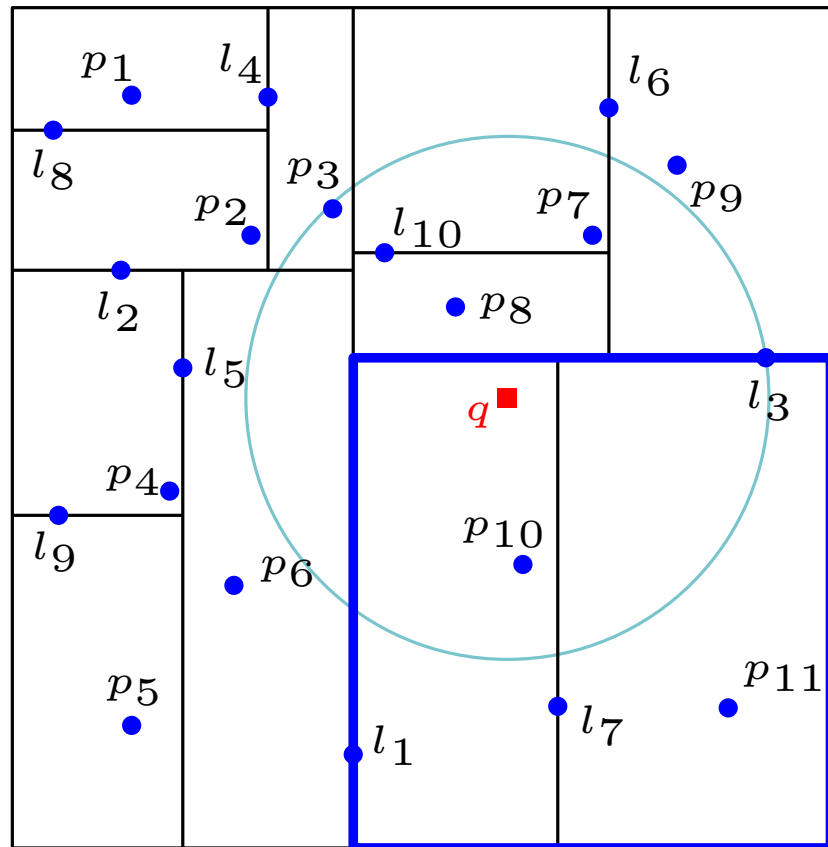


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

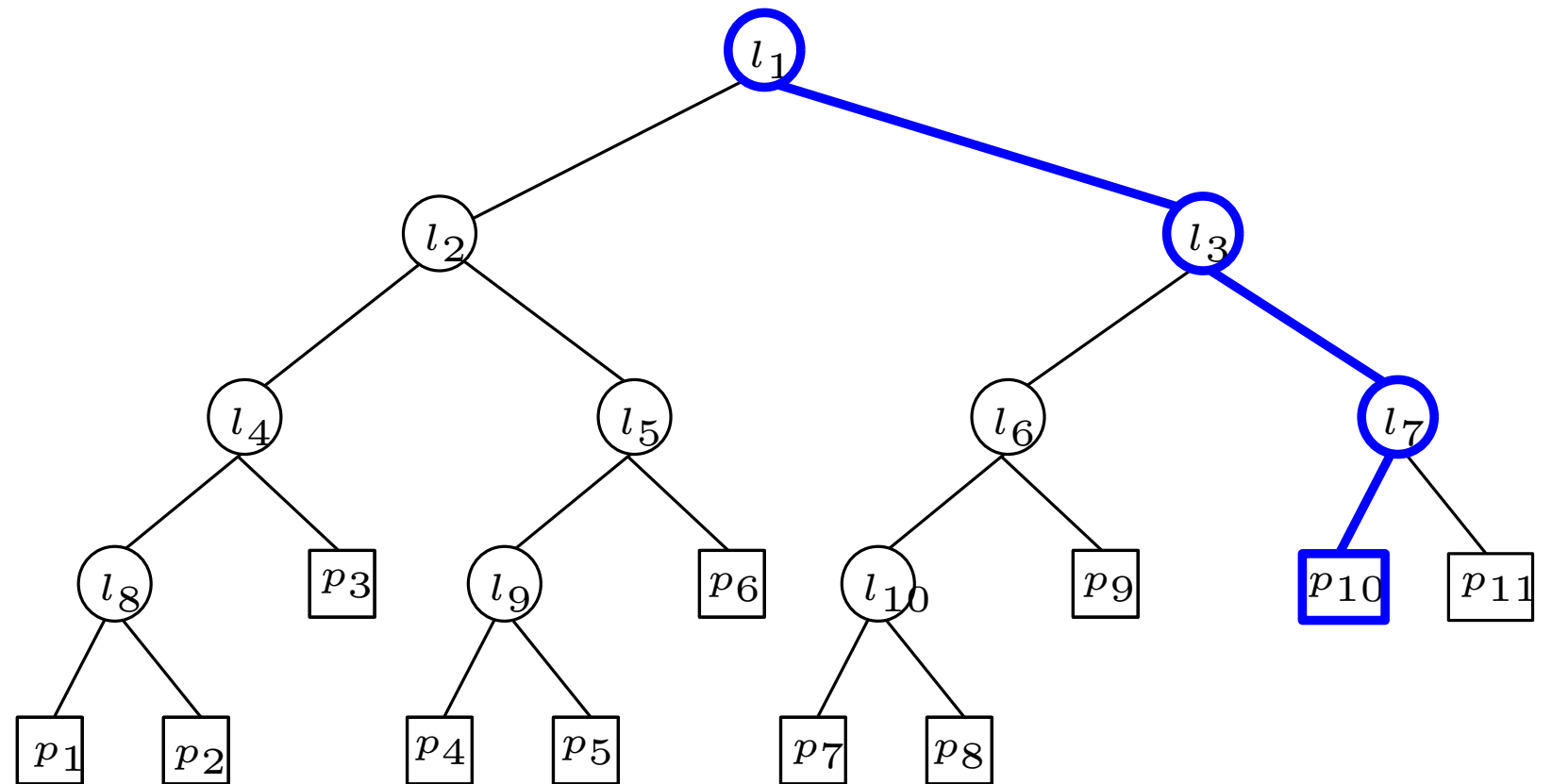
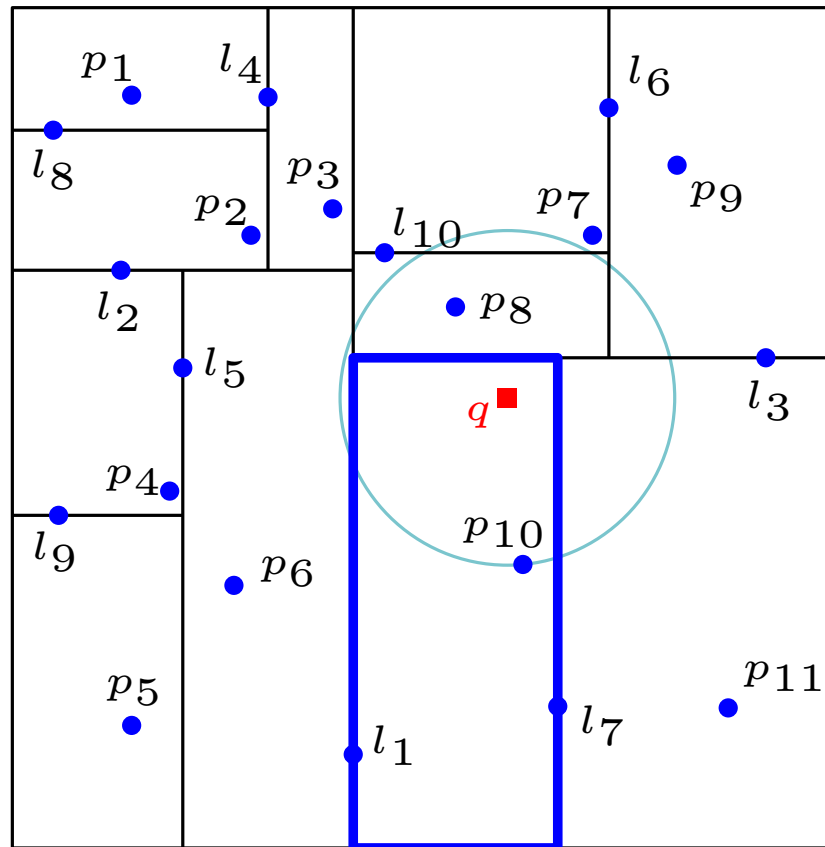


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example



$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Usage for NN search

## Strategy 2: **backtracking** search

$d_{\min} := \infty$  (dist. to pts viewed so far)

search (*node*): (*node* = *root* initially)

**if** *node* = *leaf*:

$d_{\min} := \min\{d_{\min}, \min_{p \in \text{node.batch}} d(q, p)\}$

**else:**

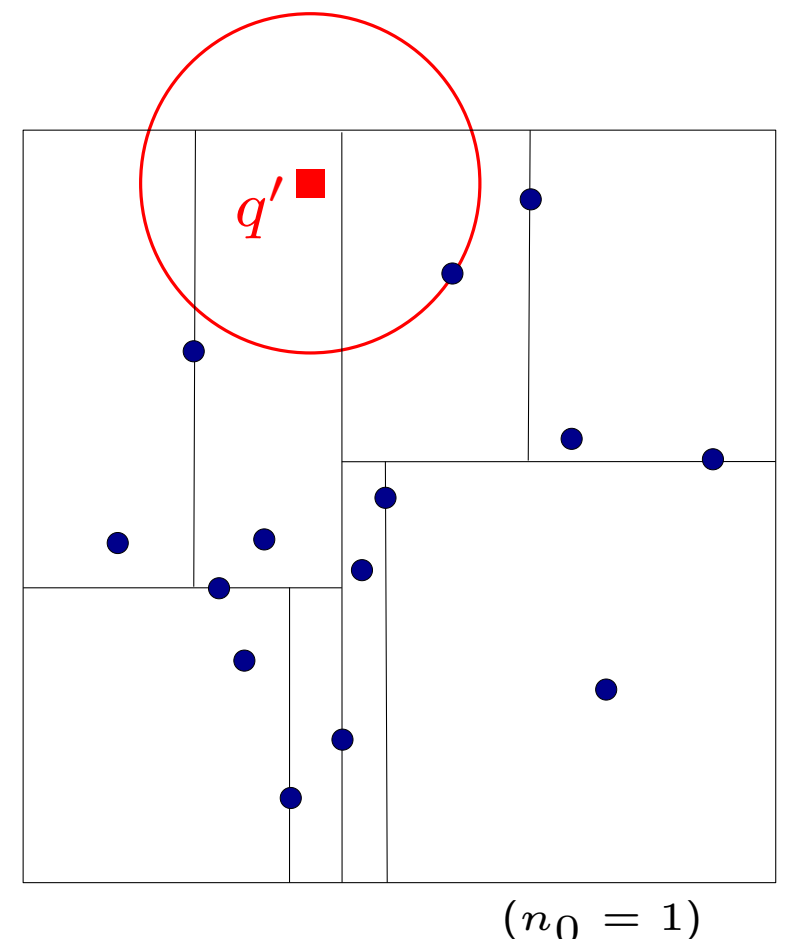
$d_{\min} := \min\{d_{\min}, d(q, \text{node.point})\}$

**if**  $B(q, d_{\min})$  intersects "left" side of *node.H*

**recurse** on *node.left*

**if**  $B(q, d_{\min})$  intersects "right" side of *node.H*

**recurse** on *node.right*



# Usage for NN search

## Strategy 2: **backtracking** search

$d_{\min} := \infty$  (dist. to pts viewed so far)

Always succeeds

search (*node*): (*node* = *root* initially)

$d_{\min} \geq d(q, \text{NN}(q)) \Rightarrow B(q, d_{\min})$   
intersects all cells containing  $\text{NN}(q)$   
in subdivision throughout search

**if** *node* = *leaf*:

$d_{\min} := \min\{d_{\min}, \min_{p \in \text{node.batch}} d(q, p)\}$

**else:**

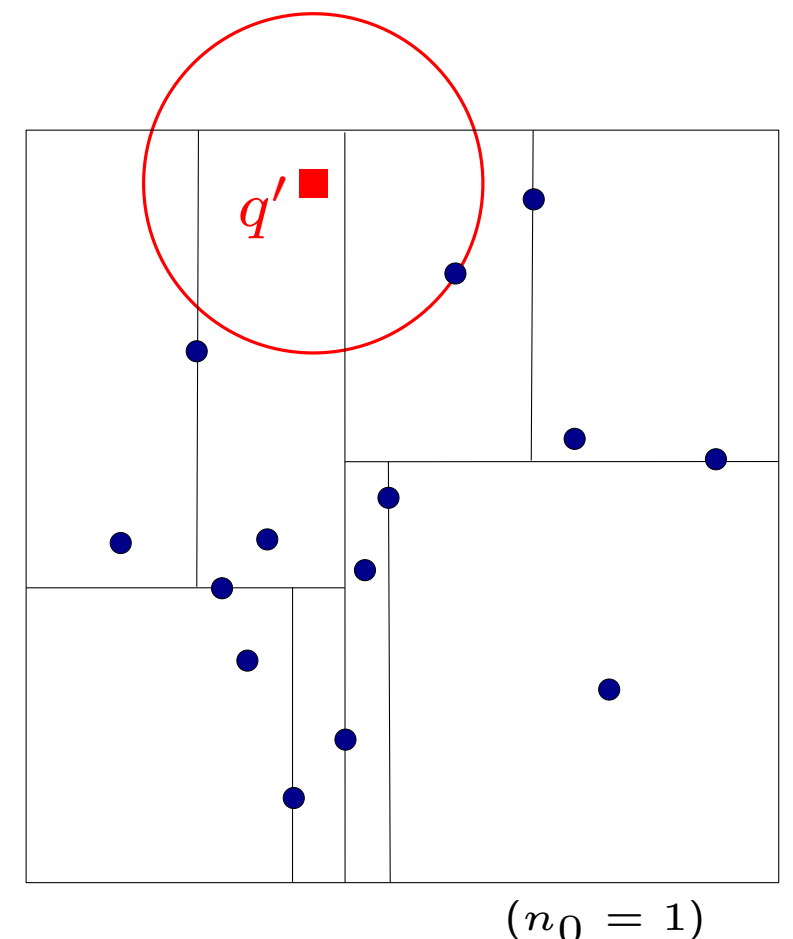
$d_{\min} := \min\{d_{\min}, d(q, \text{node.point})\}$

**if**  $B(q, d_{\min})$  intersects "left" side of *node.H*

**recurse** on *node.left*

**if**  $B(q, d_{\min})$  intersects "right" side of *node.H*

**recurse** on *node.right*



# Usage for NN search

## Strategy 2: **backtracking** search

$d_{\min} := \infty$  (dist. to pts viewed so far)

Always succeeds

search (*node*): (*node* = *root* initially)

Query time may be up to linear  
(all cells visited)

**if** *node* = *leaf*:

$d_{\min} := \min\{d_{\min}, \min_{p \in \text{node.batch}} d(q, p)\}$

**else:**

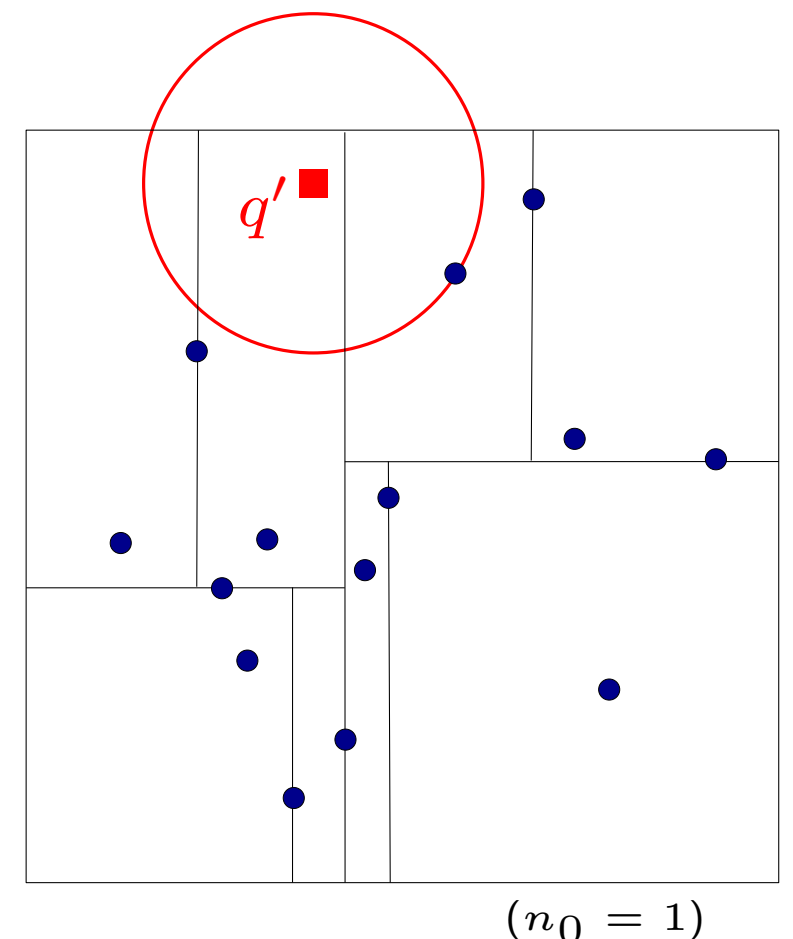
$d_{\min} := \min\{d_{\min}, d(q, \text{node.point})\}$

**if**  $B(q, d_{\min})$  intersects "left" side of *node.H*

**recurse** on *node.left*

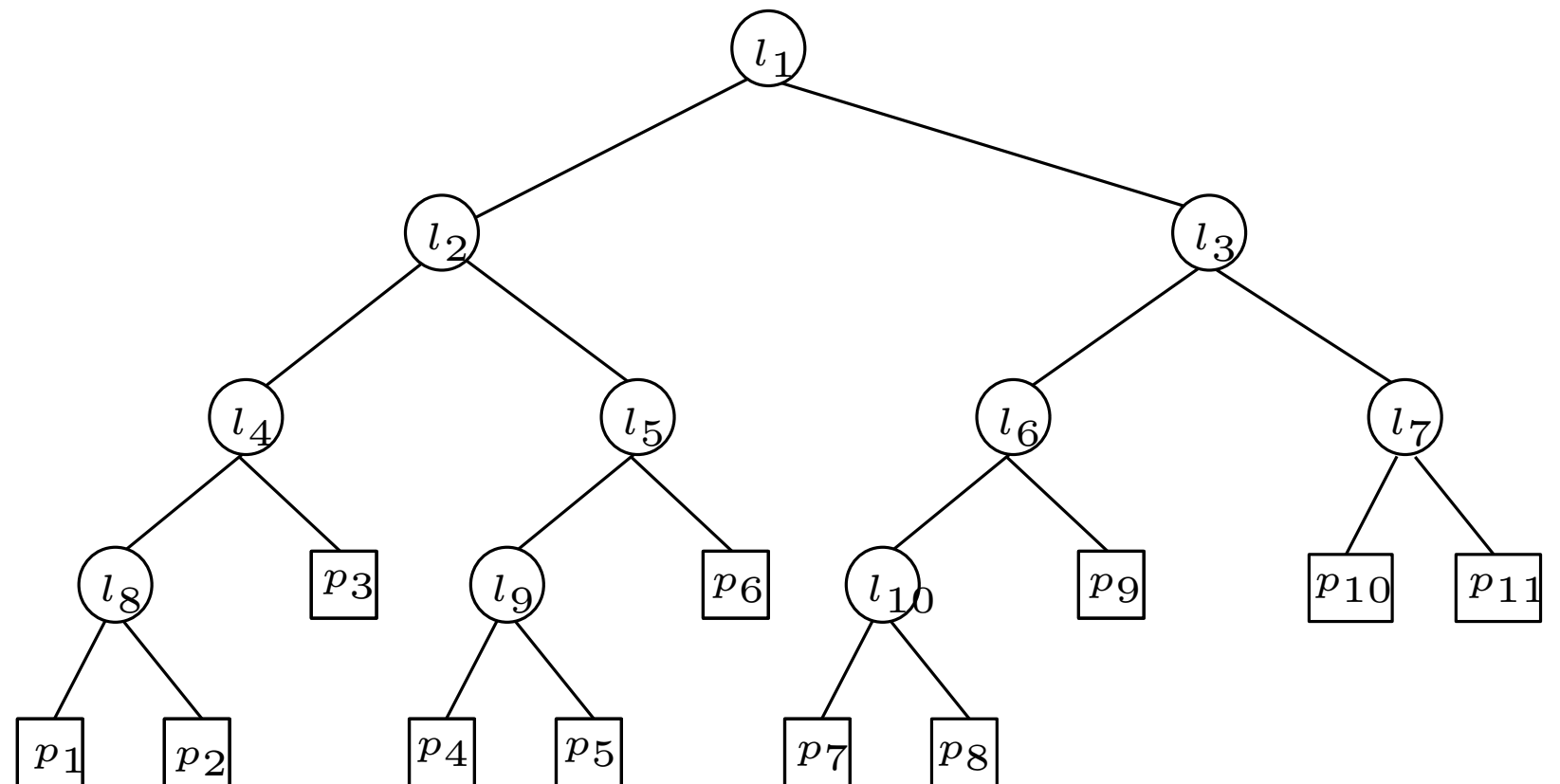
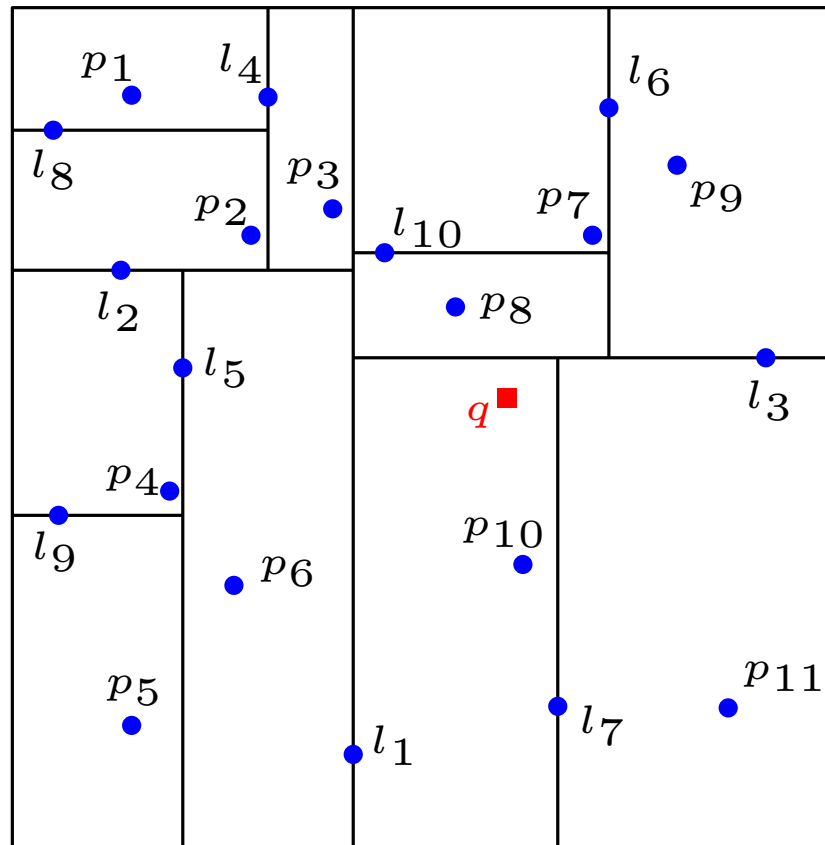
**if**  $B(q, d_{\min})$  intersects "right" side of *node.H*

**recurse** on *node.right*





# Example

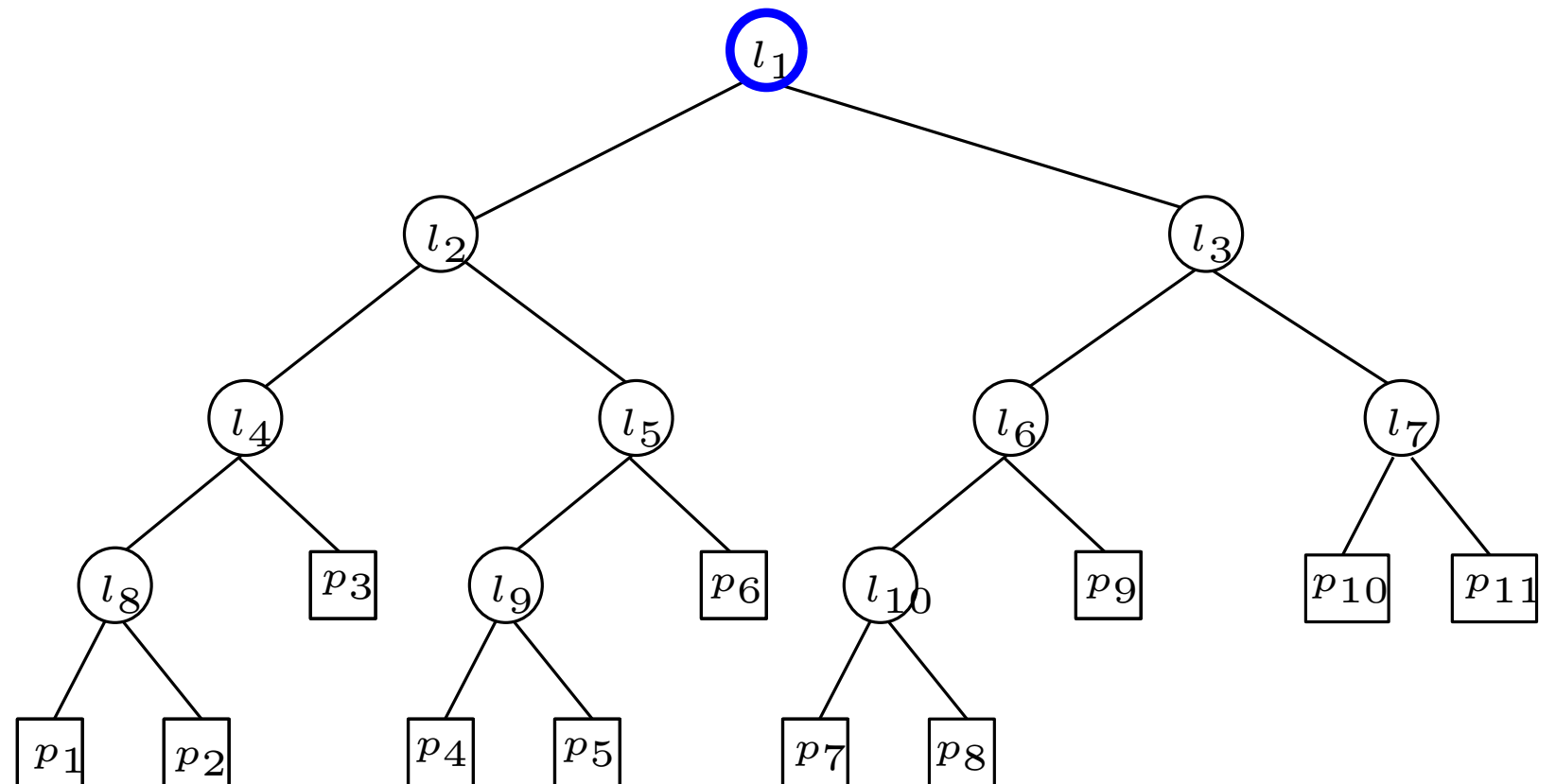
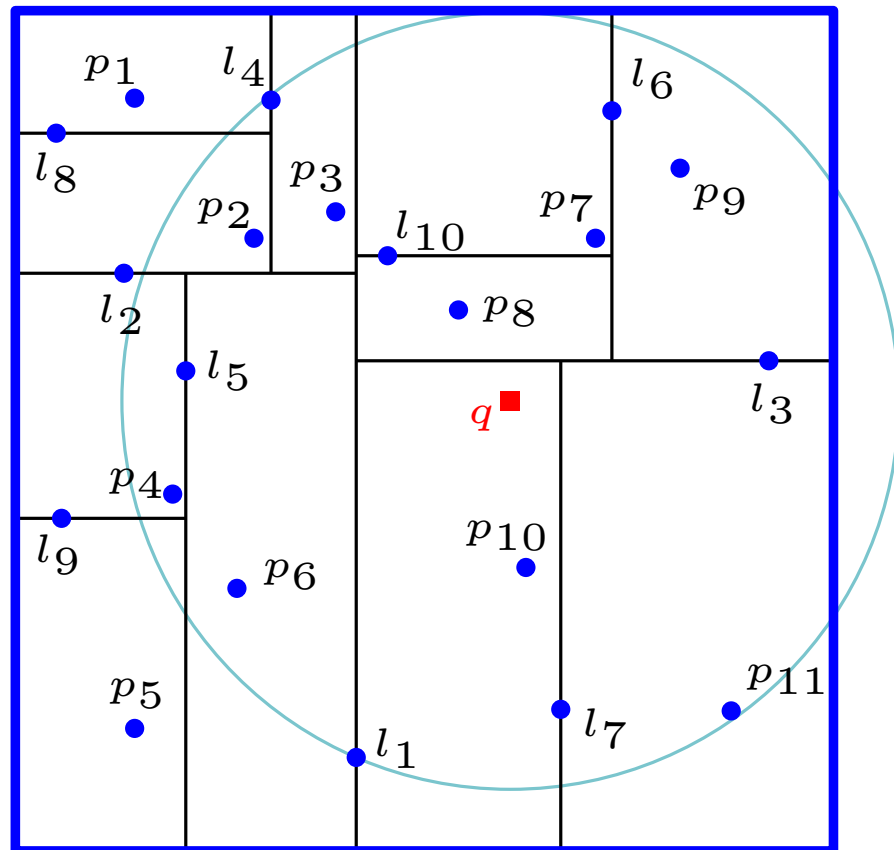


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

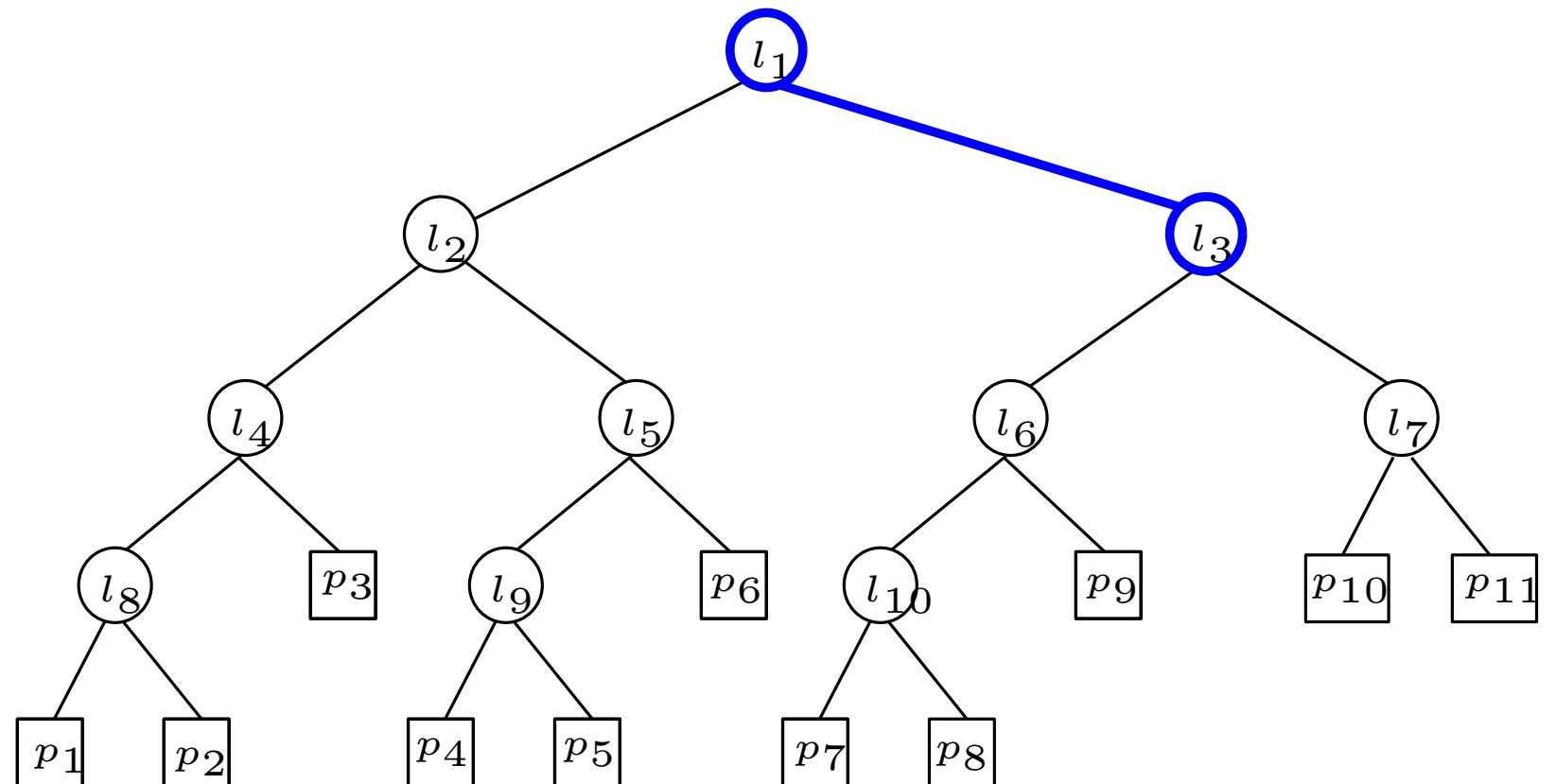
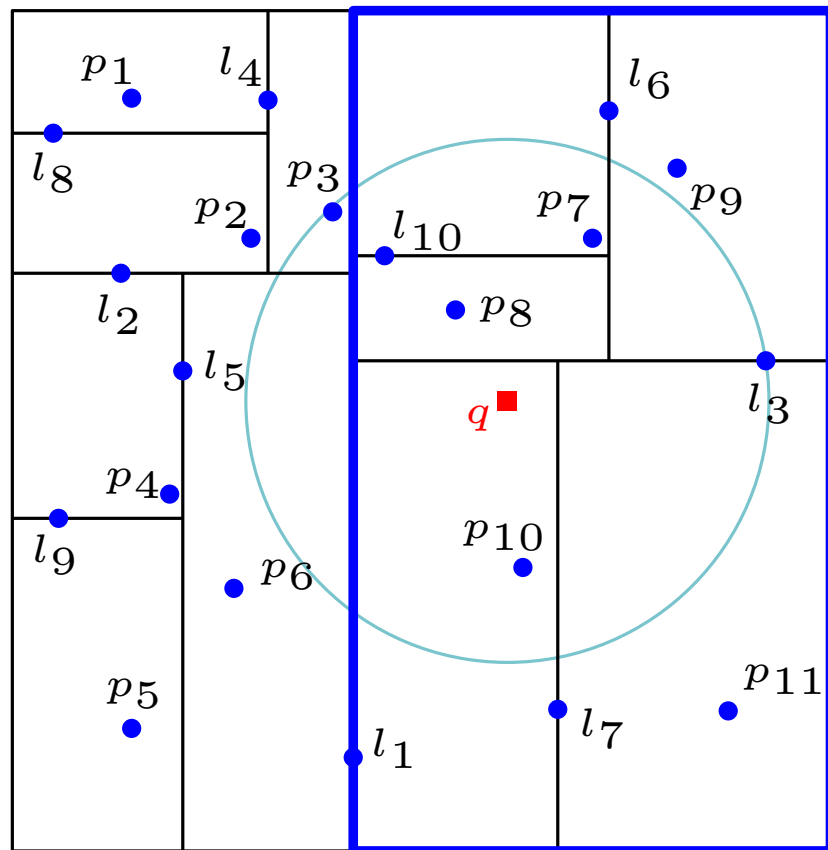


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

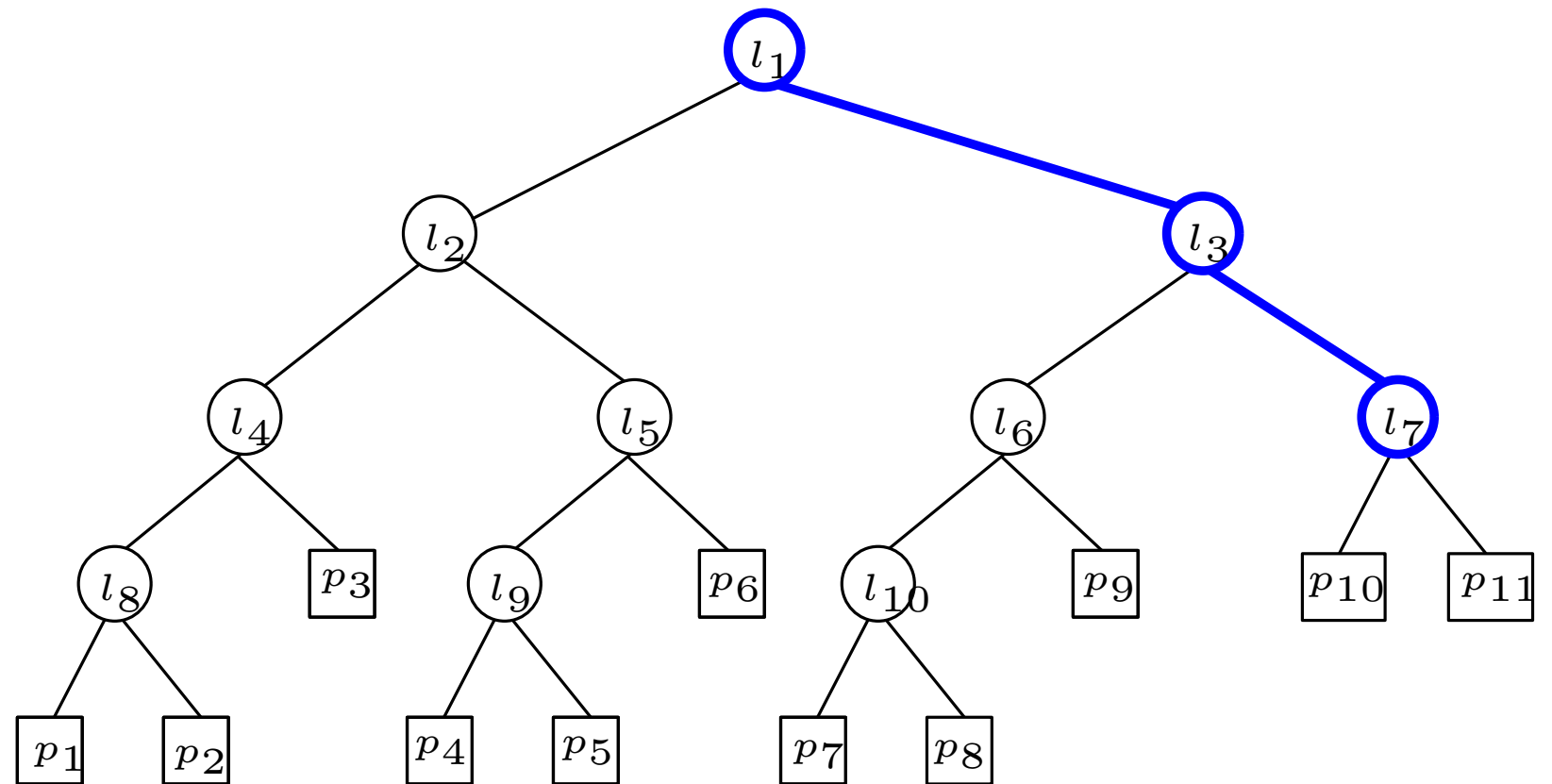
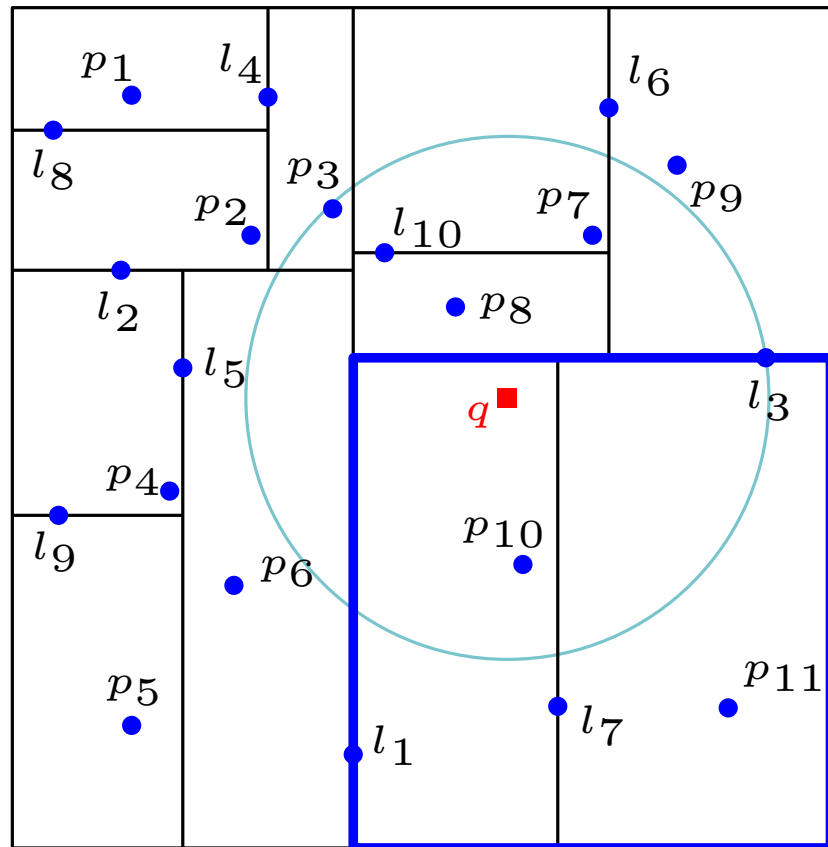


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

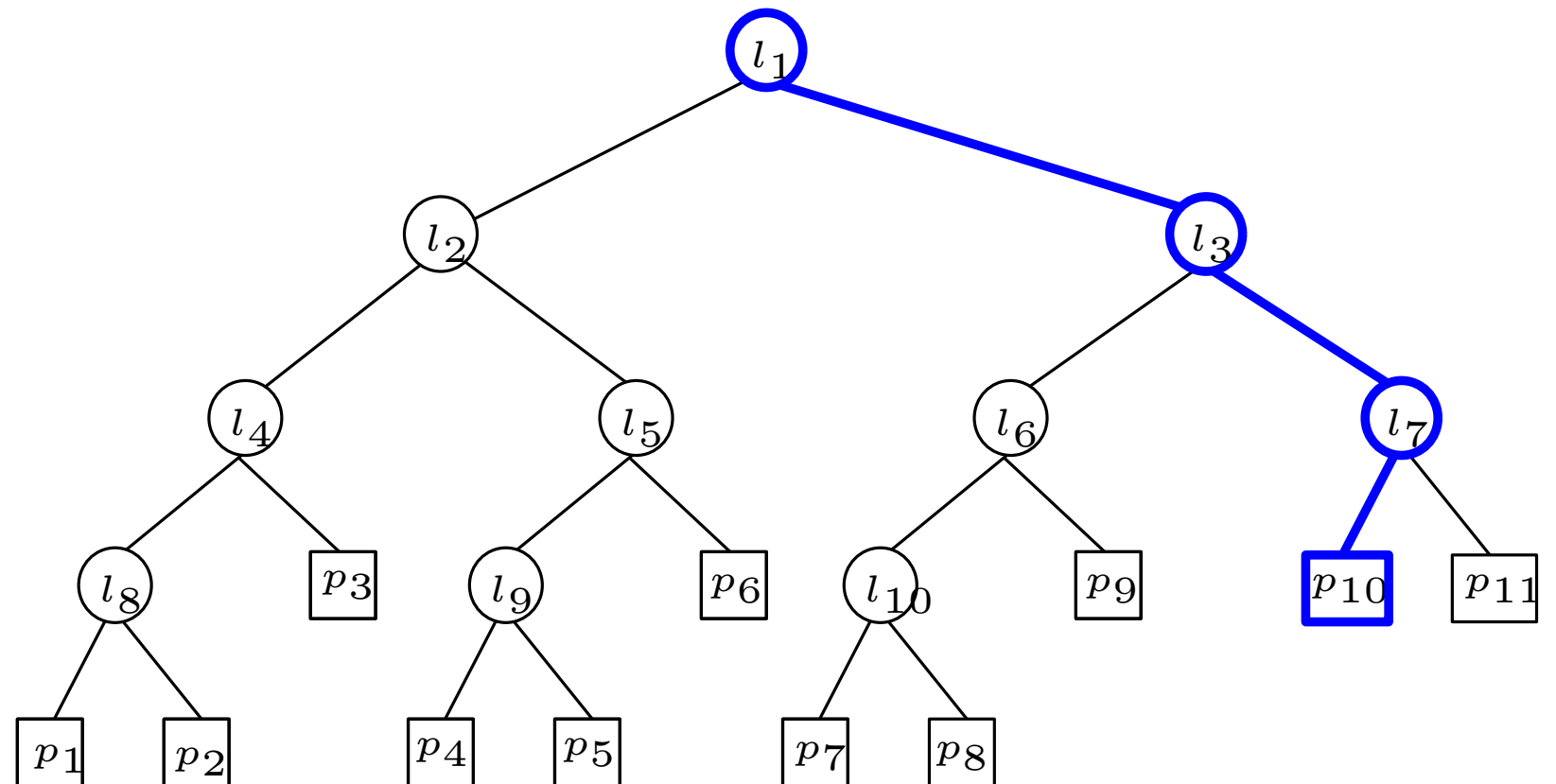
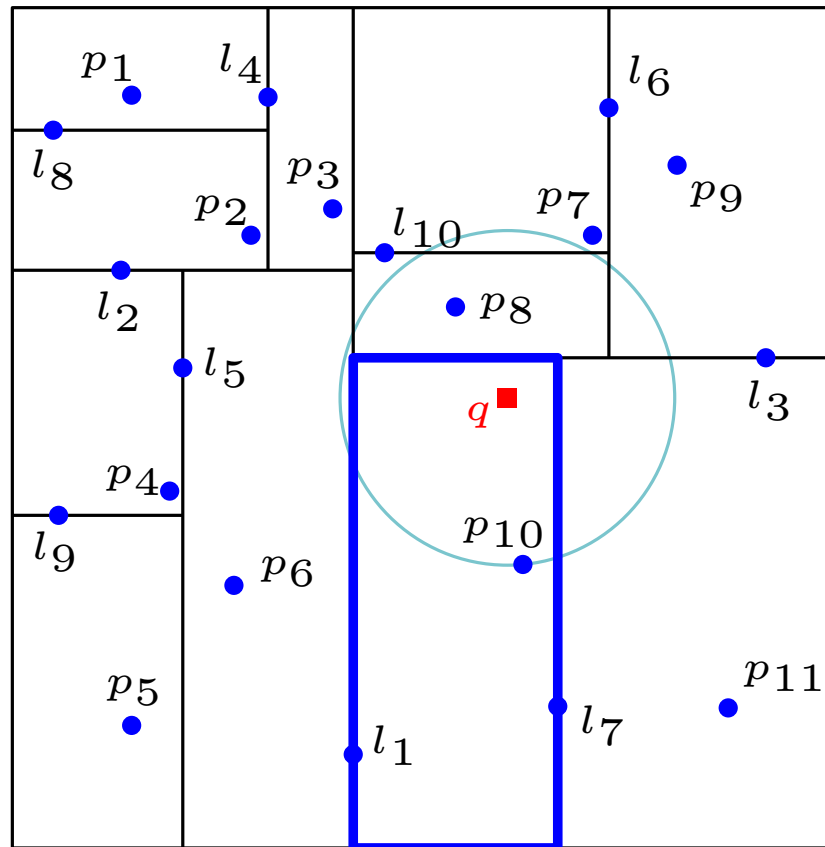


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

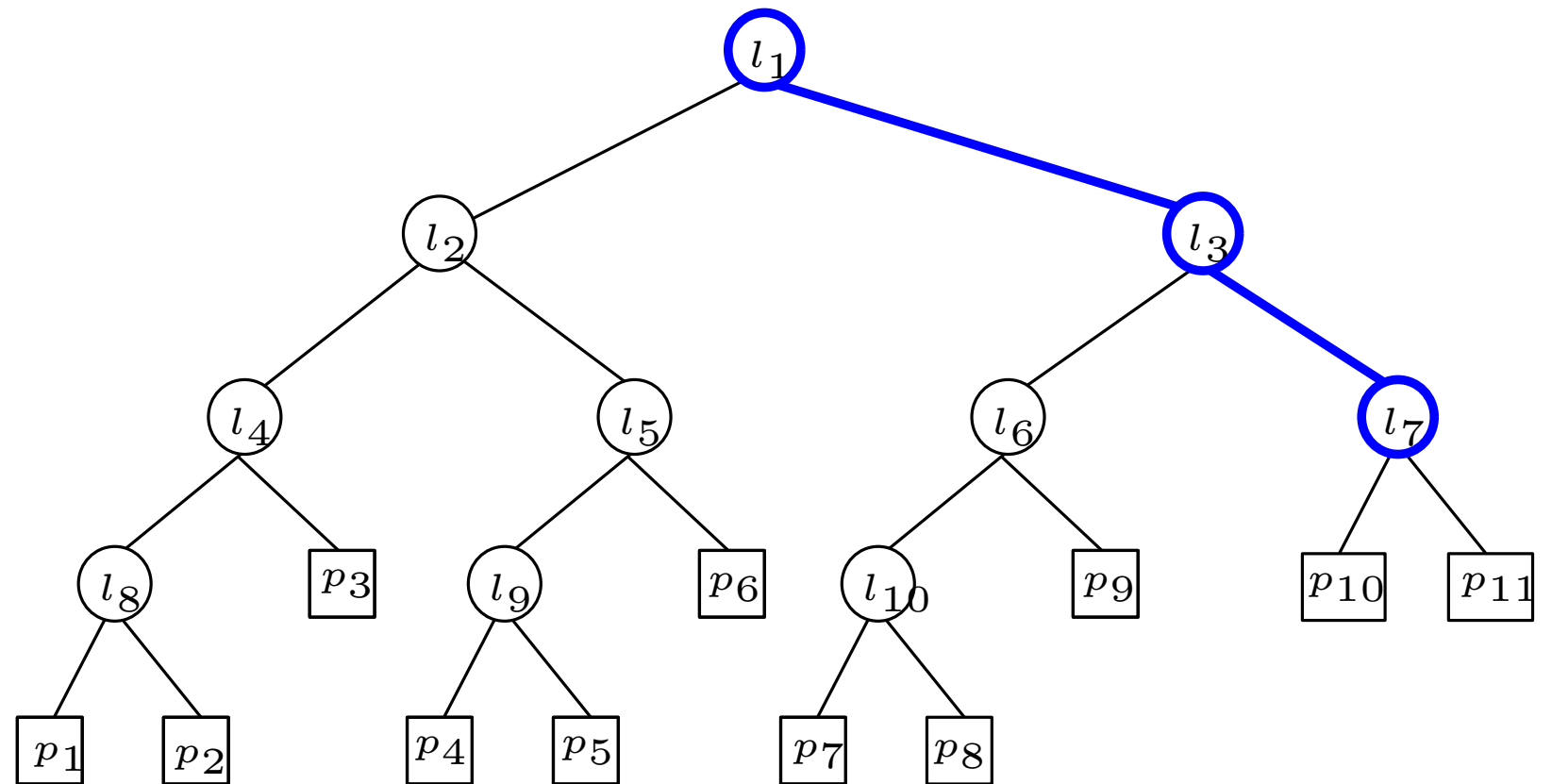
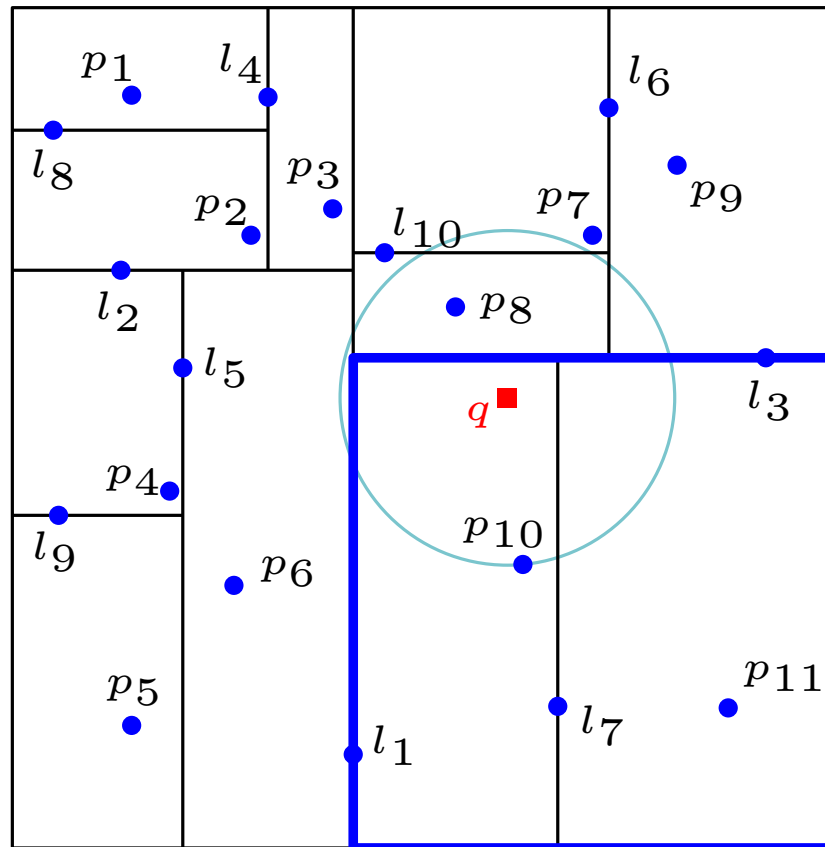


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

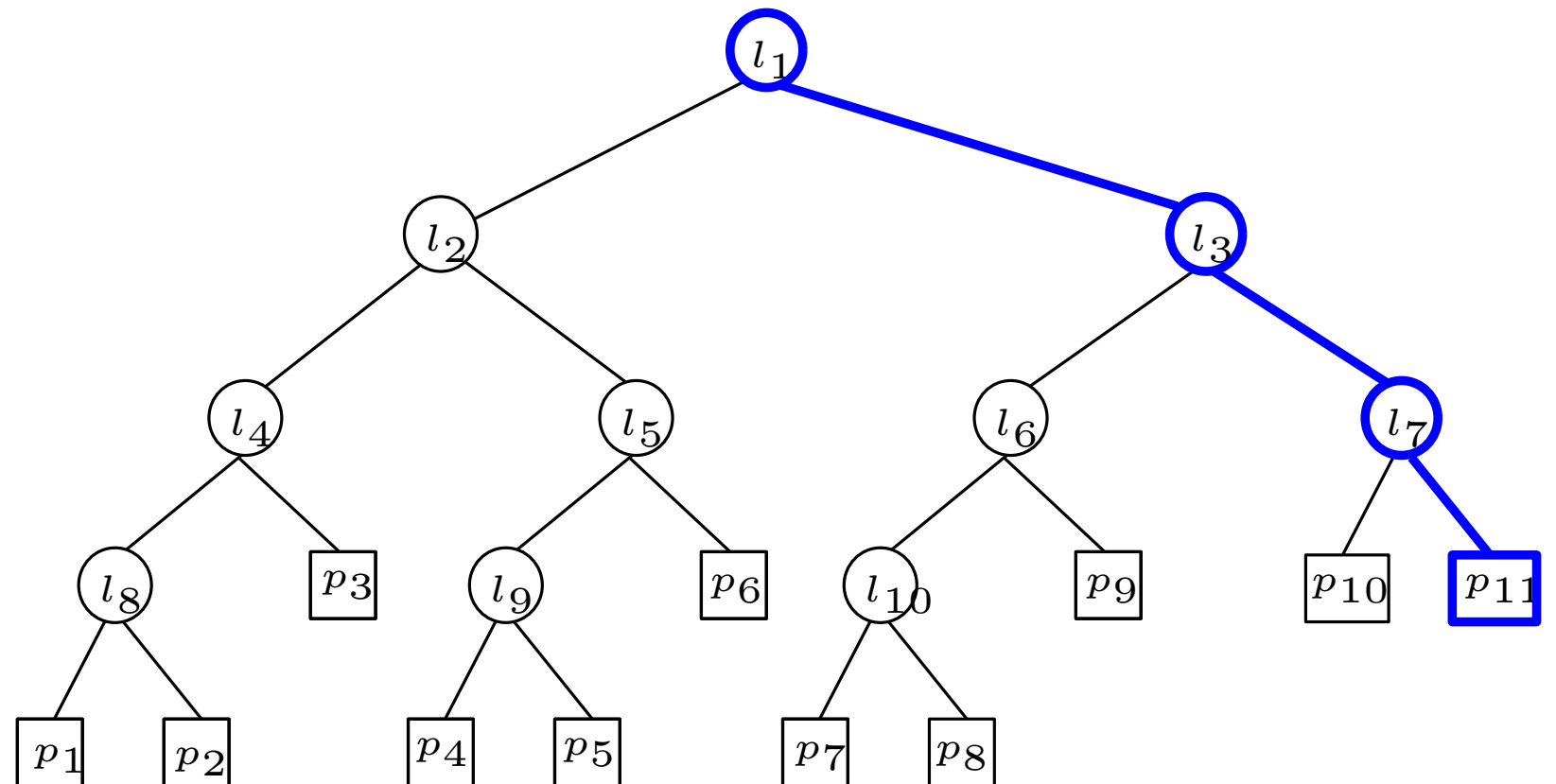
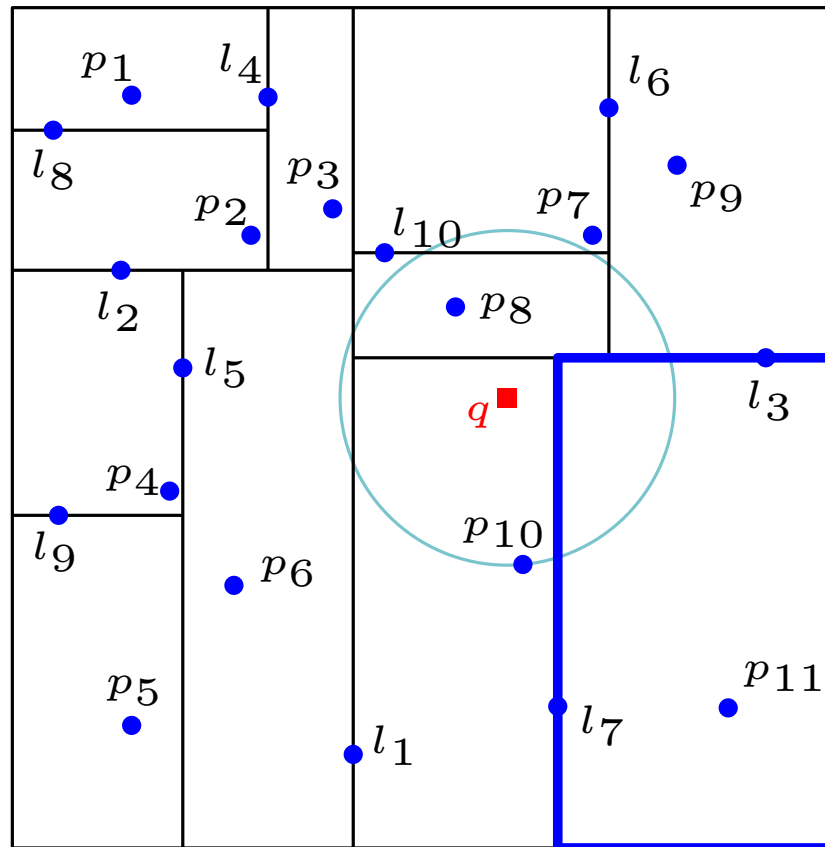


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

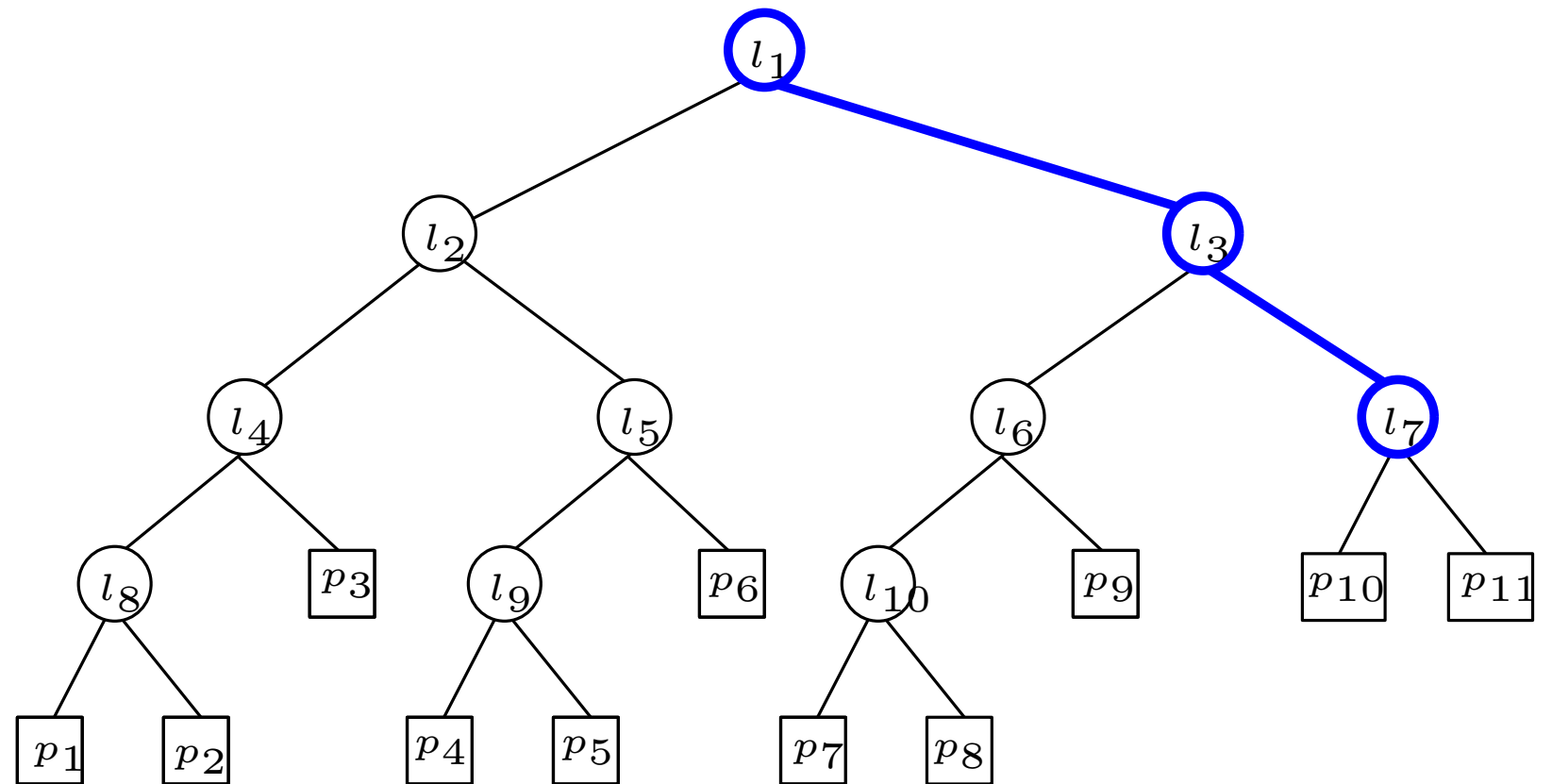
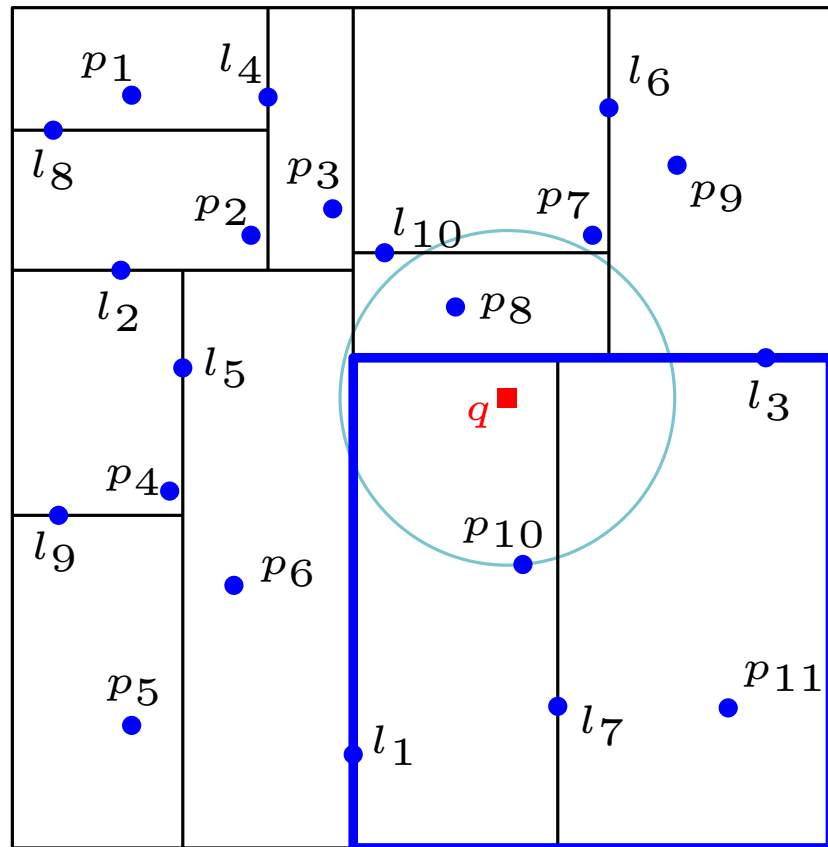


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example



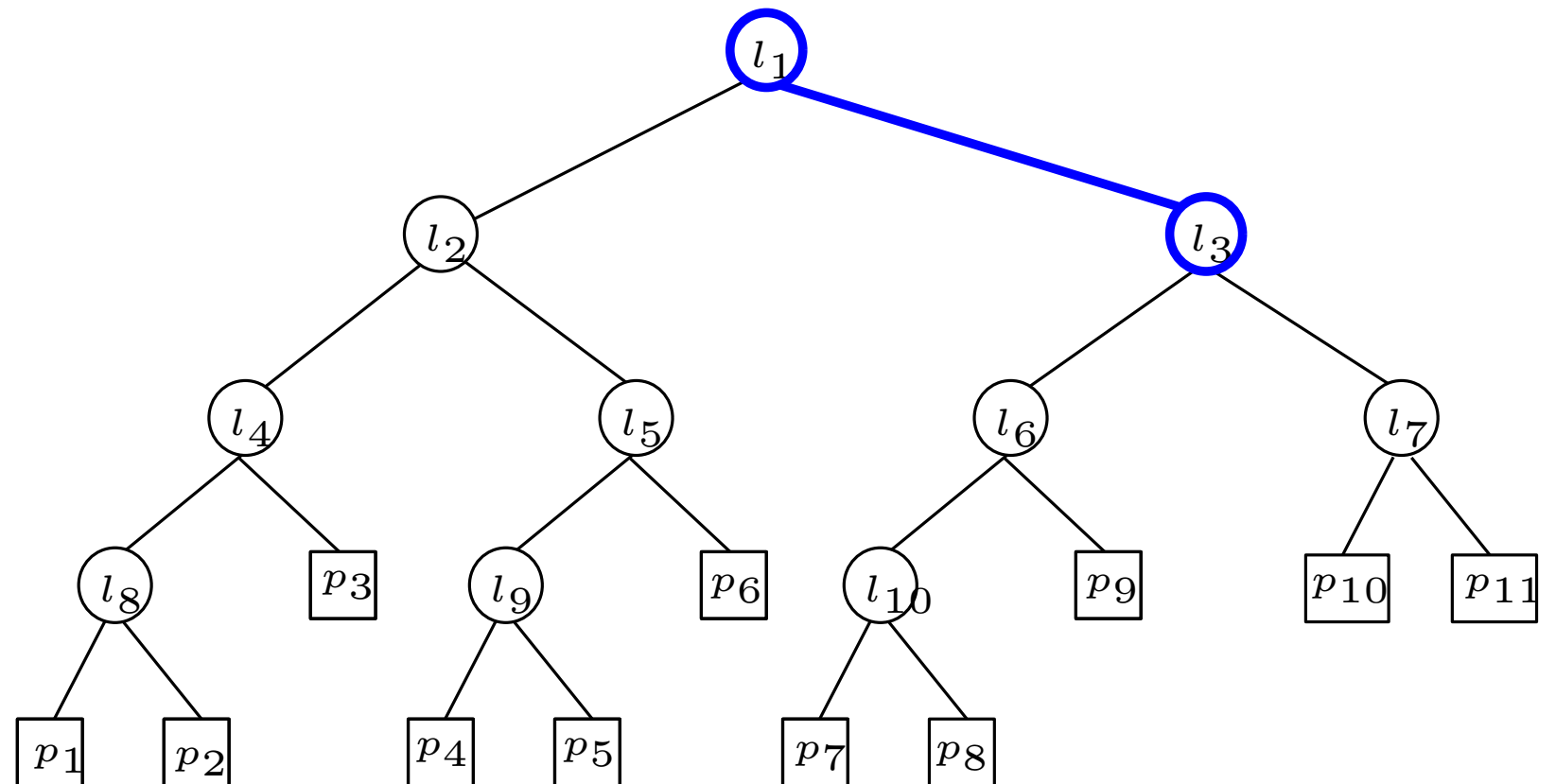
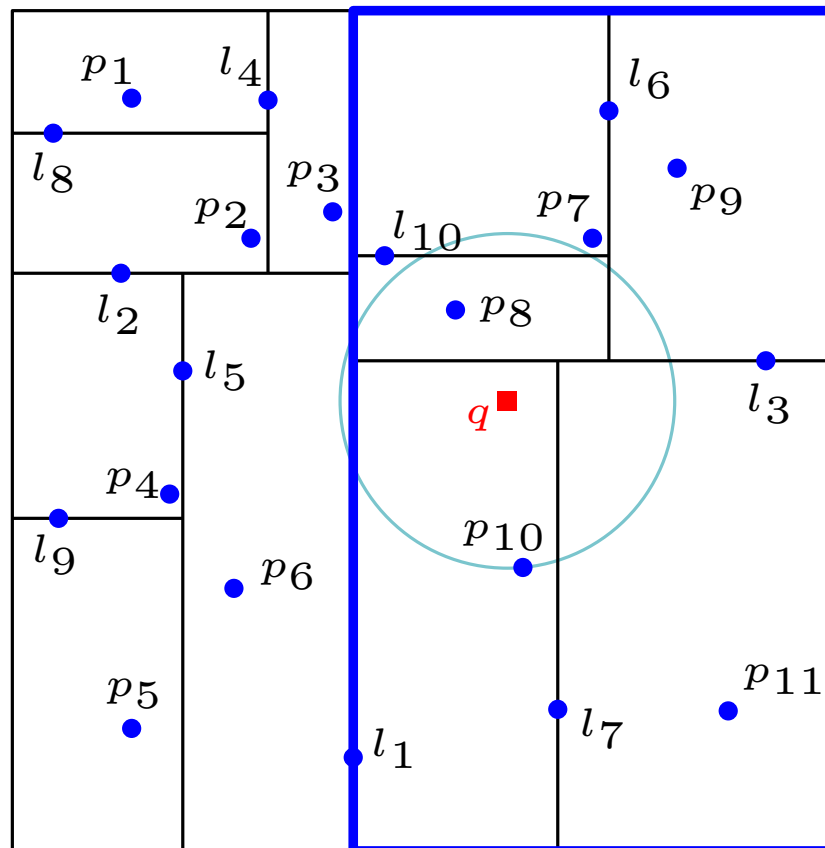
$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)



# Example

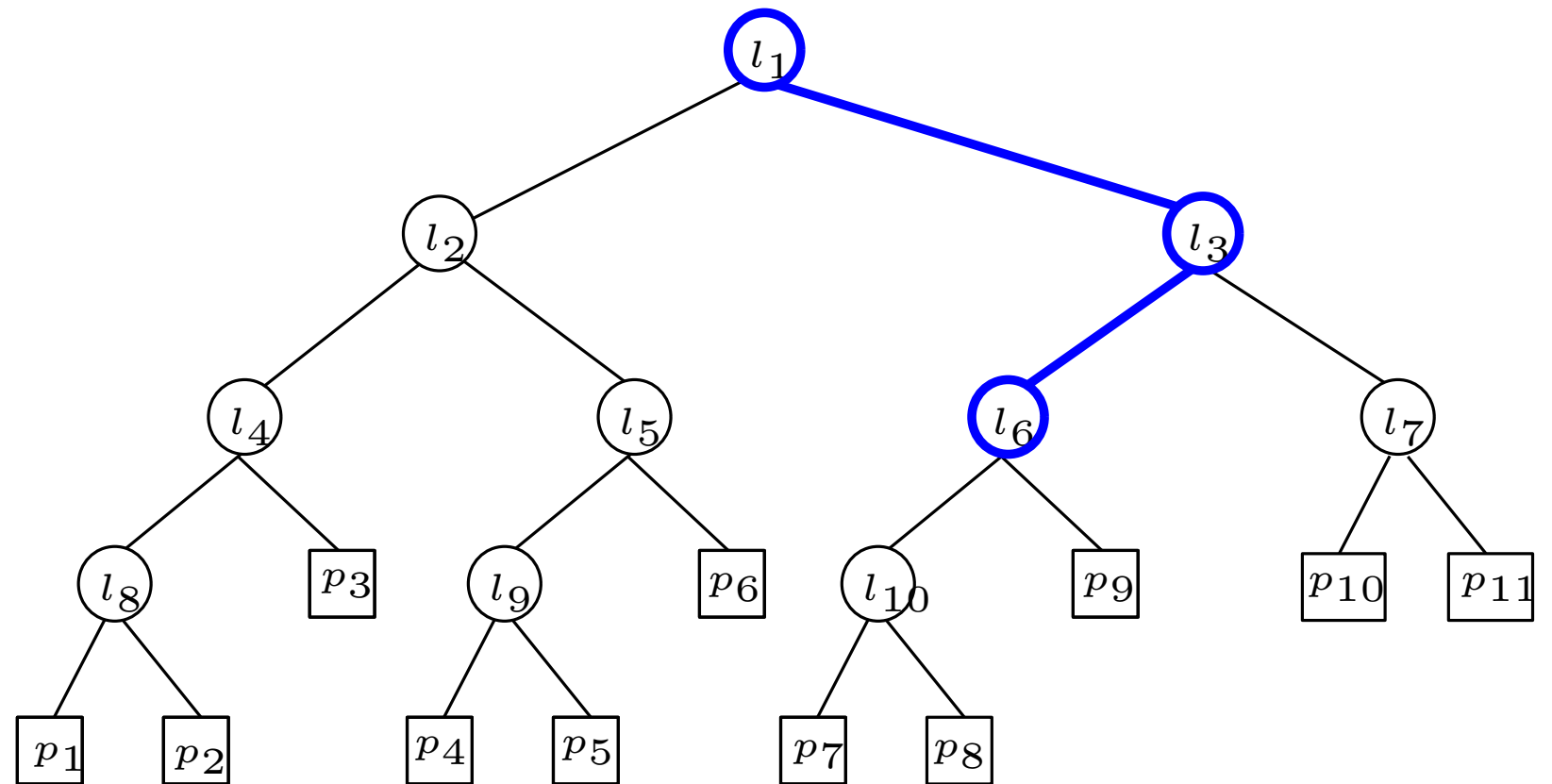
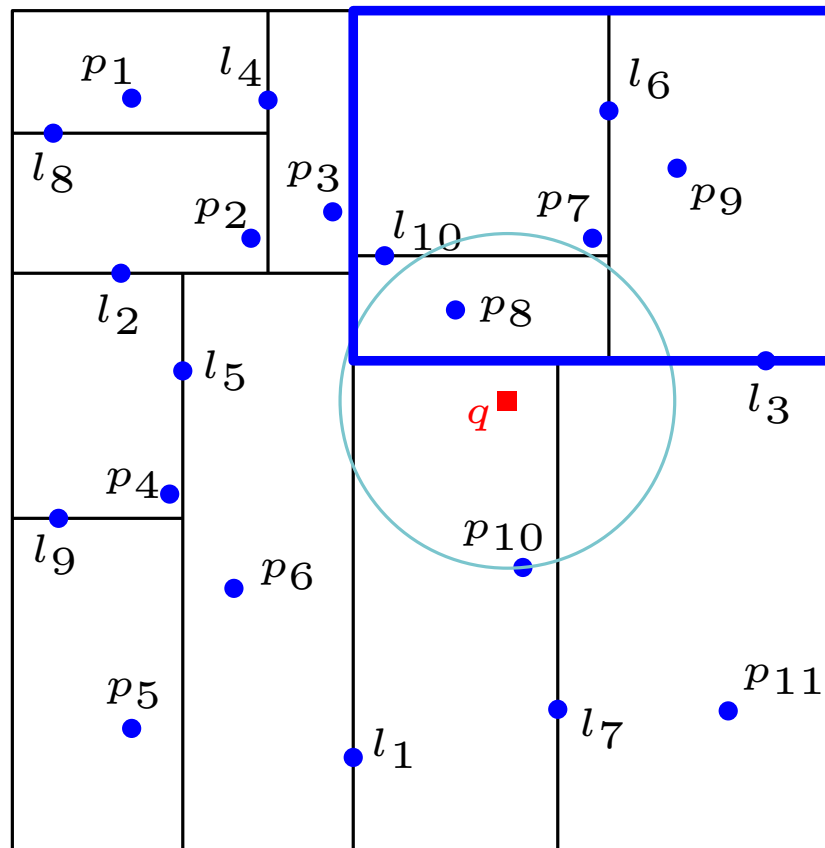


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

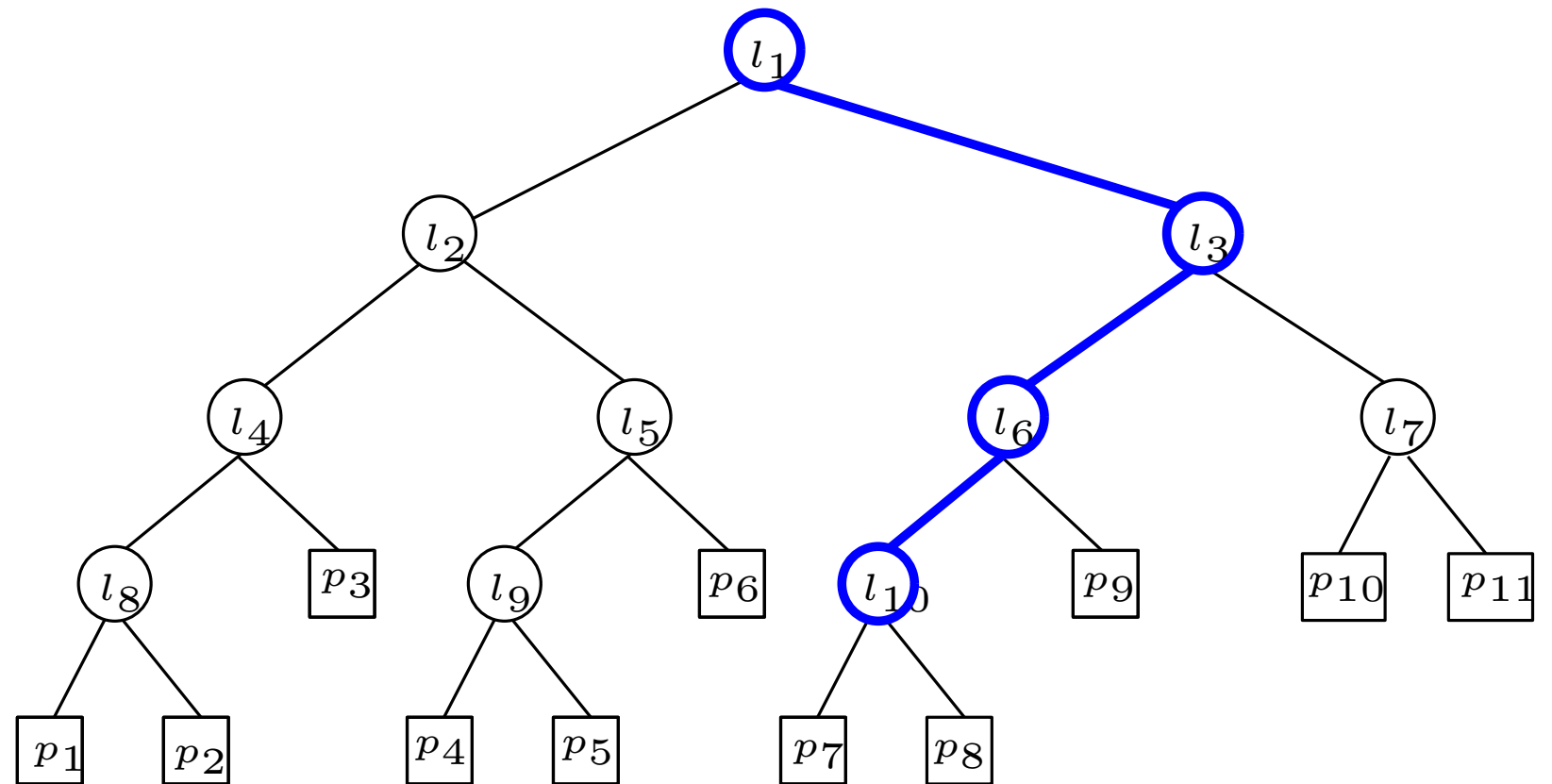
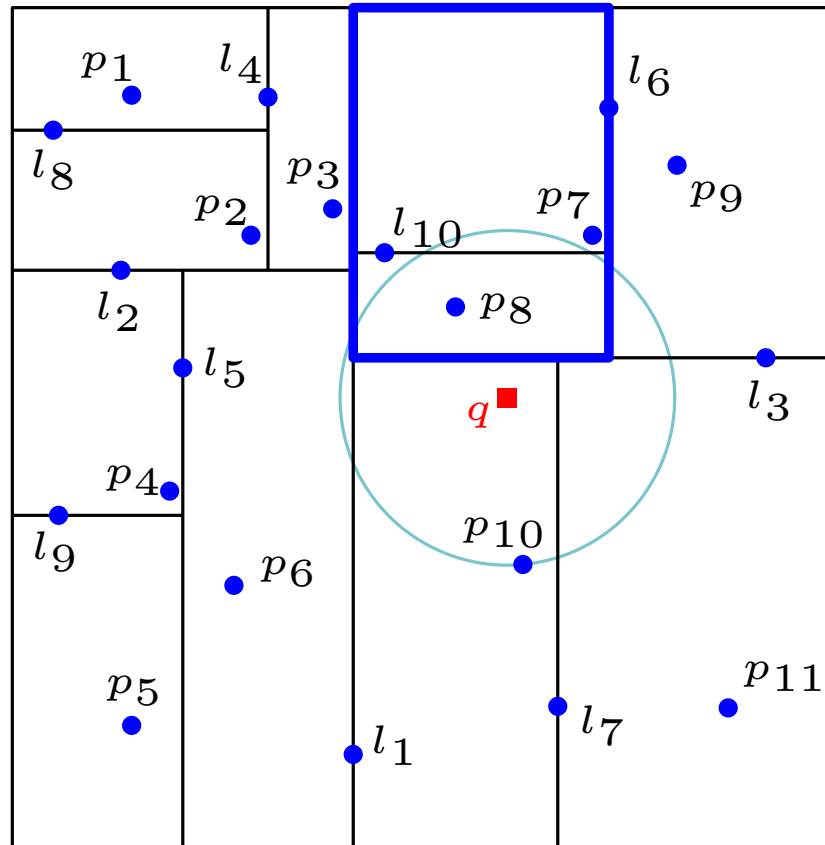


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

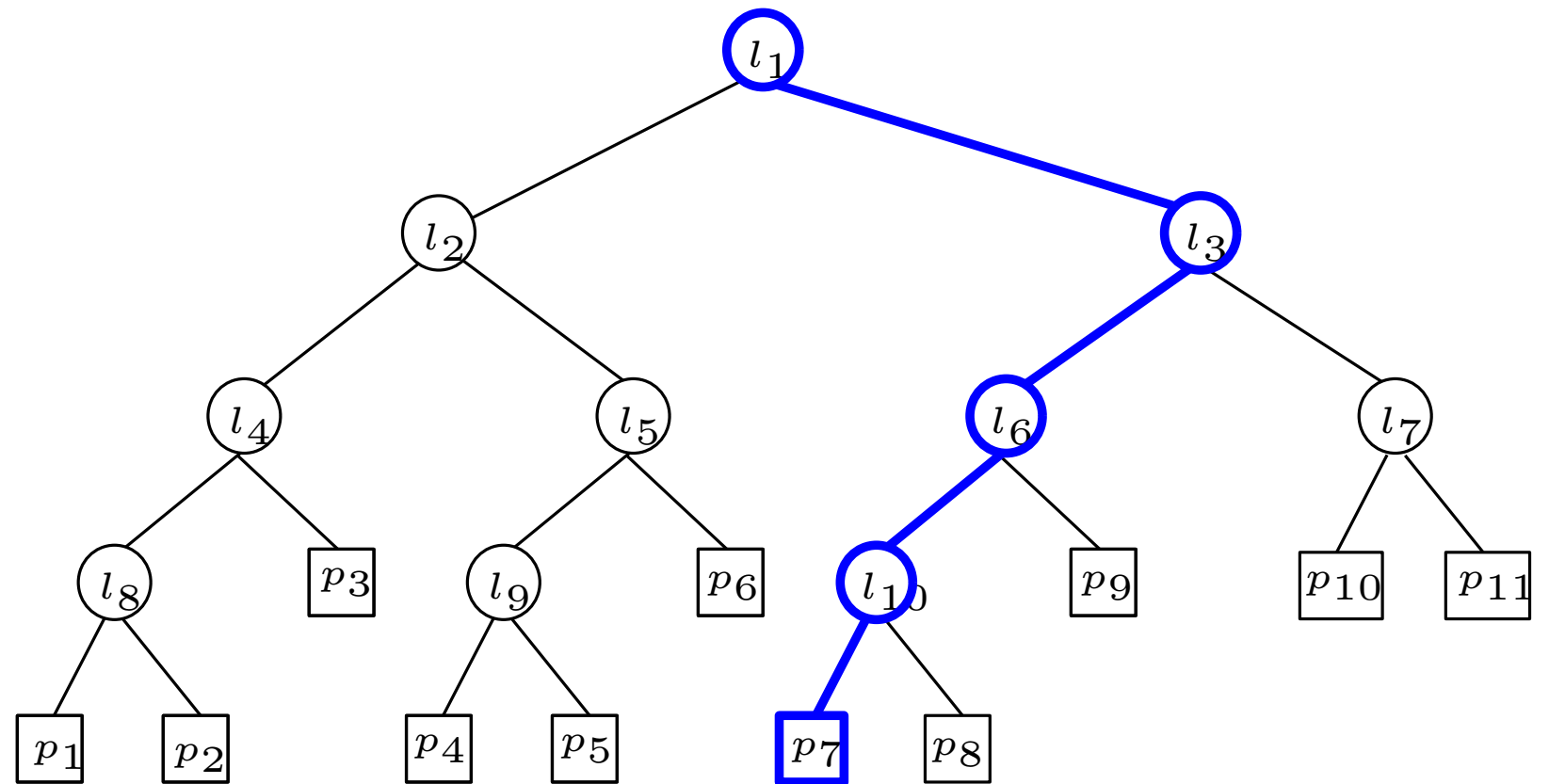
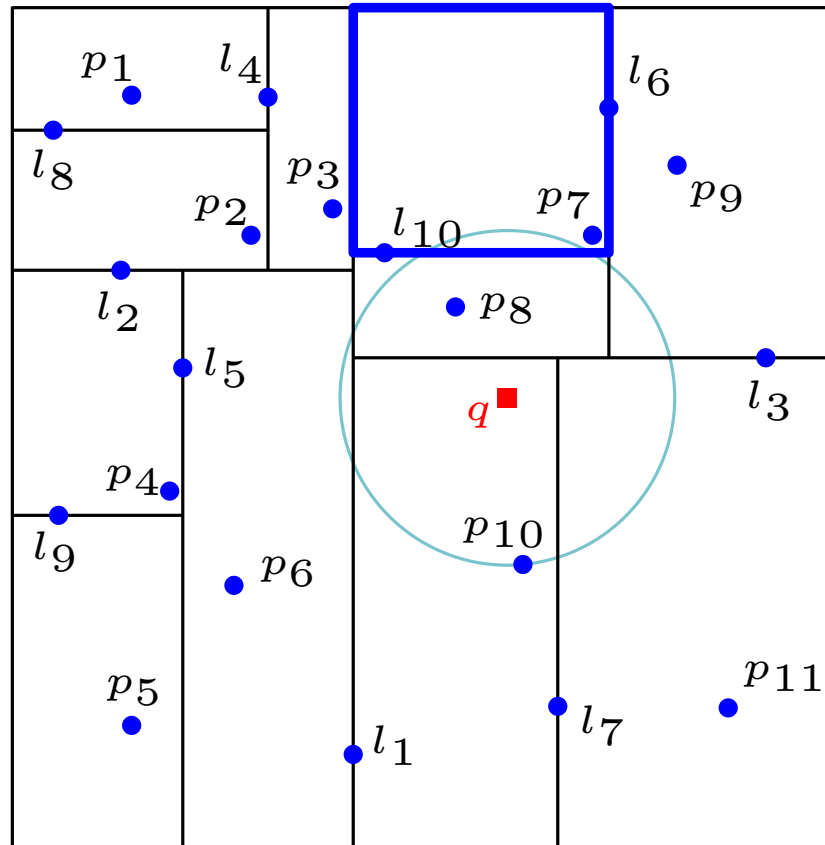


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

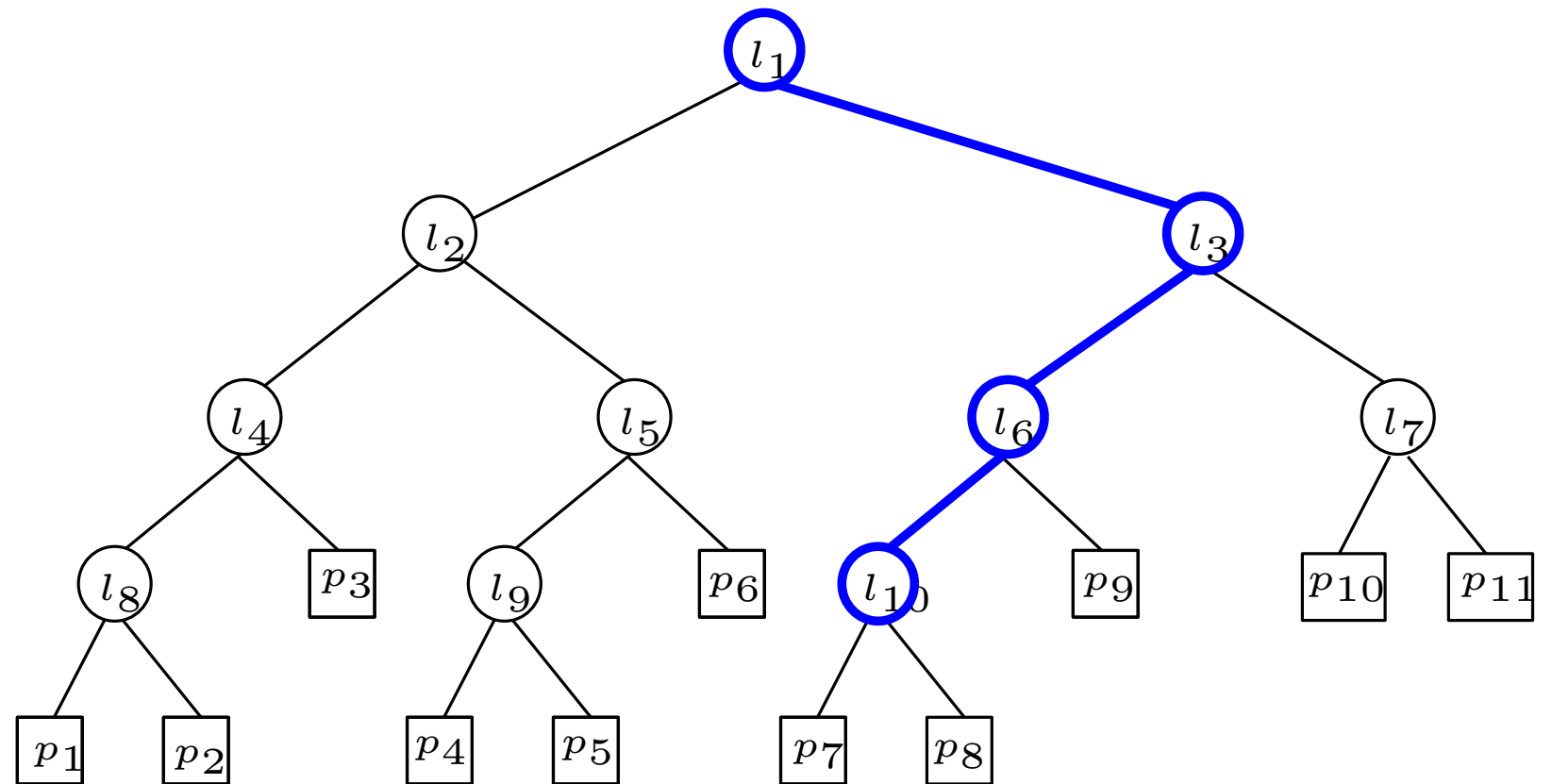
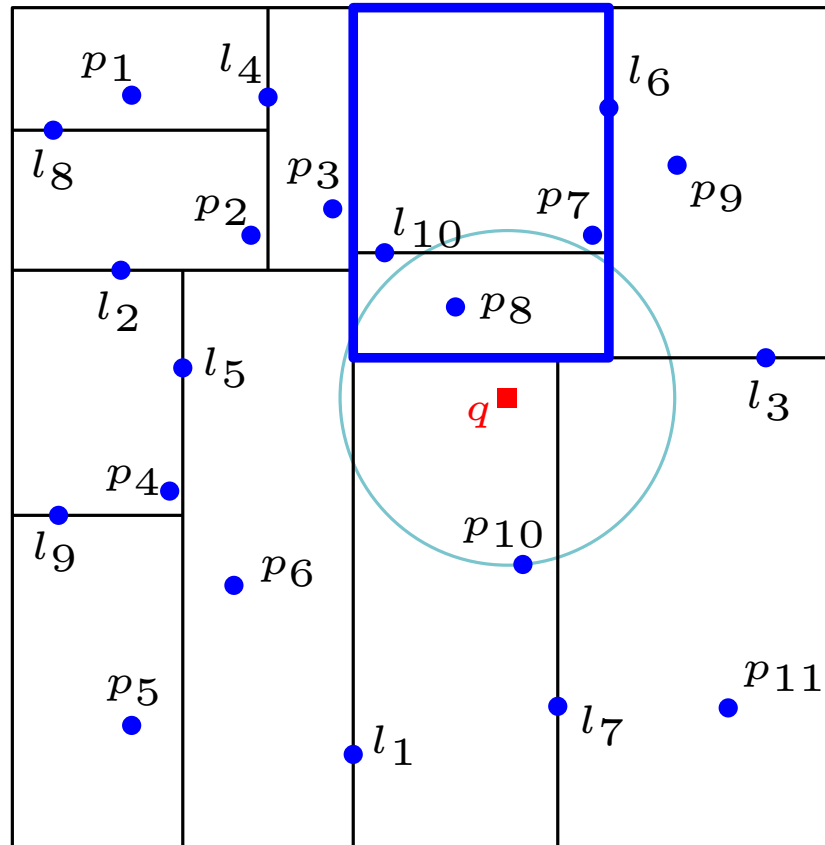


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

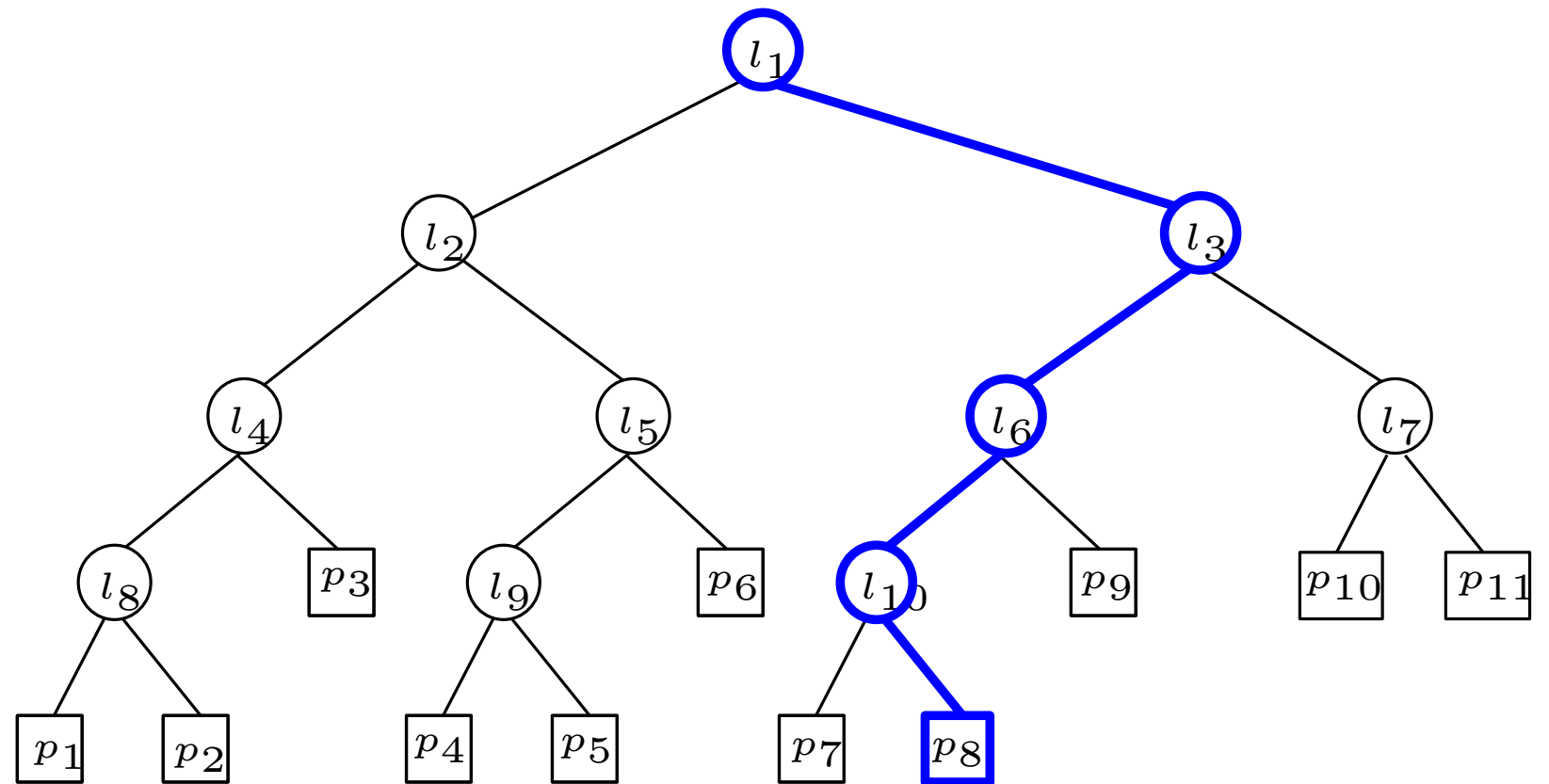
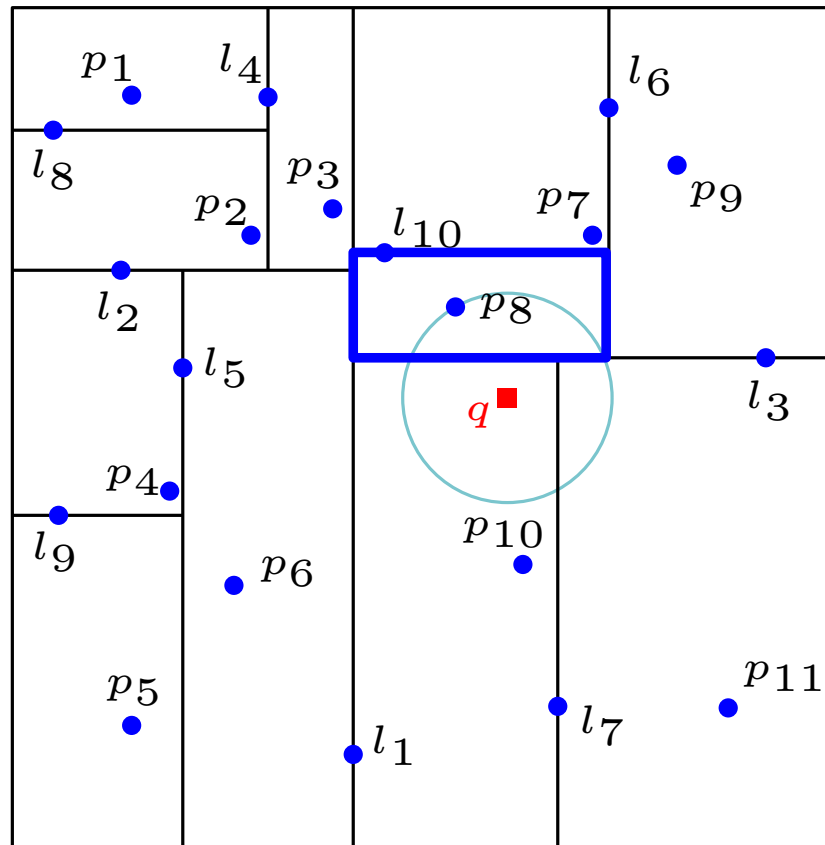


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

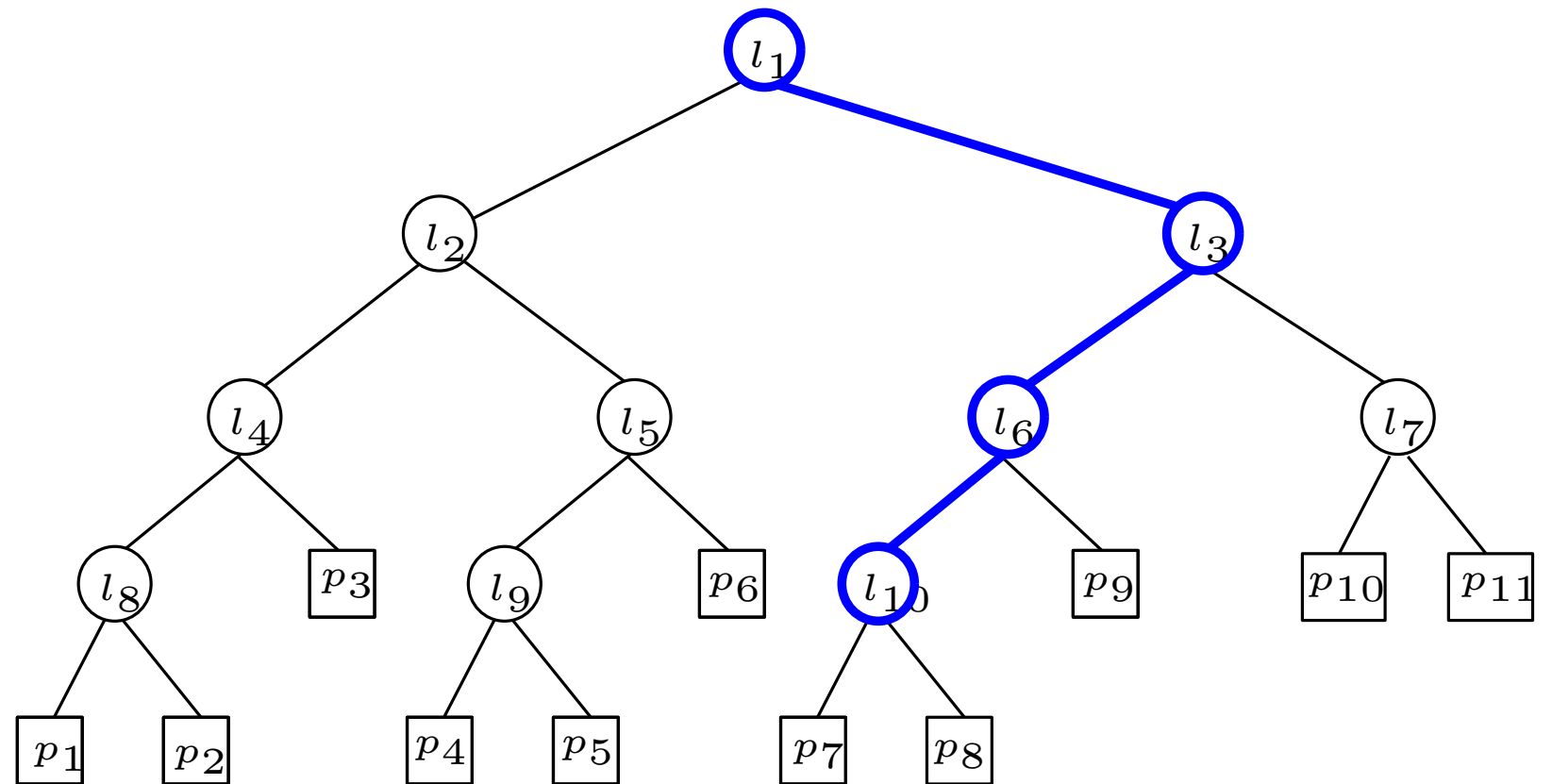
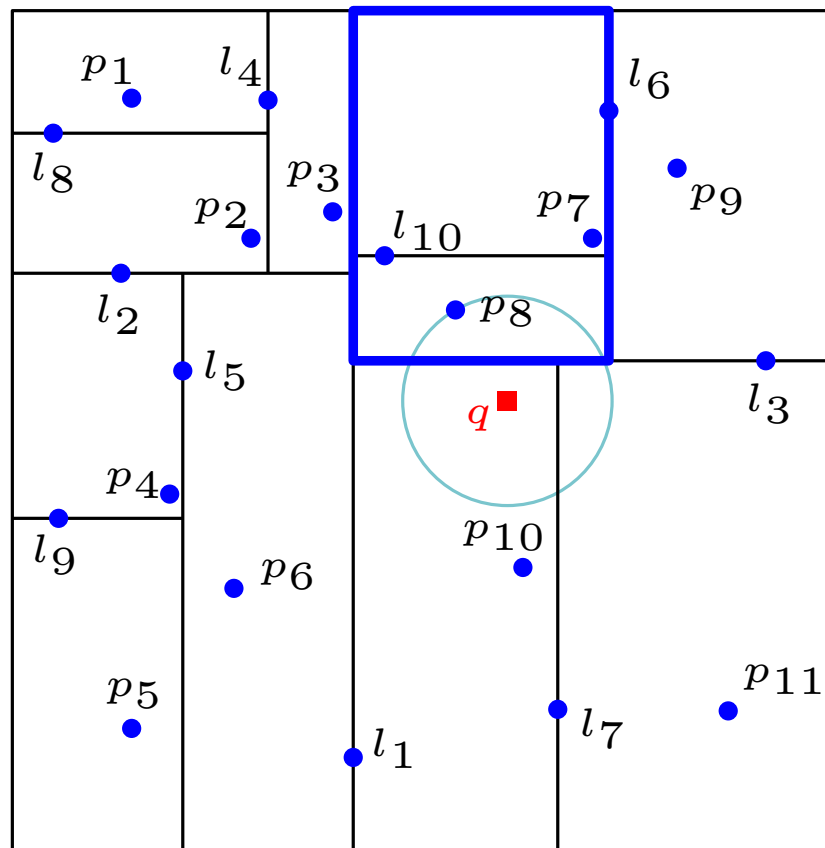


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

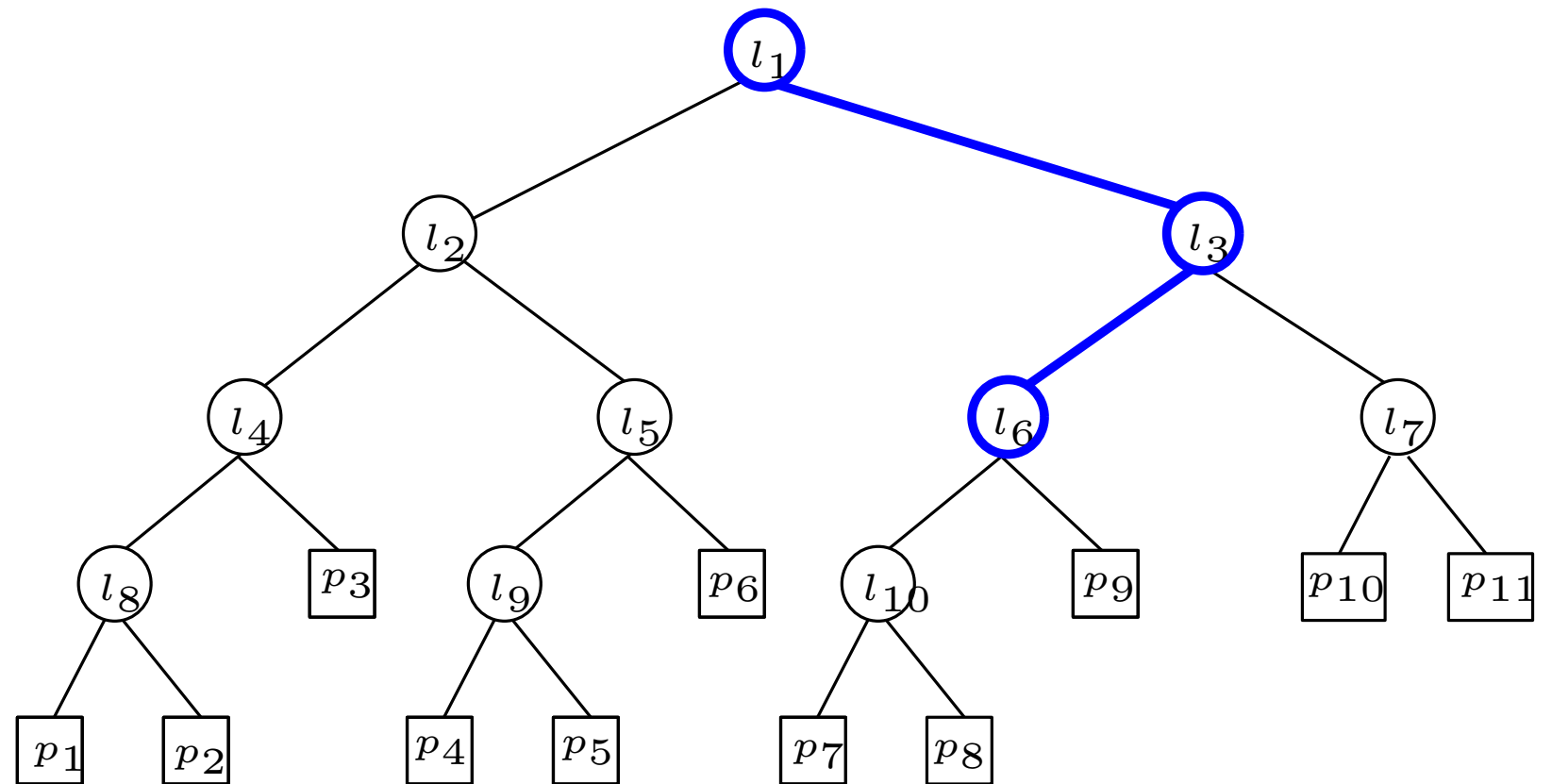
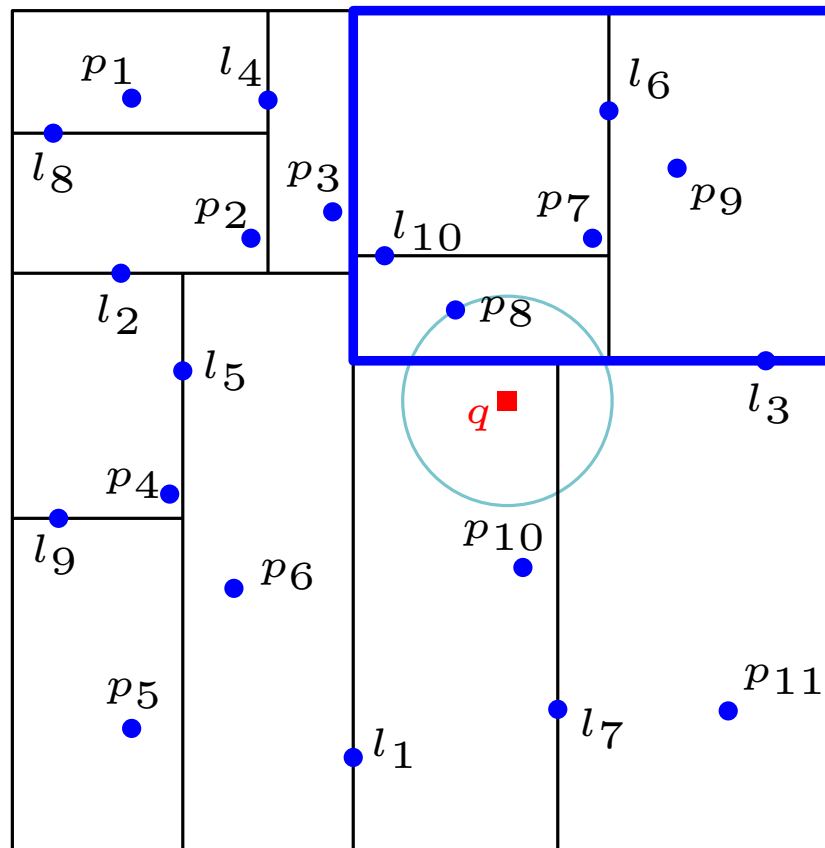


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example



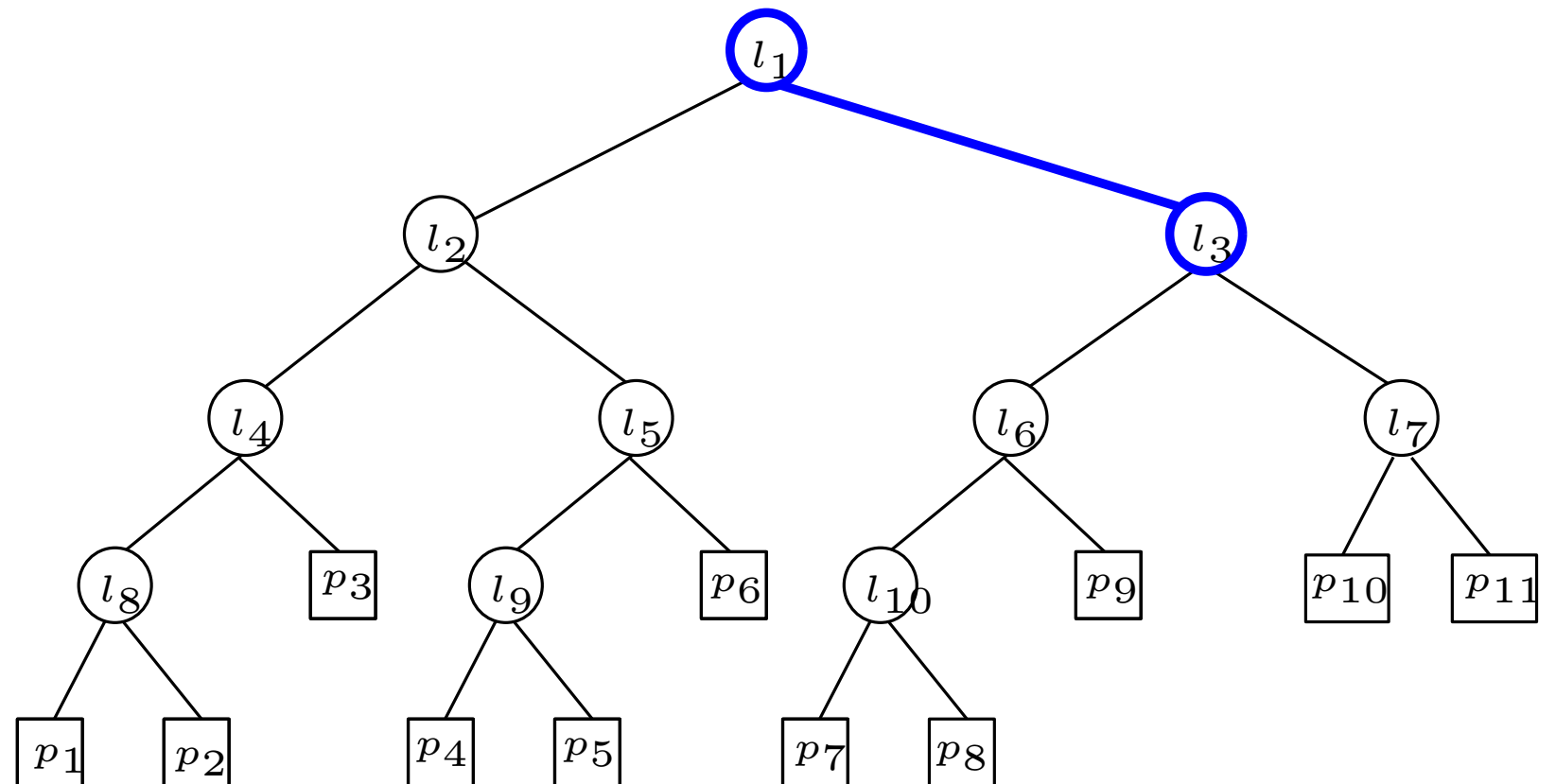
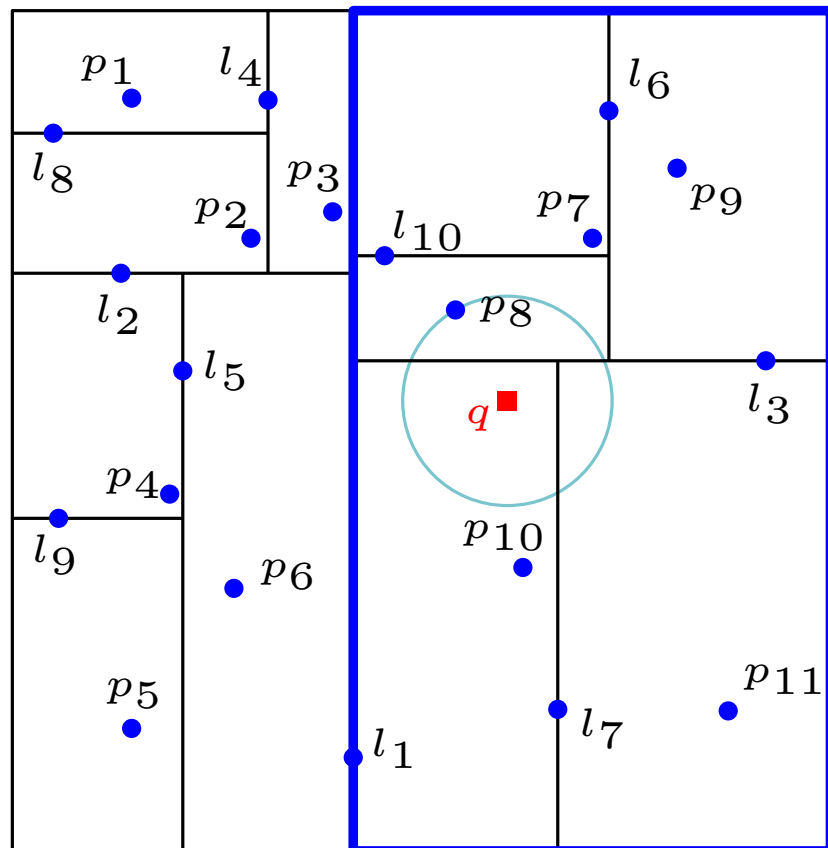
$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)



# Example

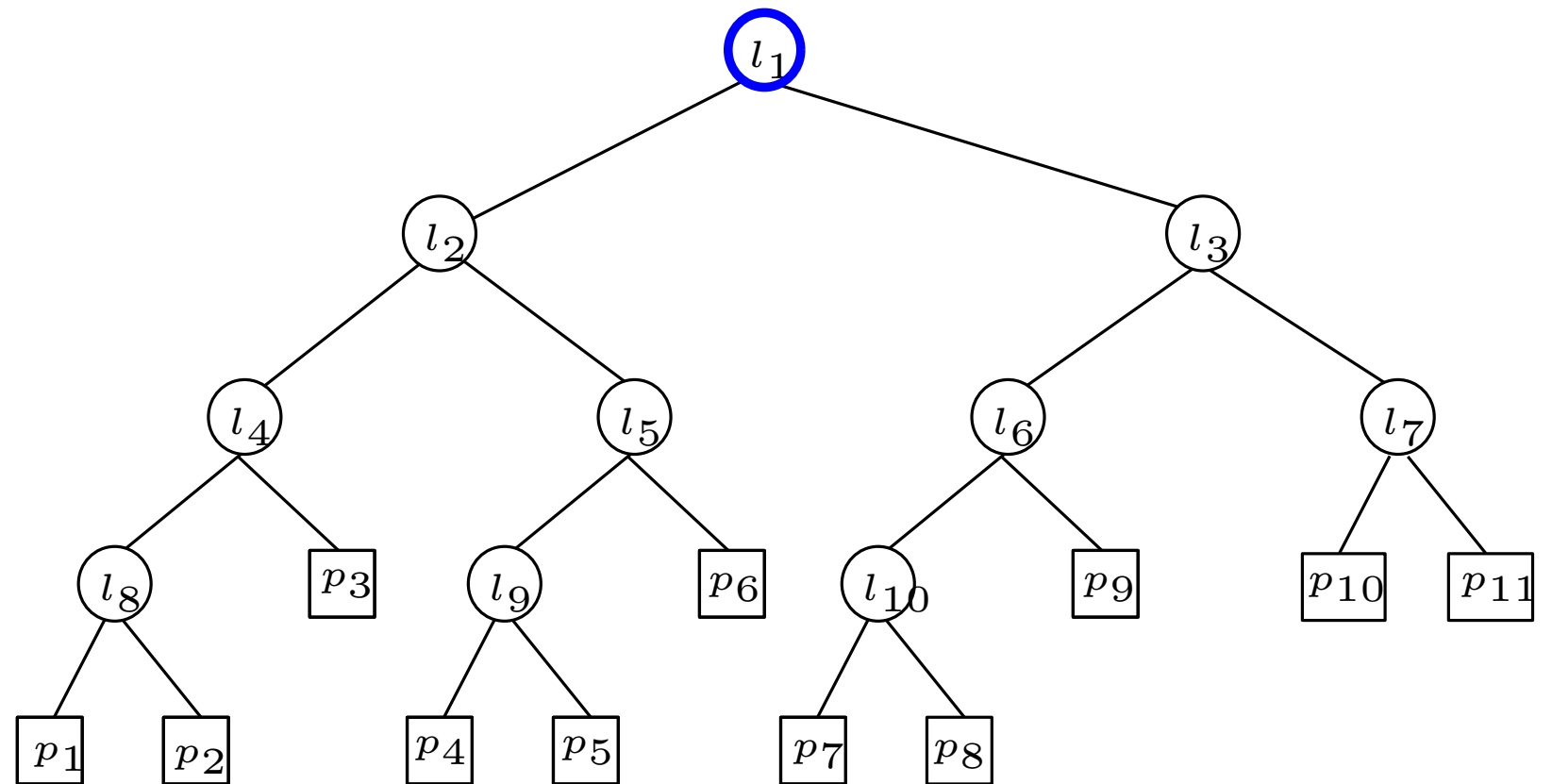
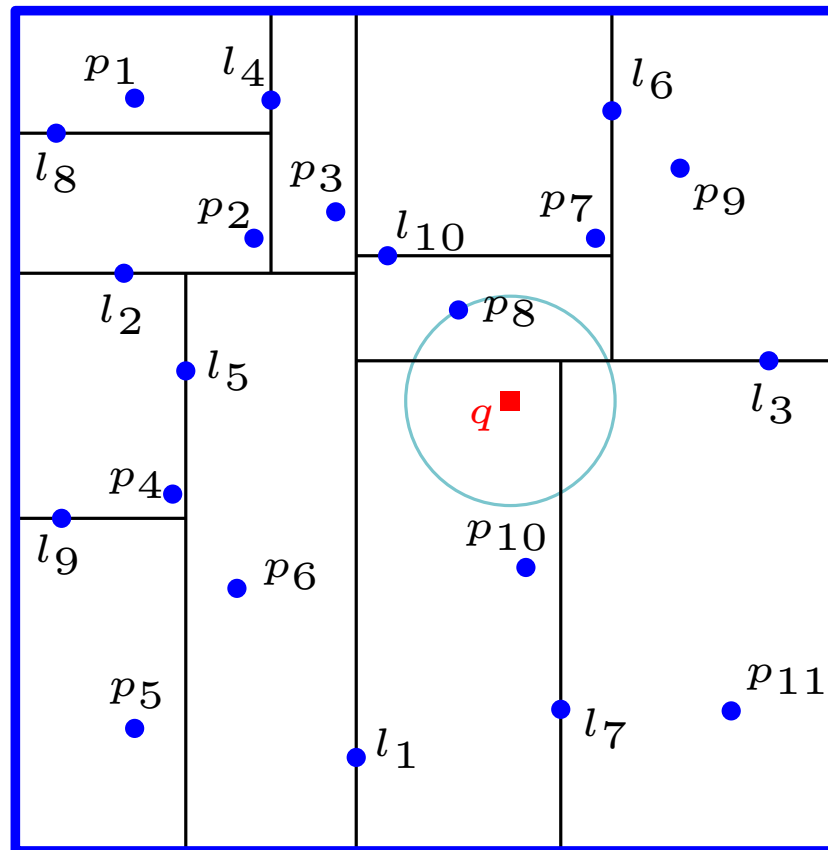


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example

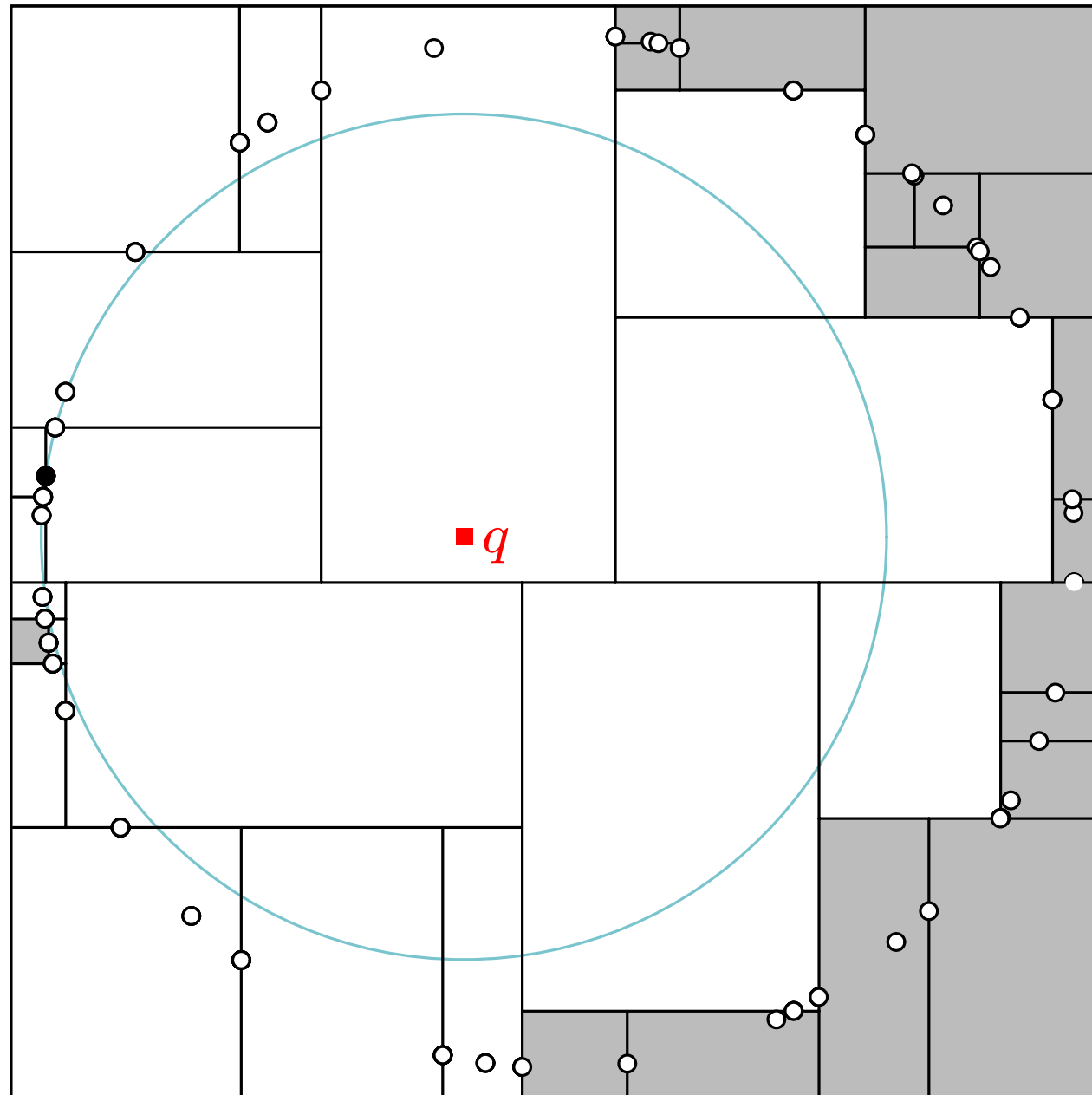


$l_i$ : data at internal node

$p_i$ : data at leaf node

(note: left-right labels are arbitrary)

# Example



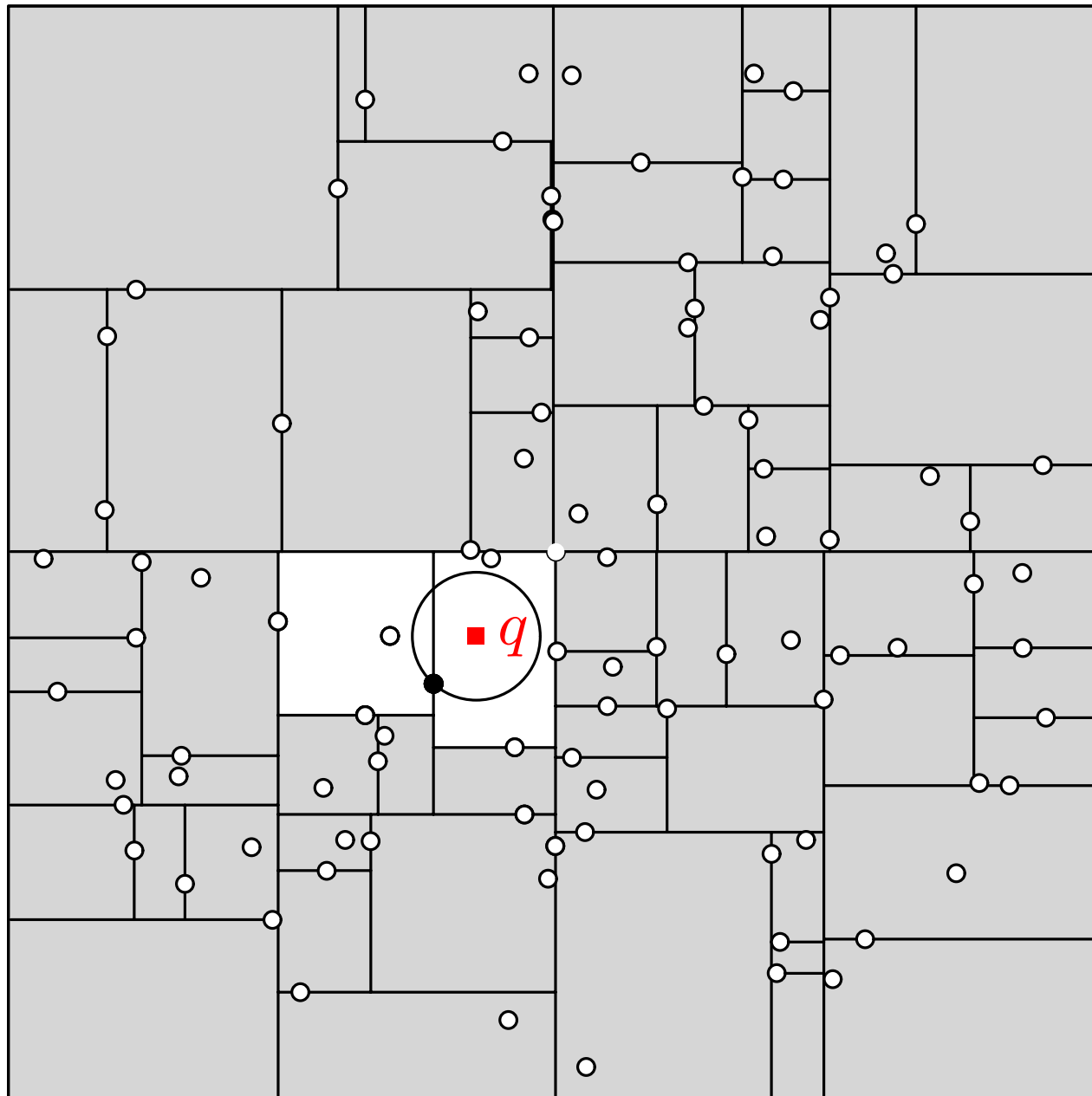
worst-case input (non-unif. distrib.):

long skinny cells



query time =  $\Omega(dn)$

# Example



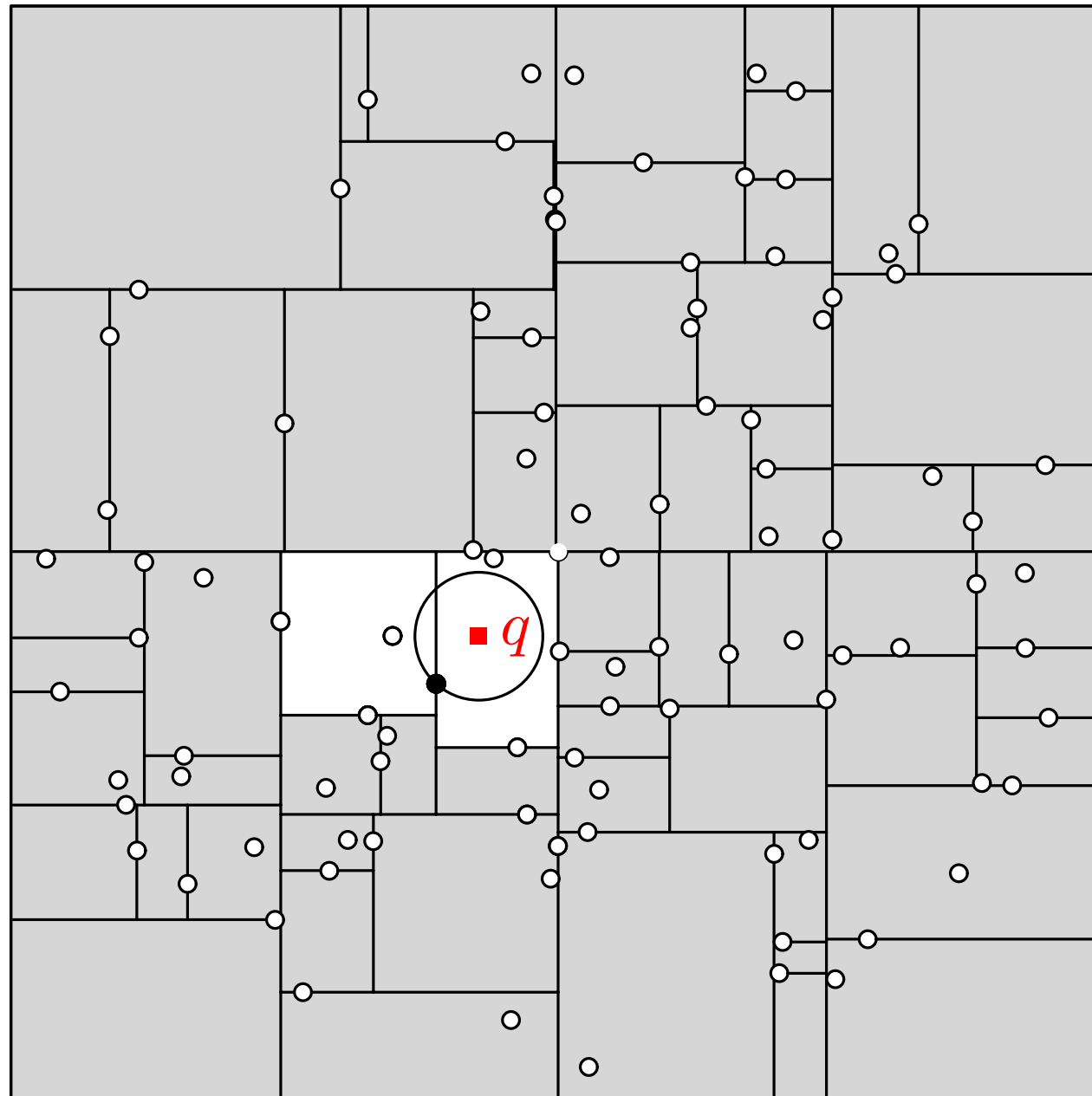
best-case input (unif. distrib.):

small fat cells



query time =  $O(c_d \log n)$

# Example



best-case input (unif. distrib.):

small fat cells



query time =  $O(c_d \log n)$

Randomness should help!

(many variants: priority search, early backtracking, random cutting hyperplanes, etc.)

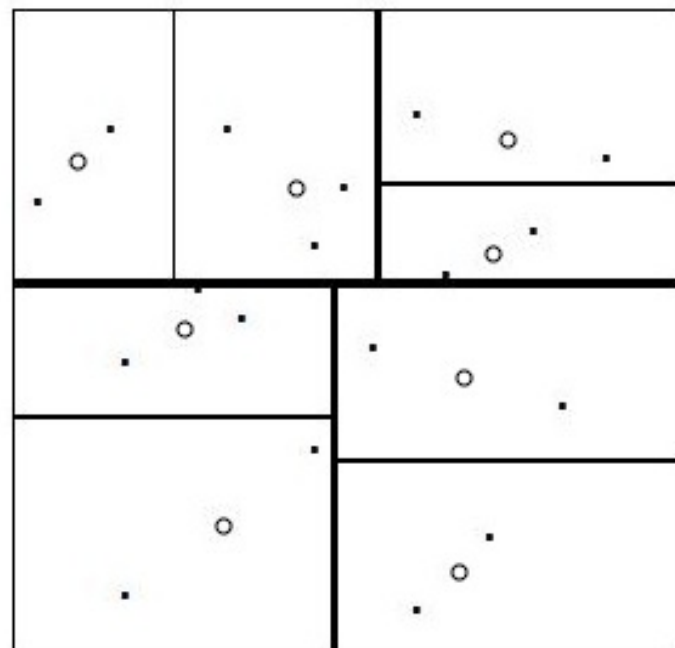
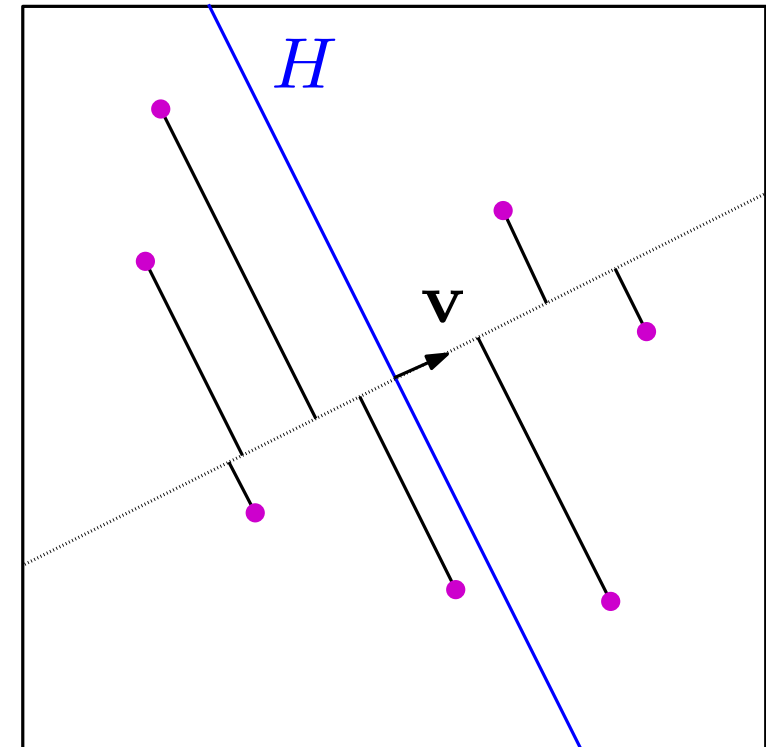
# Random Projection/Partition Trees

# Exploiting randomness: RP-trees

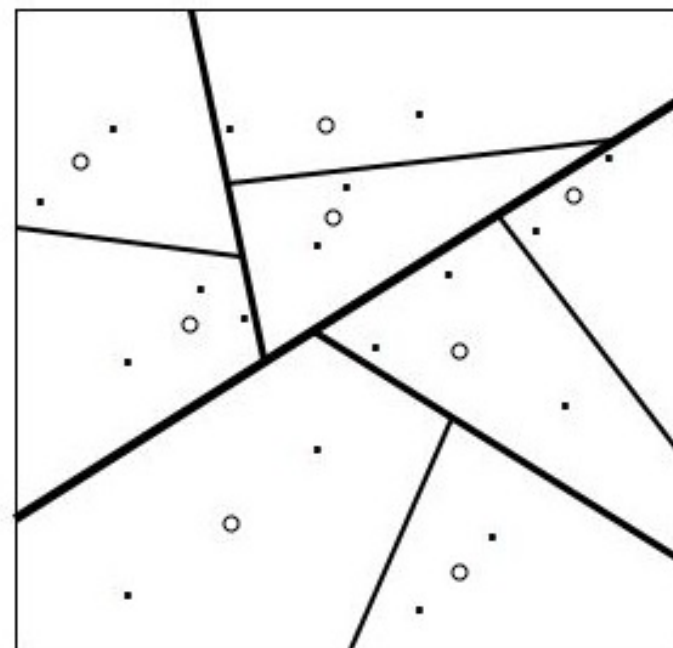
## Random Projection/Partition tree:

at each node (corresponding to some cell  $C$ ):

- choose  $\mathbf{v} \sim \text{unif}(\mathbb{S}^{d-1})$  and  $\beta \sim \text{unif}([\frac{1}{4}, \frac{3}{4}])$
- let  $H = \mathbf{v}^\perp + \text{median}_\beta \{ (P \cap C) \cdot \mathbf{v} \} \mathbf{v}$
- partition  $P \cap C$  by  $H$  (as in kd-tree)



kd-tree

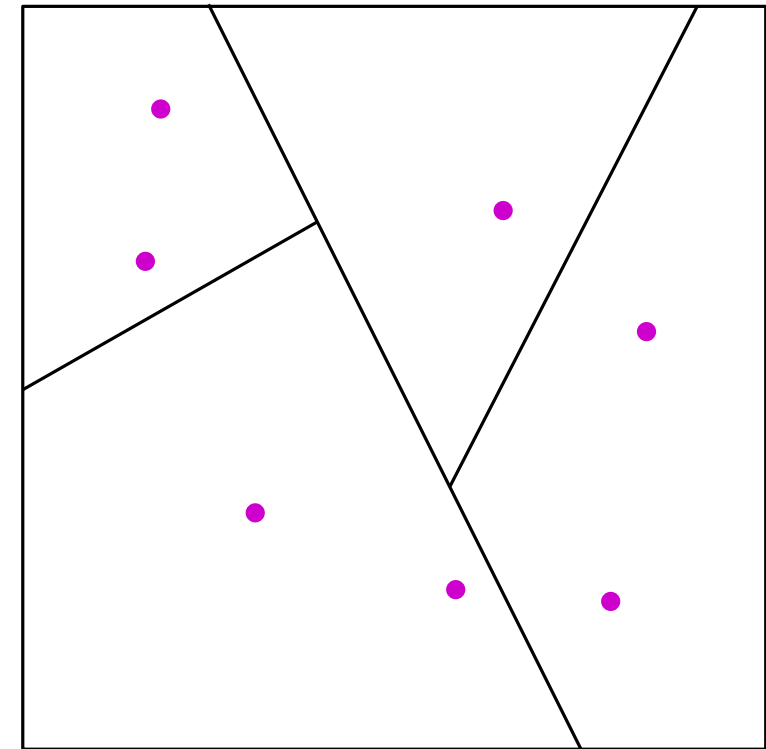


RP-tree

# Exploiting randomness: RP-trees

**Prop:** [Dasgupta, Freund'08]

There is a constant  $c > 0$  such that, for any cell  $C$  in a RP-tree built on  $P \in \mathbb{R}^d$ , with probability at least  $1/2$  (over the choice of  $\mathbf{v}, \beta$ ) all the cells lying at least  $c k \log k$  levels below  $C$  in the tree have at most  $\text{half the radius of } C$ , where  $k = \dim_2(P \cap C)$ .



**doubling dimension** of  $S \subseteq \mathbb{R}^d$ : smallest  $k \in \mathbb{N}$  such that, for every Euclidean ball  $B$ ,  $B \cap S$  can be covered by  $2^k$  Euclidean balls of half radius.

**radius** of  $S \subseteq \mathbb{R}^d$ : smallest  $r > 0$  such that  $\exists x \in C$  with  $B(x, r) \supseteq S$ .



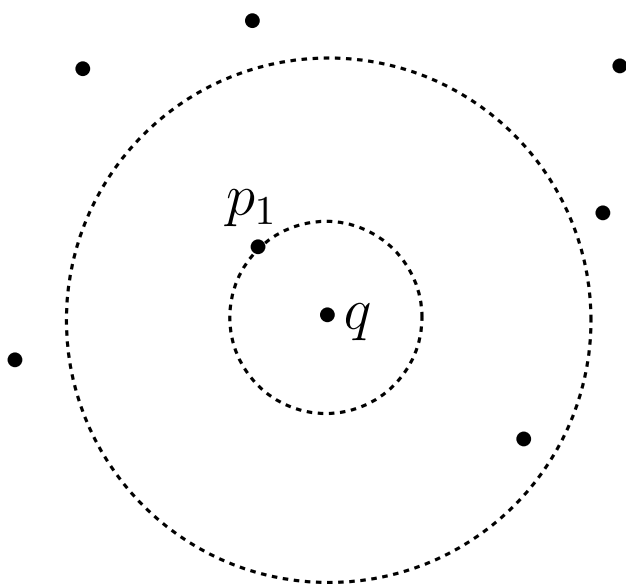
# Exploiting randomness: RP-trees

**Thm:** [Dasgupta, Sinha'13]

Let  $\ell = \log(n/n_0)$  and  $\bar{\beta} = \frac{3}{4}$ .

$$\mathbb{P}_{\mathbf{v}, \beta} [\text{defeatist search does not return } \text{NN}_P(q)] \leq \sum_{i=0}^{\ell} \Phi_{\bar{\beta}^i n} \log \frac{2e}{\Phi_{\bar{\beta}^i n}}$$

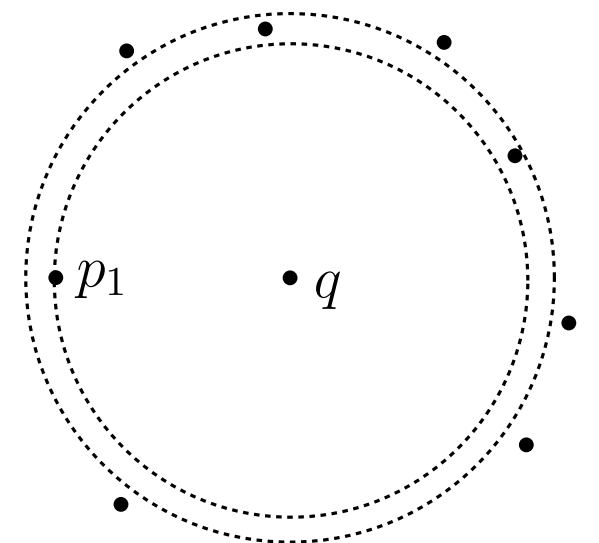
$$\Phi_m := \frac{1}{m} \sum_{i=2}^m \frac{\|q - p_1\|}{\|q - p_i\|}, \text{ where the } p_i \text{ are ordered by increasing distance to } q$$



Extreme cases:

$\Phi \approx 0$ :  $p_1$  isolated, easy to find

$\Phi \approx 1$ :  $p_1$  equidistant, hard to find



# Exploiting randomness: RP-trees

**Thm:** [Dasgupta, Sinha'13]

Suppose  $p_1, \dots, p_n \stackrel{\text{iid}}{\sim} \mu$  continuous probability measure in  $\mathbb{R}^d$  with doubling dimension  $k \geq 2$ . Then  $\exists c_0 > 0$  s.t. for any  $q \in \mathbb{R}^d$  and  $\delta < 1/e$ , with proba.  $\geq 1 - 3\delta$  over the choice of the  $p_i$ 's:

$$\mathbb{P}_{\mathbf{v}, \beta} [\text{defeatist search does not return } \text{NN}_P(q)] \leq c_0(k + \ln n_0) \left( \frac{8 \ln 1/\delta}{n_0} \right)^{1/k}$$

**doubling dimension** of  $\mu$ : smallest  $k \in \mathbb{N}$  such that, for every  $x \in \mathbb{R}^d$  and every  $r > 0$ :  $\mu(B(x, 2r)) \leq 2^k \mu(B(x, r))$ .

# Exploiting randomness: RP-trees

**Thm:** [Dasgupta, Sinha'13]

Suppose  $p_1, \dots, p_n \stackrel{\text{iid}}{\sim} \mu$  continuous probability measure in  $\mathbb{R}^d$  with doubling dimension  $k \geq 2$ . Then  $\exists c_0 > 0$  s.t. for any  $q \in \mathbb{R}^d$  and  $\delta < 1/e$ , with proba.  $\geq 1 - 3\delta$  over the choice of the  $p_i$ 's:

$$\mathbb{P}_{\mathbf{v}, \beta} [\text{defeatist search does not return } \text{NN}_P(q)] \leq c_0(k + \ln n_0) \left( \frac{8 \ln 1/\delta}{n_0} \right)^{1/k}$$

**doubling dimension** of  $\mu$ : smallest  $k \in \mathbb{N}$  such that, for every  $x \in \mathbb{R}^d$  and every  $r > 0$ :  $\mu(B(x, 2r)) \leq 2^k \mu(B(x, r))$ .

→ take  $n_0 \propto (k \ln k)^k \ln 1/\delta$  to make  $\mathbb{P}_{\mathbf{v}, \beta} [\dots]$  an arbitrarily small constant

→ query time:  $O(d((k \ln k)^k + \log n))$  ← sensitive to intrinsic dim.  
requires to know  $k$

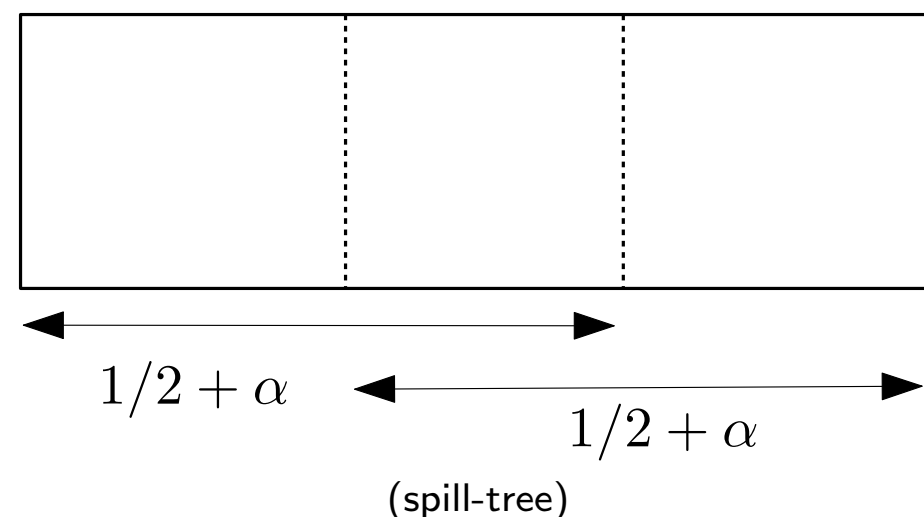
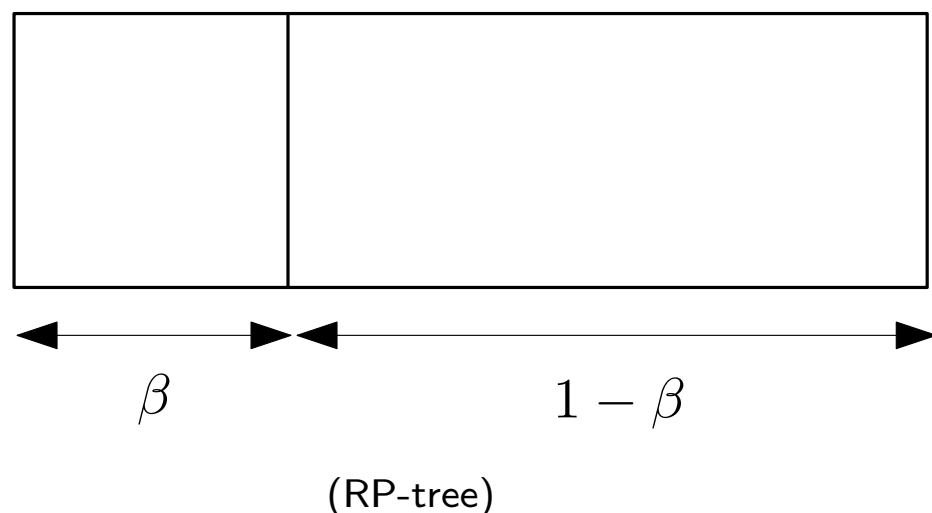
# Exploiting randomness: RP-trees

**Thm:** [Dasgupta, Sinha'13]

Suppose  $p_1, \dots, p_n \stackrel{\text{iid}}{\sim} \mu$  continuous probability measure in  $\mathbb{R}^d$  with doubling dimension  $k \geq 2$ . Then  $\exists c_0 > 0$  s.t. for any  $q \in \mathbb{R}^d$  and  $\delta < 1/e$ , with proba.  $\geq 1 - 3\delta$  over the choice of the  $p_i$ 's:

$$\mathbb{P}_{\mathbf{v}, \beta} [\text{defeatist search does not return } \text{NN}_P(q)] \leq c_0(k + \ln n_0) \left( \frac{8 \ln 1/\delta}{n_0} \right)^{1/k}$$

Variant: **spill-trees** (overlapping splits)



# Exploiting randomness: RP-trees

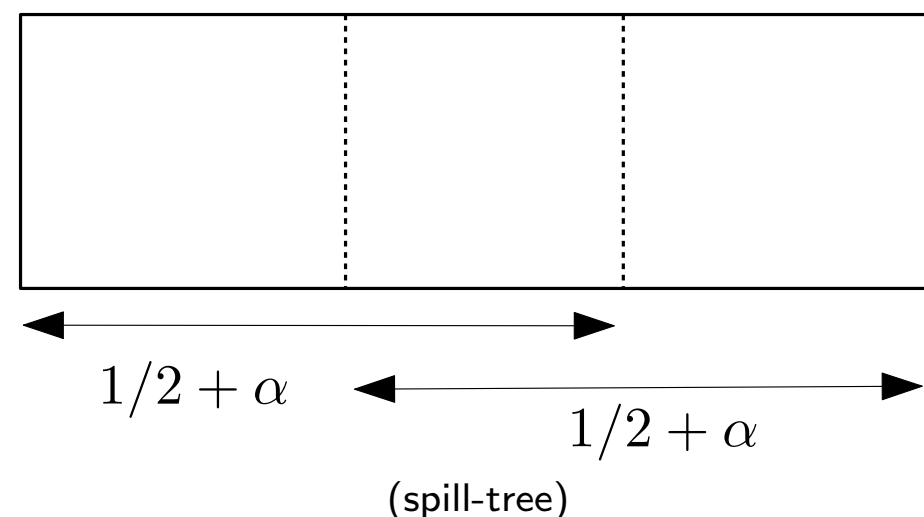
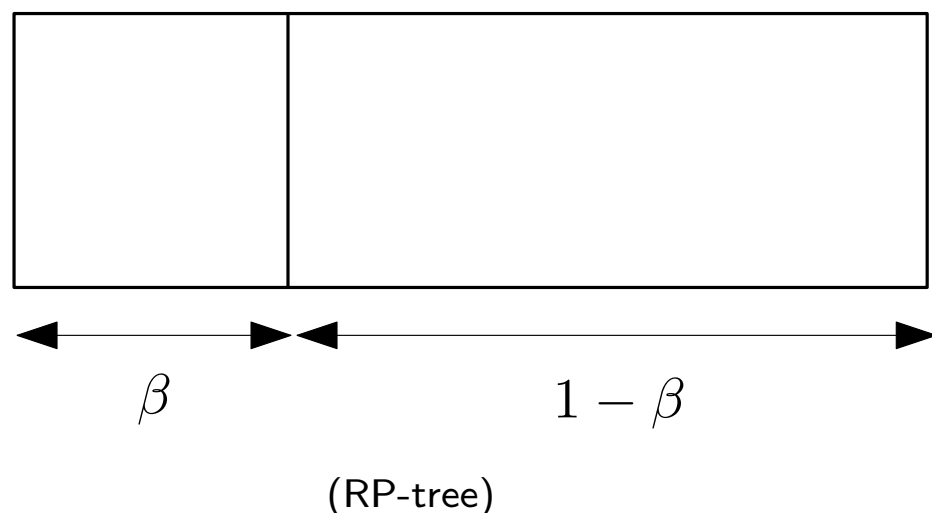
**Thm:** [Dasgupta, Sinha'13]

Suppose  $p_1, \dots, p_n \stackrel{\text{iid}}{\sim} \mu$  continuous probability measure in  $\mathbb{R}^d$  with doubling dimension  $k \geq 2$ . Then  $\exists c_0 > 0$  s.t. for any  $q \in \mathbb{R}^d$  and  $\delta < 1/e$ , with proba.  $\geq 1 - 3\delta$  over the choice of the  $p_i$ 's:

$$\mathbb{P}_{\mathbf{v}, \beta} [\text{defeatist search does not return } \text{NN}_P(q)] \leq c_0(k + \ln n_0) \left( \frac{8 \ln 1/\delta}{n_0} \right)^{1/k}$$

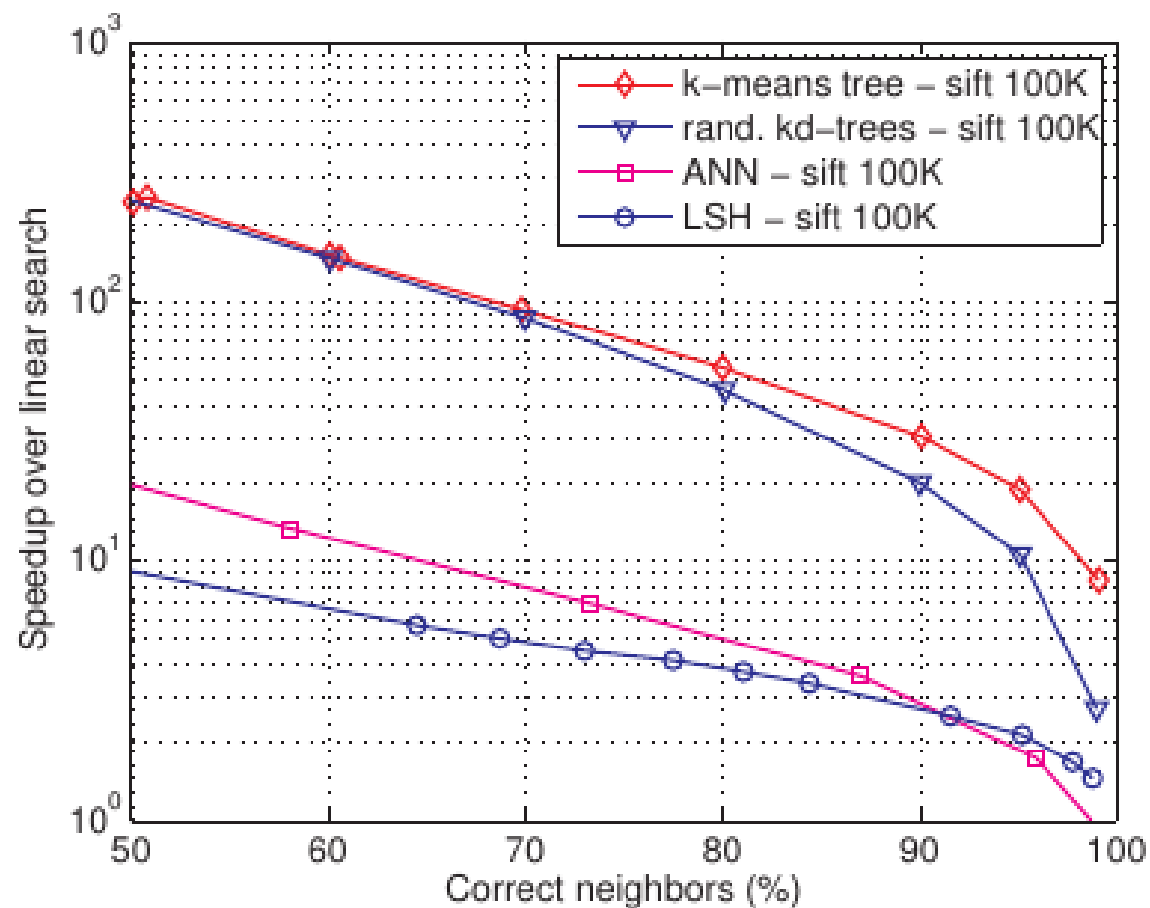
Variant: **spill-trees** (overlapping splits)

similar behavior



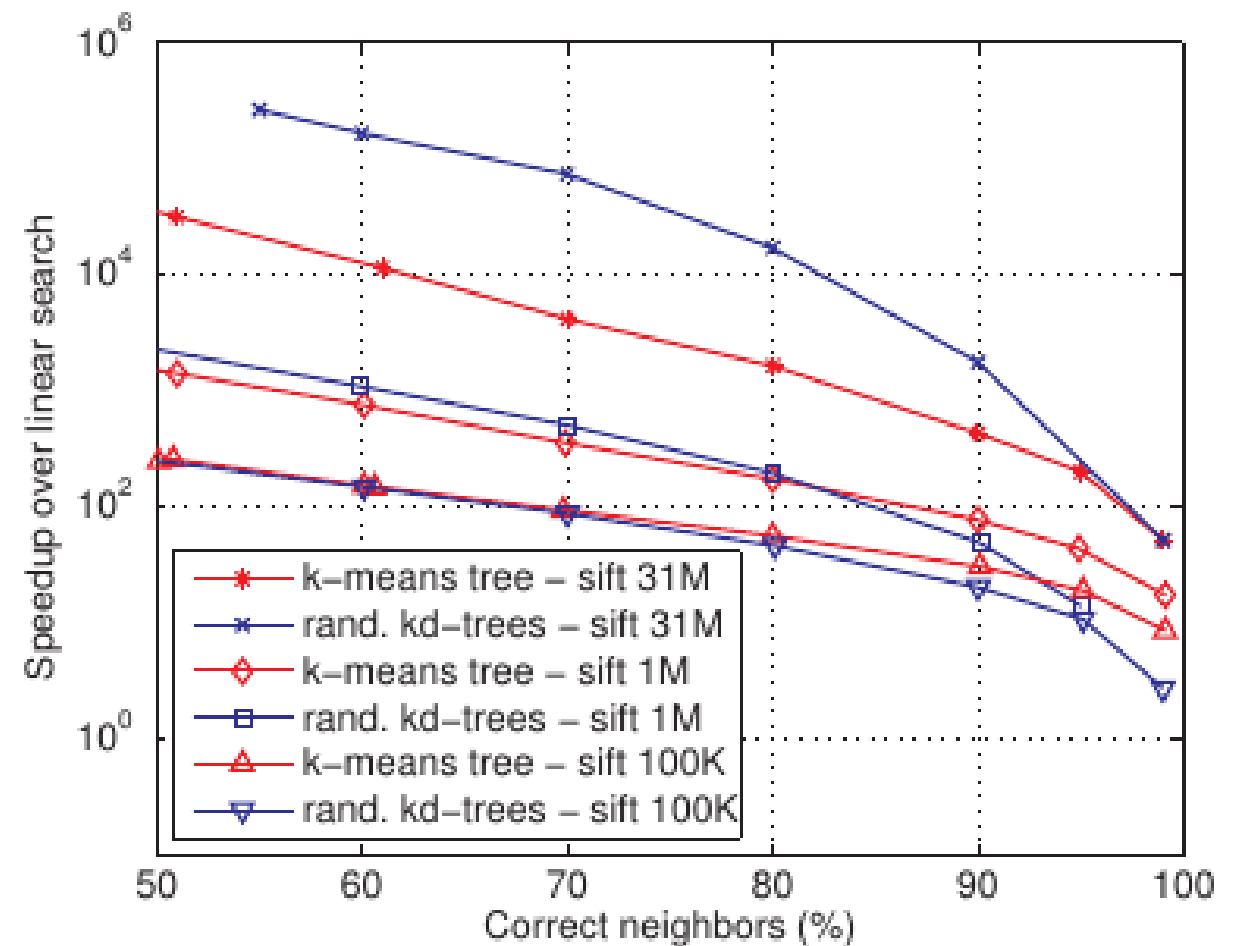
# Benchmarking

contenders



(a)

effect of size on winners

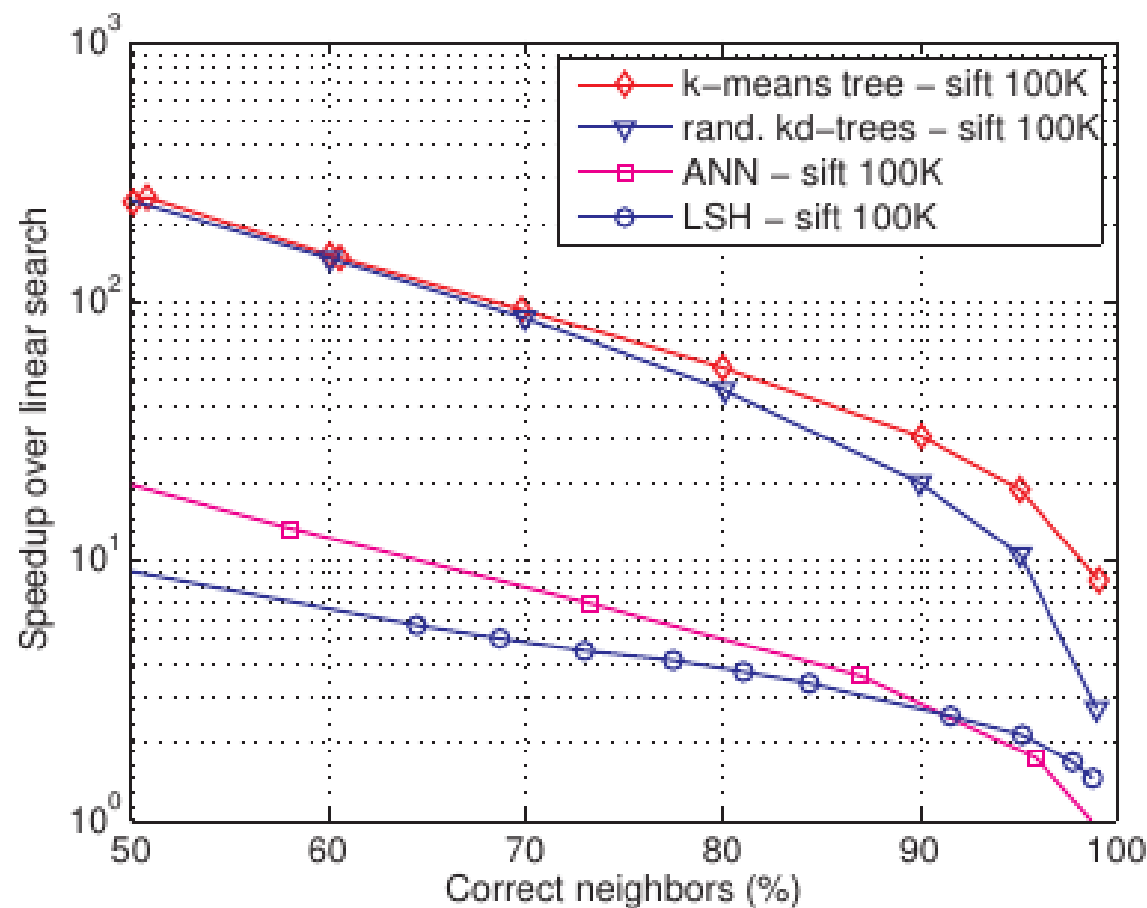


(b)

RP-trees vs. other methods on data sets of 100k, 1M and 31M features

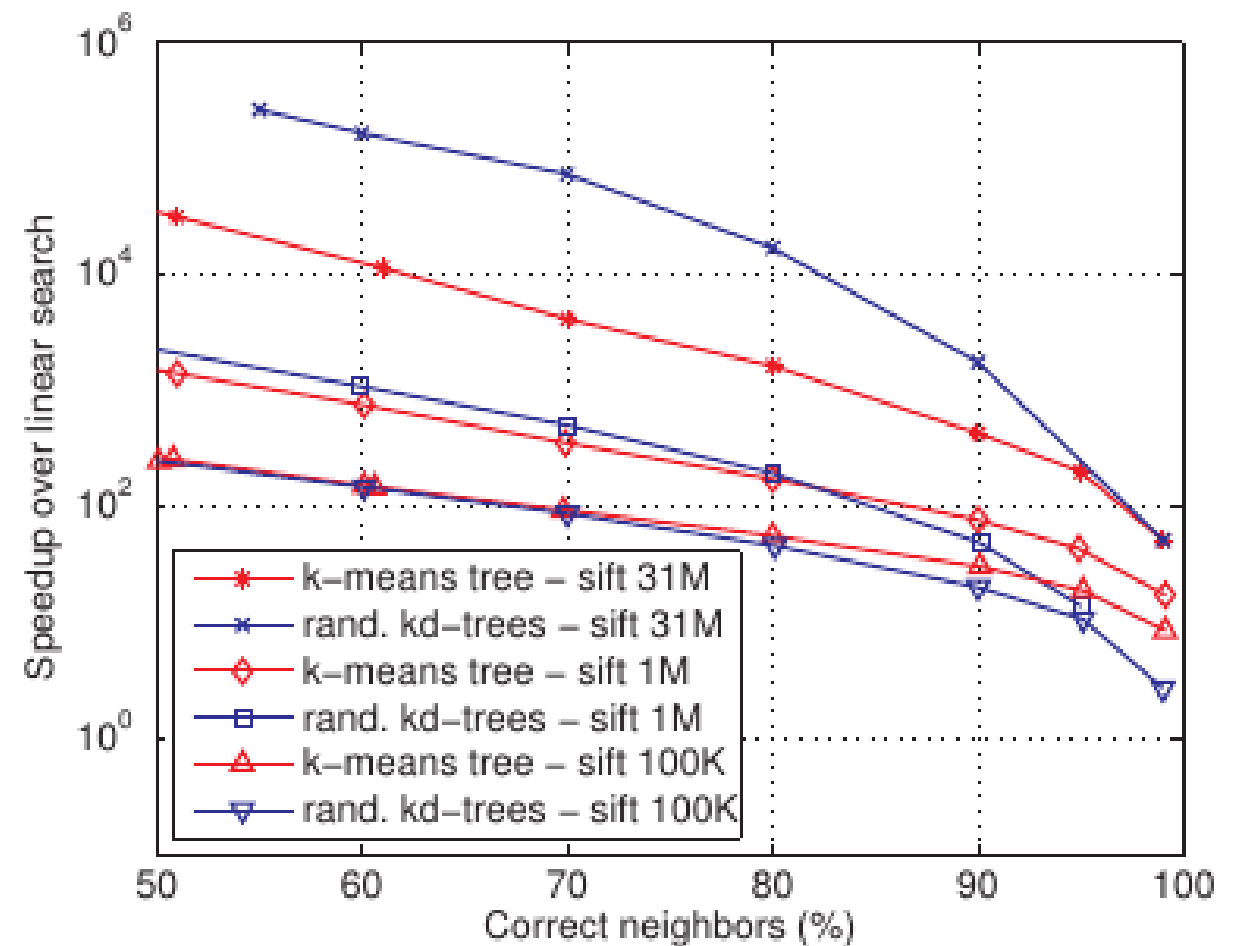
# Benchmarking

contenders



(a)

effect of size on winners



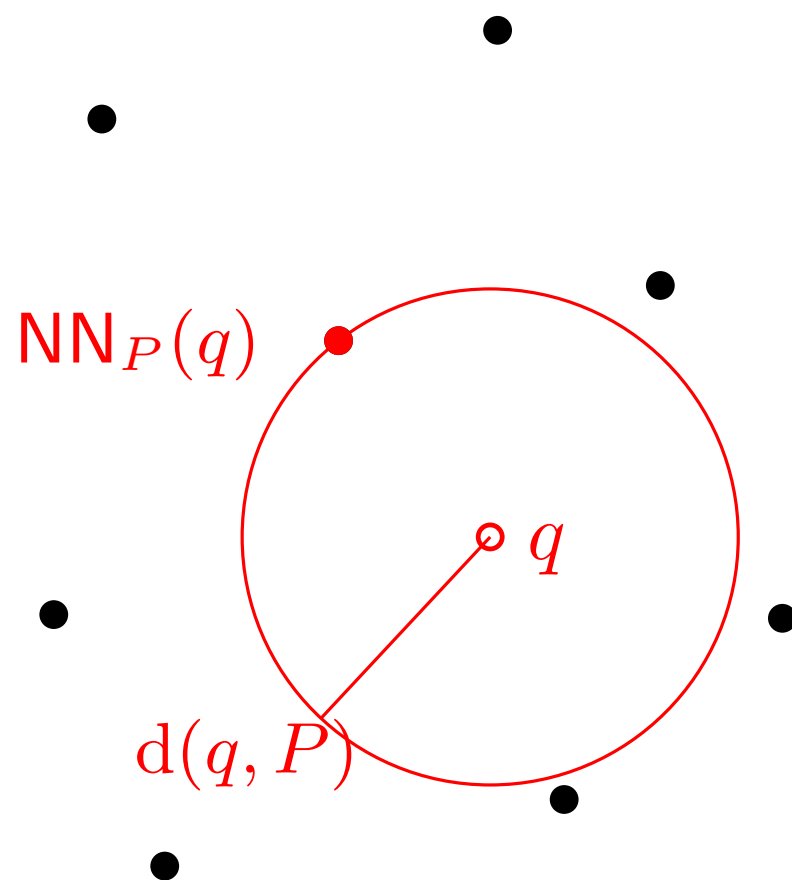
(b)

Random kd-trees (RP-trees, spill-trees) are fast, scalable and reliable on data with (low-dimensional) intrinsic structure

# Locality-Sensitive Hashing



# Back to the NN problem



pre-processing input:  $P$

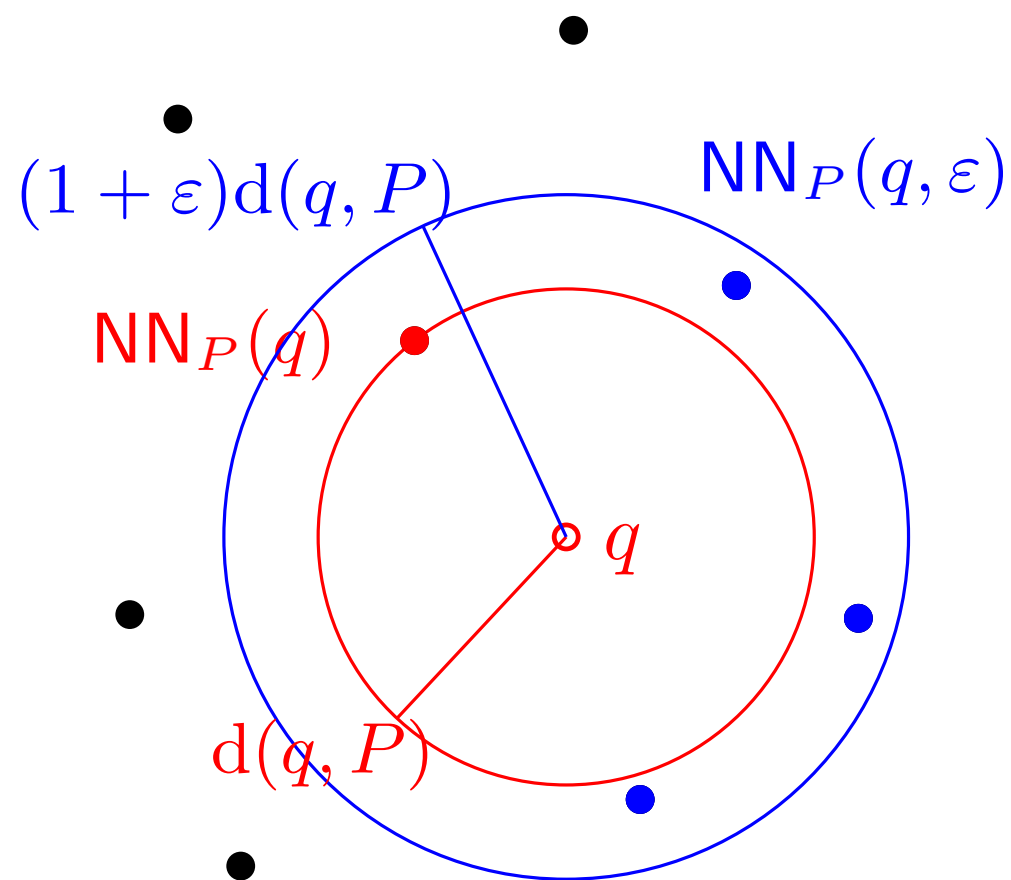
query input:  $q$

goal: find  $p \in NN_P(q)$

**Curse of Dimensionality:** every DS for NN-search has either exponential size or exponential query time (in  $d$ ) in the worst case.

→ holds in theory and in practice for exact NN queries [Weber et al. '98]

# Back to the $\varepsilon$ -NN problem



pre-processing input:  $P, \varepsilon$

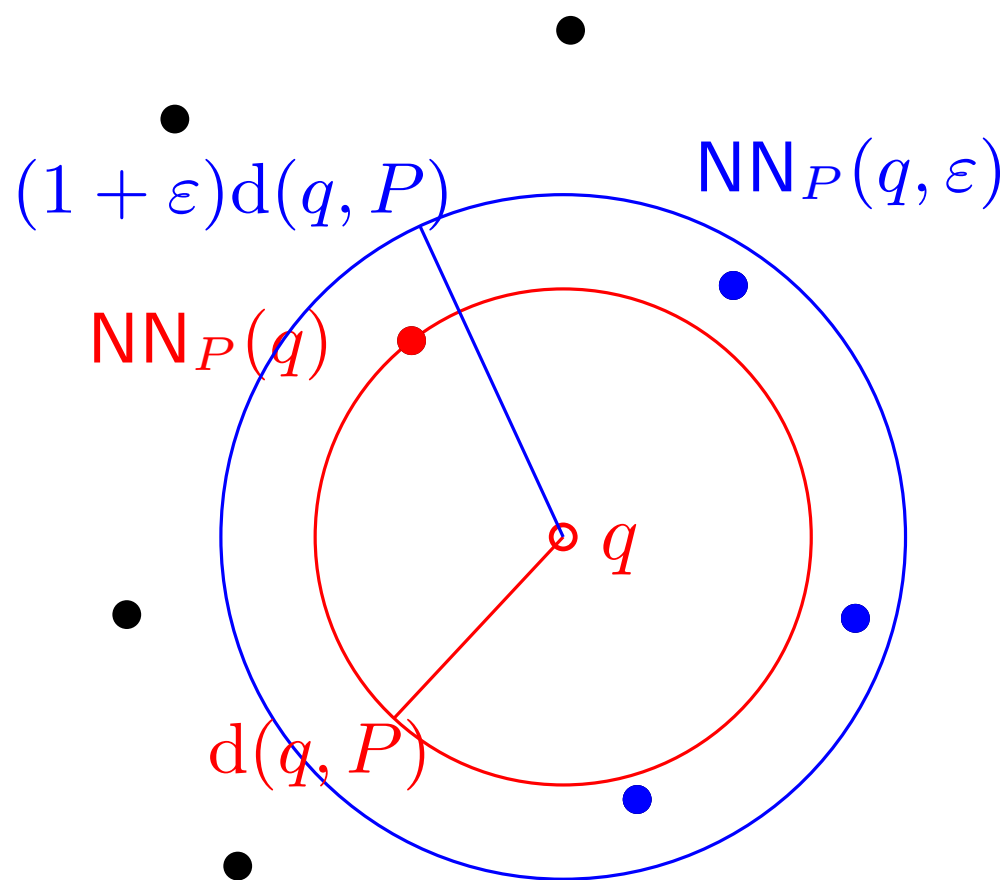
query input:  $q$

goal: find  $p \in \text{NN}_P(q, \varepsilon)$

**Curse of Dimensionality:** every DS for NN-search has either exponential size or exponential query time (in  $d$ ) in the worst case.

- holds in theory and in practice for exact NN queries [Weber et al. '98]
- still holds for approximate queries in decision tree model [Arya et al. '98]

# Back to the $\varepsilon$ -NN problem



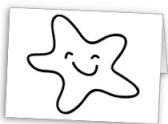
pre-processing input:  $P, \varepsilon$

query input:  $q$

goal: find  $p \in \text{NN}_P(q, \varepsilon)$

**Curse of Dimensionality:** every DS for NN-search has either exponential size or exponential query time (in  $d$ ) in the worst case.

- holds in theory and in practice for exact NN queries [Weber et al. '98]
- still holds for approximate queries in decision tree model [Arya et al. '98]
- no longer true in Real-RAM model thanks to LSH [Indyk, Motwani '98]



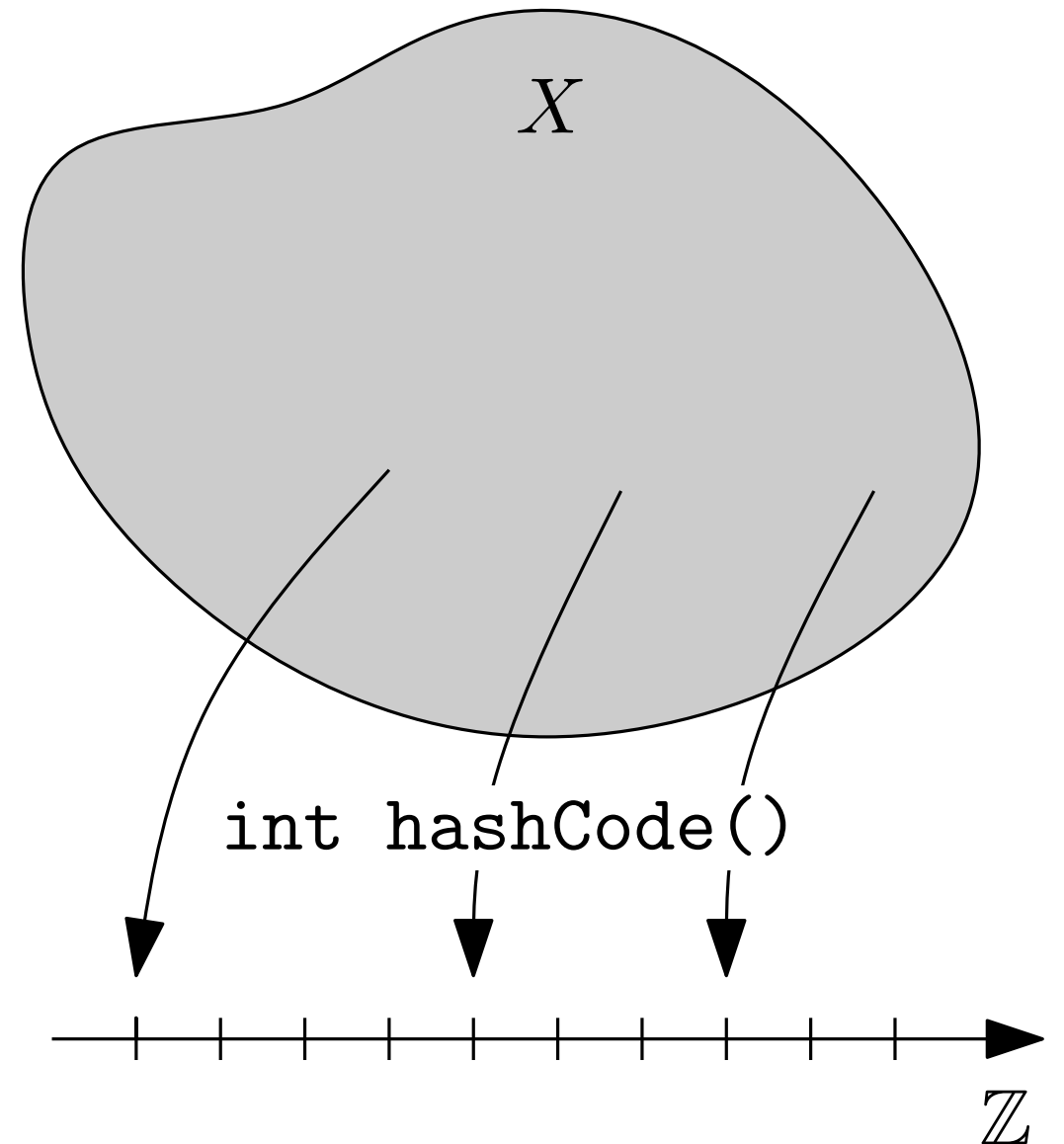
# Locality-Sensitive Hashing

Comparing elements via hashing:

$\text{hashCode} : X \rightarrow \mathbb{Z}$

$x = y \Rightarrow \text{hashCode}(x) = \text{hashCode}(y)$

$x \neq y \Rightarrow \text{hashCode}(x) \neq \text{hashCode}(y)$   
(no collisions hypothesis)



# Locality-Sensitive Hashing

Comparing elements via hashing:

$$\text{hashCode} : X \rightarrow \mathbb{Z}$$

$$x = y \Rightarrow \text{hashCode}(x) = \text{hashCode}(y)$$

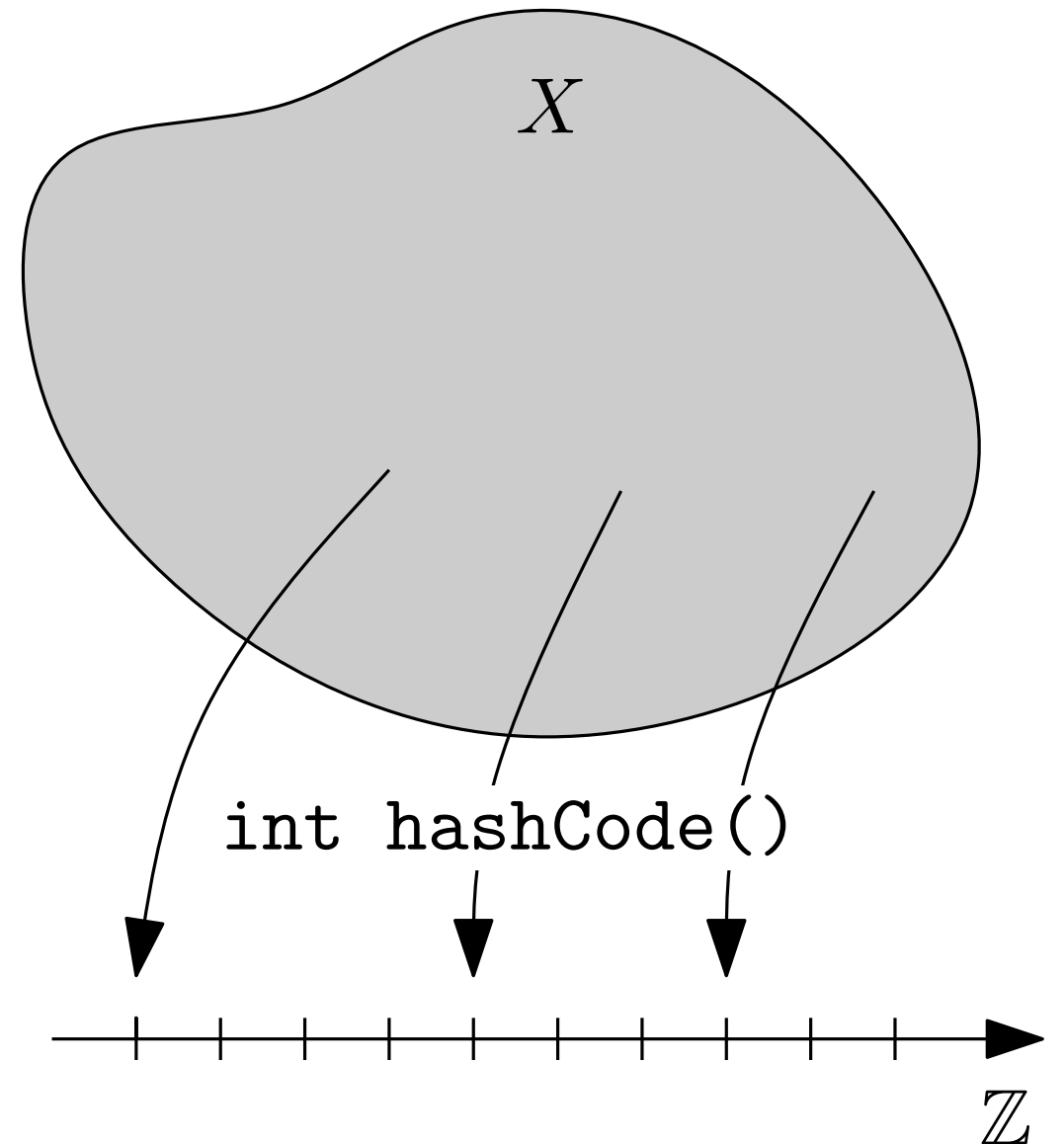
$$x \neq y \Rightarrow \text{hashCode}(x) \neq \text{hashCode}(y)$$

(no collisions hypothesis)

Metric case  $(X, d)$ : given  $r > 0$ ,

$$d(x, y) \leq r \Rightarrow \text{hashCode}(x) = \text{hashCode}(y)$$

$$d(x, y) > r \Rightarrow \text{hashCode}(x) \neq \text{hashCode}(y)$$



# Locality-Sensitive Hashing

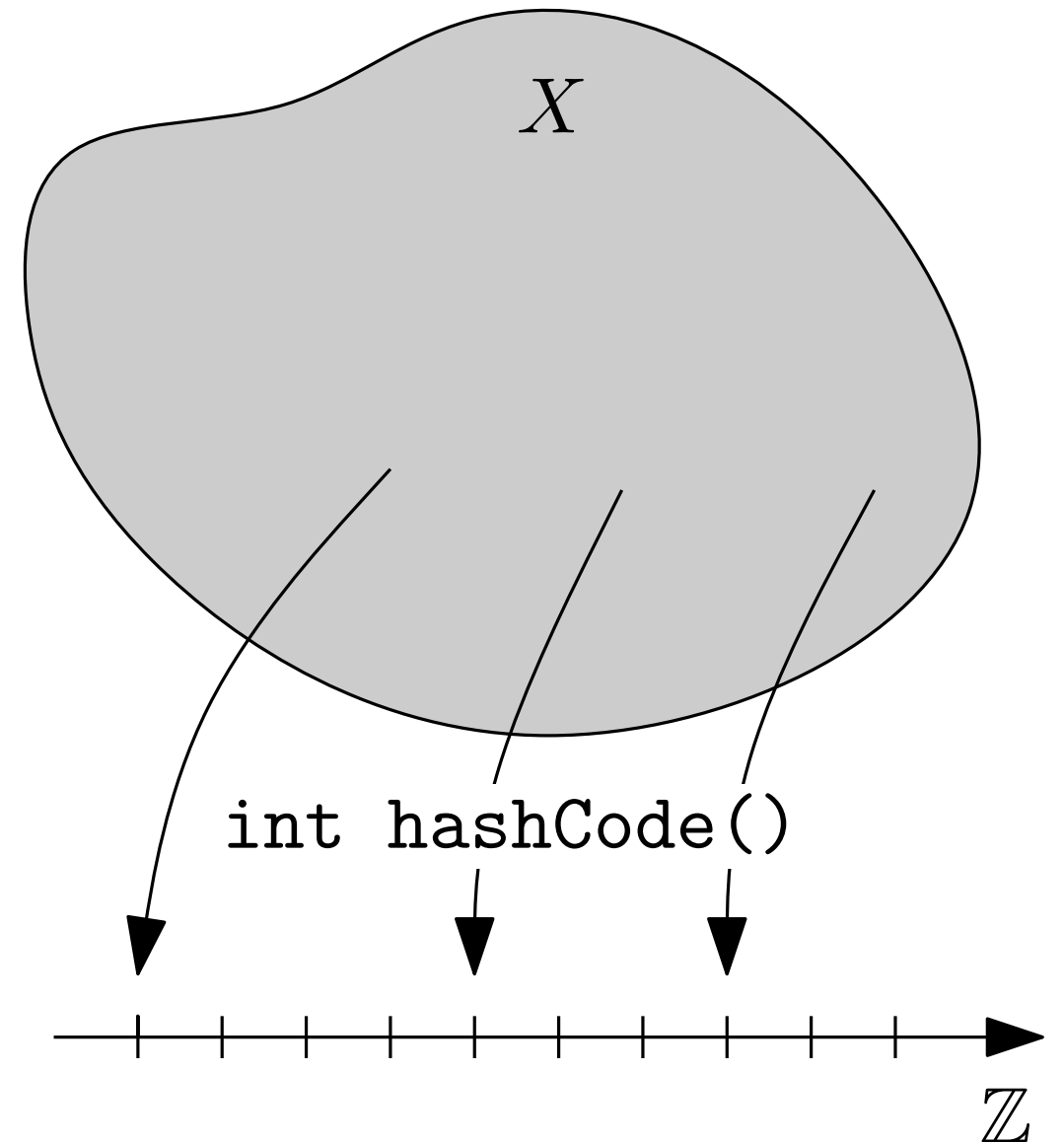
Comparing elements via hashing:

$$\text{hashCode} : X \rightarrow \mathbb{Z}$$

$$x = y \Rightarrow \text{hashCode}(x) = \text{hashCode}(y)$$

$$x \neq y \Rightarrow \text{hashCode}(x) \neq \text{hashCode}(y)$$

(no collisions hypothesis)



Metric case  $(X, d)$ : given  $r > 0$ ,

$$d(x, y) \leq r \Rightarrow \text{hashCode}(x) = \text{hashCode}(y)$$

$$d(x, y) > r \Rightarrow \text{hashCode}(x) \neq \text{hashCode}(y)$$

too good to be true  $\rightarrow$  allow for some slack

# Locality-Sensitive Hashing

**Def:** Given  $r_1 < r_2$ ,  $p_1 > p_2$  and  $\mathcal{U} \subset \mathbb{N}$ , a family  $\mathcal{F}$  of hash functions  $f : (X, d) \rightarrow \mathcal{U}$  is  $(r_1, r_2, p_1, p_2)$ -**sensitive** if  $\forall x, y \in X$ ,

- $d(x, y) \leq r_1 \Rightarrow \mathbb{P}[f(x) = f(y)] \geq p_1$
- $d(x, y) \geq r_2 \Rightarrow \mathbb{P}[f(x) = f(y)] \leq p_2$

(probability is over a random choice of function according to a given probability distribution over  $\mathcal{F}$ )

# Locality-Sensitive Hashing

**Def:** Given  $r_1 < r_2$ ,  $p_1 > p_2$  and  $\mathcal{U} \subset \mathbb{N}$ , a family  $\mathcal{F}$  of hash functions  $f : (X, d) \rightarrow \mathcal{U}$  is  $(r_1, r_2, p_1, p_2)$ -**sensitive** if  $\forall x, y \in X$ ,

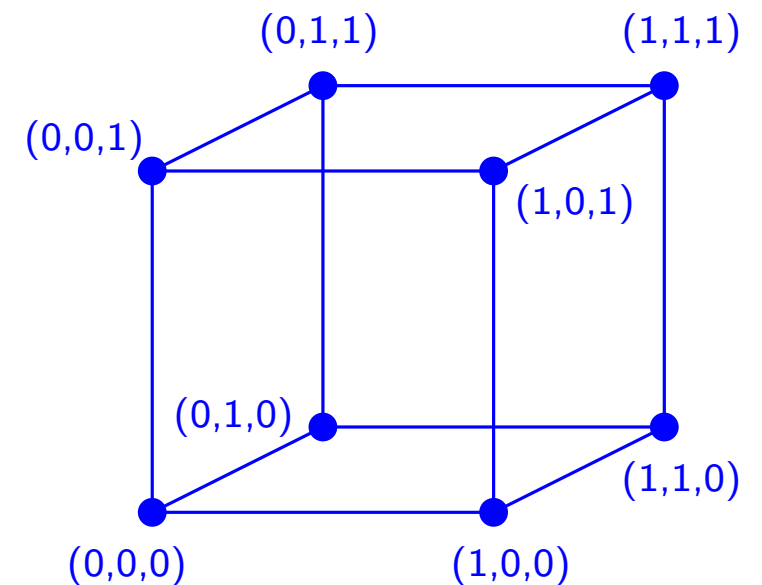
- $d(x, y) \leq r_1 \Rightarrow \mathbb{P}[f(x) = f(y)] \geq p_1$
- $d(x, y) \geq r_2 \Rightarrow \mathbb{P}[f(x) = f(y)] \leq p_2$

(probability is over a random choice of function according to a given probability distribution over  $\mathcal{F}$ )

**Example 1:**  $(X, d) = (\{0, 1\}^d, d_{\mathcal{H}})$

→ take  $\mathcal{F} = \{f_i\}_{i=1}^d$  where  $f_i(b_1 \cdots b_d) = b_i$  | unif. proba. on  $\mathcal{F}$

→  $\mathcal{F}$  is  $(r, r(1 + \varepsilon), 1 - \frac{r}{d}, 1 - \frac{r(1+\varepsilon)}{d})$ -sensitive for all  $r \geq 1$  and  $\varepsilon \geq 0$ .





# Locality-Sensitive Hashing

**Def:** Given  $r_1 < r_2$ ,  $p_1 > p_2$  and  $\mathcal{U} \subset \mathbb{N}$ , a family  $\mathcal{F}$  of hash functions  $f : (X, d) \rightarrow \mathcal{U}$  is  $(r_1, r_2, p_1, p_2)$ -**sensitive** if  $\forall x, y \in X$ ,

- $d(x, y) \leq r_1 \Rightarrow \mathbb{P}[f(x) = f(y)] \geq p_1$
- $d(x, y) \geq r_2 \Rightarrow \mathbb{P}[f(x) = f(y)] \leq p_2$

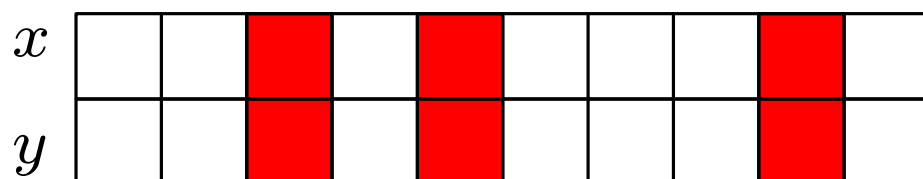
(probability is over a random choice of function according to a given probability distribution over  $\mathcal{F}$ )

**Example 1:**  $(X, d) = (\{0, 1\}^d, d_{\mathcal{H}})$

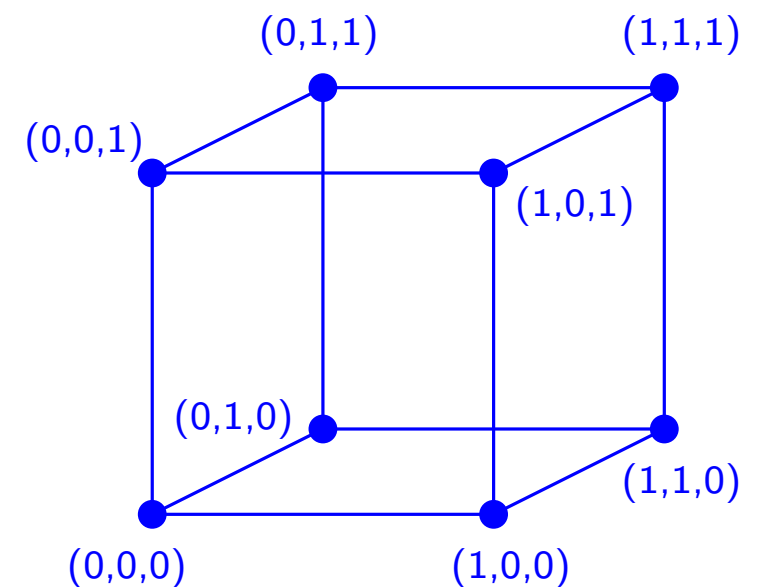
→ take  $\mathcal{F} = \{f_i\}_{i=1}^d$  where  $f_i(b_1 \cdots b_d) = b_i$  | unif. proba. on  $\mathcal{F}$

→  $\mathcal{F}$  is  $(r, r(1 + \varepsilon), 1 - \frac{r}{d}, 1 - \frac{r(1+\varepsilon)}{d})$ -sensitive for all  $r \geq 1$  and  $\varepsilon \geq 0$ .

proof:  $\forall x, y, \mathbb{P}_f[f(x) = f(y)] = \frac{d - d_{\mathcal{H}}(x, y)}{d} = 1 - \frac{d_{\mathcal{H}}(x, y)}{d}$



$d_{\mathcal{H}}(x, y)$  bits differ  $\Rightarrow d - d_{\mathcal{H}}(x, y)$  functions make  $x$  and  $y$  collide



# Locality-Sensitive Hashing

**Def:** Given  $r_1 < r_2$ ,  $p_1 > p_2$  and  $\mathcal{U} \subset \mathbb{N}$ , a family  $\mathcal{F}$  of hash functions  $f : (X, d) \rightarrow \mathcal{U}$  is  $(r_1, r_2, p_1, p_2)$ -**sensitive** if  $\forall x, y \in X$ ,

- $d(x, y) \leq r_1 \Rightarrow \mathbb{P}[f(x) = f(y)] \geq p_1$
- $d(x, y) \geq r_2 \Rightarrow \mathbb{P}[f(x) = f(y)] \leq p_2$

(probability is over a random choice of function according to a given probability distribution over  $\mathcal{F}$ )

**Example 2:**  $(X, d) = (\mathbb{R}^d, \|\cdot\|_2)$

→ take  $\mathcal{F} = \{f_{\mathbf{v}, b}\}_{\mathbf{v} \in \mathbb{R}^d}^{b \in [0, r]}$  where  $f_{\mathbf{v}, b}(x) = \lfloor \frac{x \cdot \mathbf{v} + b}{r} \rfloor$

→ choose  $\mathbf{v} = (v_1, \dots, v_d)$  with  $v_i \sim \mathcal{N}(0, 1)$ , and  $b$  uniformly in  $[0, r]$

→  $\mathcal{F}$  is  $(r, r(1 + \varepsilon), p_1, p_2)$  sensitive for  $p_1 = g(1)$  and  $p_2 = g(1 + \varepsilon)$ ,

where  $g(\kappa) = 1 - 2\text{cdf}(-r/\kappa) - \frac{2}{\sqrt{2\pi}r/\kappa}(1 - e^{-r^2/2\kappa^2})$

 cumulative density func. of normal distrib.

# Locality-Sensitive Hashing

**Def:** Given  $r_1 < r_2$ ,  $p_1 > p_2$  and  $\mathcal{U} \subset \mathbb{N}$ , a family  $\mathcal{F}$  of hash functions  $f : (X, d) \rightarrow \mathcal{U}$  is  $(r_1, r_2, p_1, p_2)$ -**sensitive** if  $\forall x, y \in X$ ,

- $d(x, y) \leq r_1 \Rightarrow \mathbb{P}[f(x) = f(y)] \geq p_1$
- $d(x, y) \geq r_2 \Rightarrow \mathbb{P}[f(x) = f(y)] \leq p_2$

(probability is over a random choice of function according to a given probability distribution over  $\mathcal{F}$ )

**Lemma** [Johnson, Lindenstrauss 84]:

*For any dimensions  $0 < k < d$  there is a probability distribution  $\mu$  over the projections  $\mathbb{R}^d \rightarrow \mathbb{R}^k$  such that, given any set  $P$  of  $n$  points in  $\mathbb{R}^d$  and any  $\varepsilon \in (0, 1)$  with  $k > 10 \ln n / \varepsilon^2$ , a projection  $\pi : \mathbb{R}^d \rightarrow \mathbb{R}^k$  sampled at random from  $\mu$  satisfies w.h.p.*

$$\forall p, q \in P, (1 - \varepsilon)\|p - q\| \leq \|\pi(p) - \pi(q)\| \leq (1 + \varepsilon)\|p - q\|$$

# Locality-Sensitive Hashing

**Def:** Given  $r_1 < r_2$ ,  $p_1 > p_2$  and  $\mathcal{U} \subset \mathbb{N}$ , a family  $\mathcal{F}$  of hash functions  $f : (X, d) \rightarrow \mathcal{U}$  is  $(r_1, r_2, p_1, p_2)$ -**sensitive** if  $\forall x, y \in X$ ,

- $d(x, y) \leq r_1 \Rightarrow \mathbb{P}[f(x) = f(y)] \geq p_1$
- $d(x, y) \geq r_2 \Rightarrow \mathbb{P}[f(x) = f(y)] \leq p_2$

(probability is over a random choice of function according to a given probability distribution over  $\mathcal{F}$ )

**Lemma** [Johnson, Lindenstrauss 84]:

*For any dimensions  $0 < k < d$  there is a probability distribution  $\mu$  over the projections  $\mathbb{R}^d \rightarrow \mathbb{R}^k$  such that, given any set  $P$  of  $n$  points in  $\mathbb{R}^d$  and any  $\varepsilon \in (0, 1)$  with  $k > 10 \ln n / \varepsilon^2$ , a projection  $\pi : \mathbb{R}^d \rightarrow \mathbb{R}^k$  sampled at random from  $\mu$  satisfies w.h.p.*

$$\forall p, q \in P, (1 - \varepsilon)\|p - q\| \leq \|\pi(p) - \pi(q)\| \leq (1 + \varepsilon)\|p - q\|$$

Proof idea: take  $\pi(x) = \sum_{i=1}^k (x \cdot \mathbf{v}_i) \mathbf{v}_i$ , where the  $\mathbf{v}_i$  are random vectors in  $\mathbb{R}^d$  whose coordinates are chosen i.i.d. according to  $\mathcal{N}(0, 1)$ . Then exploit the concentration of measure on the sphere to show that  $\pi$  induces a low distortion.

# Locality-Sensitive Hashing

**Def:** Given  $r_1 < r_2$ ,  $p_1 > p_2$  and  $\mathcal{U} \subset \mathbb{N}$ , a family  $\mathcal{F}$  of hash functions  $f : (X, d) \rightarrow \mathcal{U}$  is  $(r_1, r_2, p_1, p_2)$ -**sensitive** if  $\forall x, y \in X$ ,

- $d(x, y) \leq r_1 \Rightarrow \mathbb{P}[f(x) = f(y)] \geq p_1$
- $d(x, y) \geq r_2 \Rightarrow \mathbb{P}[f(x) = f(y)] \leq p_2$

(probability is over a random choice of function according to a given probability distribution over  $\mathcal{F}$ )

→ **General idea:**

- choose  $k$ -dimensional vector of random functions  $(f_1, \dots, f_k) \in \mathcal{F}^k$
- pre-process  $P$  by hashing its points into the corresponding hash table
- given  $q \in X$ , hash  $q$  and choose collision with smallest distance

# Locality-Sensitive Hashing

**Def:** Given  $r_1 < r_2$ ,  $p_1 > p_2$  and  $\mathcal{U} \subset \mathbb{N}$ , a family  $\mathcal{F}$  of hash functions  $f : (X, d) \rightarrow \mathcal{U}$  is  $(r_1, r_2, p_1, p_2)$ -**sensitive** if  $\forall x, y \in X$ ,

- $d(x, y) \leq r_1 \Rightarrow \mathbb{P}[f(x) = f(y)] \geq p_1$
- $d(x, y) \geq r_2 \Rightarrow \mathbb{P}[f(x) = f(y)] \leq p_2$

(probability is over a random choice of function according to a given probability distribution over  $\mathcal{F}$ )

→ **General idea:**

- choose  $k$ -dimensional vector of random functions  $(f_1, \dots, f_k) \in \mathcal{F}^k$
- pre-process  $P$  by hashing its points into the corresponding hash table
- given  $q \in X$ , hash  $q$  and choose collision with smallest distance

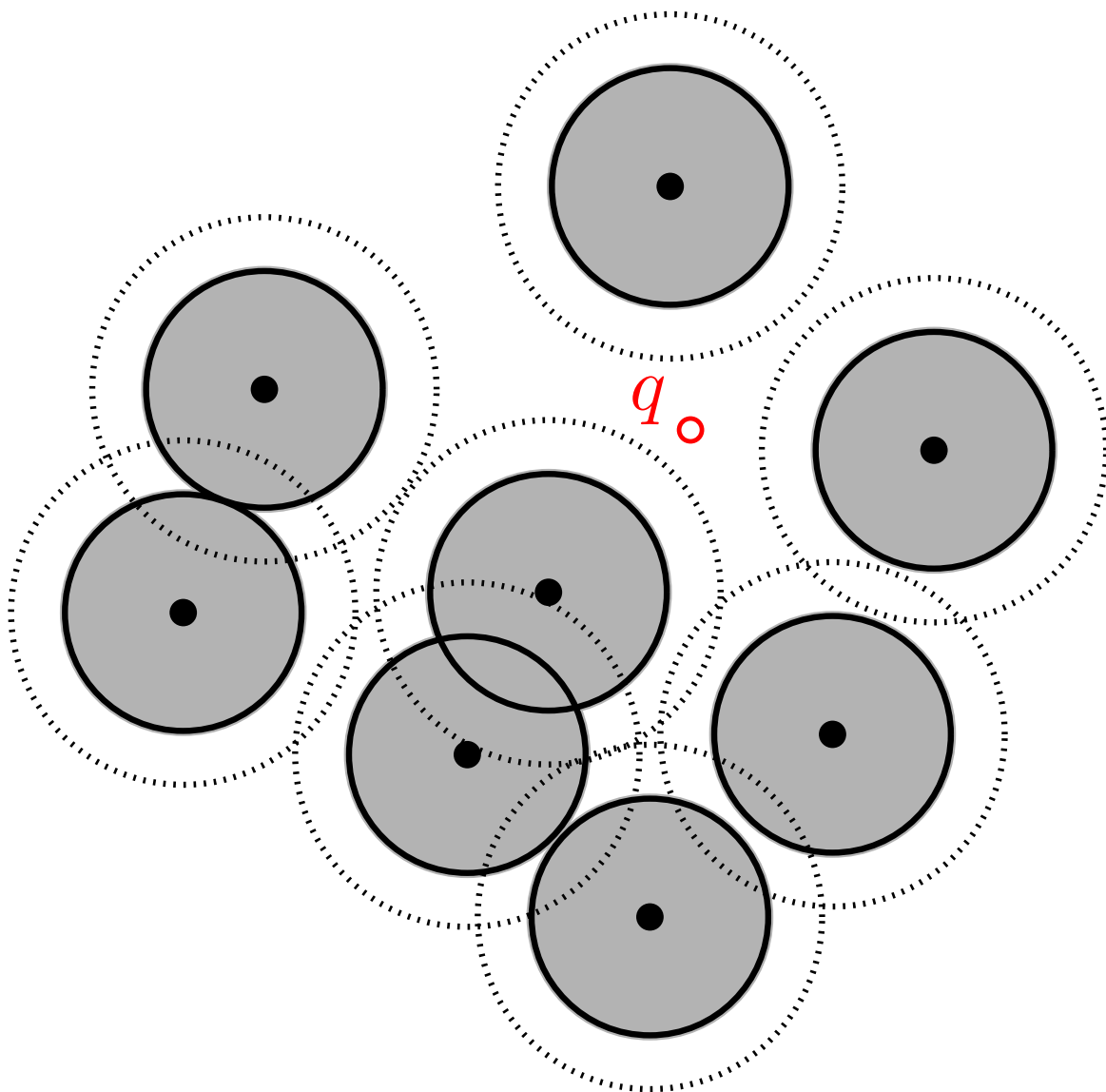
→ *Technical detail:* family works only for fixed  $r_1, r_2$

→ fix  $r_1 = r$  and  $r_2 = r(1 + \varepsilon)$ , and solve  $(r, \varepsilon)$ -NN query

# The $(r, \varepsilon)$ -NN problem (PLEB)

**Goal:** pre-process  $P$  such that, for any query point  $q$ ,

- if  $d(q, P) \leq r$  then answer YES and return some  $p \in \text{NN}_P(q, r, \varepsilon)$ ,
- if  $d(q, P) > r(1 + \varepsilon)$  then answer NO,
- else  $(r < d(q, P) \leq r(1 + \varepsilon))$  give any of the above answers.



# The $(r, \varepsilon)$ -NN problem (PLEB)

**Goal:** pre-process  $P$  such that, for any query point  $q$ ,

- if  $d(q, P) \leq r$  then answer YES and return some  $p \in \text{NN}_P(q, r, \varepsilon)$ ,
- if  $d(q, P) > r(1 + \varepsilon)$  then answer NO,
- else  $(r < d(q, P) \leq r(1 + \varepsilon))$  give any of the above answers.

**Step 1:** boost the sensitivity of the hash family

$$\mathcal{G} = \{g = (f_1, \dots, f_k) \in \mathcal{F}^k \mid f_1, \dots, f_k \text{ chosen randomly in } \mathcal{F}\}$$



# The $(r, \varepsilon)$ -NN problem (PLEB)

**Goal:** pre-process  $P$  such that, for any query point  $q$ ,

- if  $d(q, P) \leq r$  then answer YES and return some  $p \in \text{NN}_P(q, r, \varepsilon)$ ,
- if  $d(q, P) > r(1 + \varepsilon)$  then answer NO,
- else  $(r < d(q, P) \leq r(1 + \varepsilon))$  give any of the above answers.

**Step 1:** boost the sensitivity of the hash family

$$\mathcal{G} = \{g = (f_1, \dots, f_k) \in \mathcal{F}^k \mid f_1, \dots, f_k \text{ chosen randomly in } \mathcal{F}\}$$

$$\rightarrow \forall x, y, d(x, y) \leq r \Rightarrow \mathbb{P}[g(x) = g(y)] \geq p_1^k \text{ (coords. are independent)}$$

$$d(x, y) > r(1 + \varepsilon) \Rightarrow \mathbb{P}[g(x) = g(y)] \leq p_2^k$$

# The $(r, \varepsilon)$ -NN problem (PLEB)

**Goal:** pre-process  $P$  such that, for any query point  $q$ ,

- if  $d(q, P) \leq r$  then answer YES and return some  $p \in \text{NN}_P(q, r, \varepsilon)$ ,
- if  $d(q, P) > r(1 + \varepsilon)$  then answer NO,
- else ( $r < d(q, P) \leq r(1 + \varepsilon)$ ) give any of the above answers.

**Step 2:** pre-process the data points

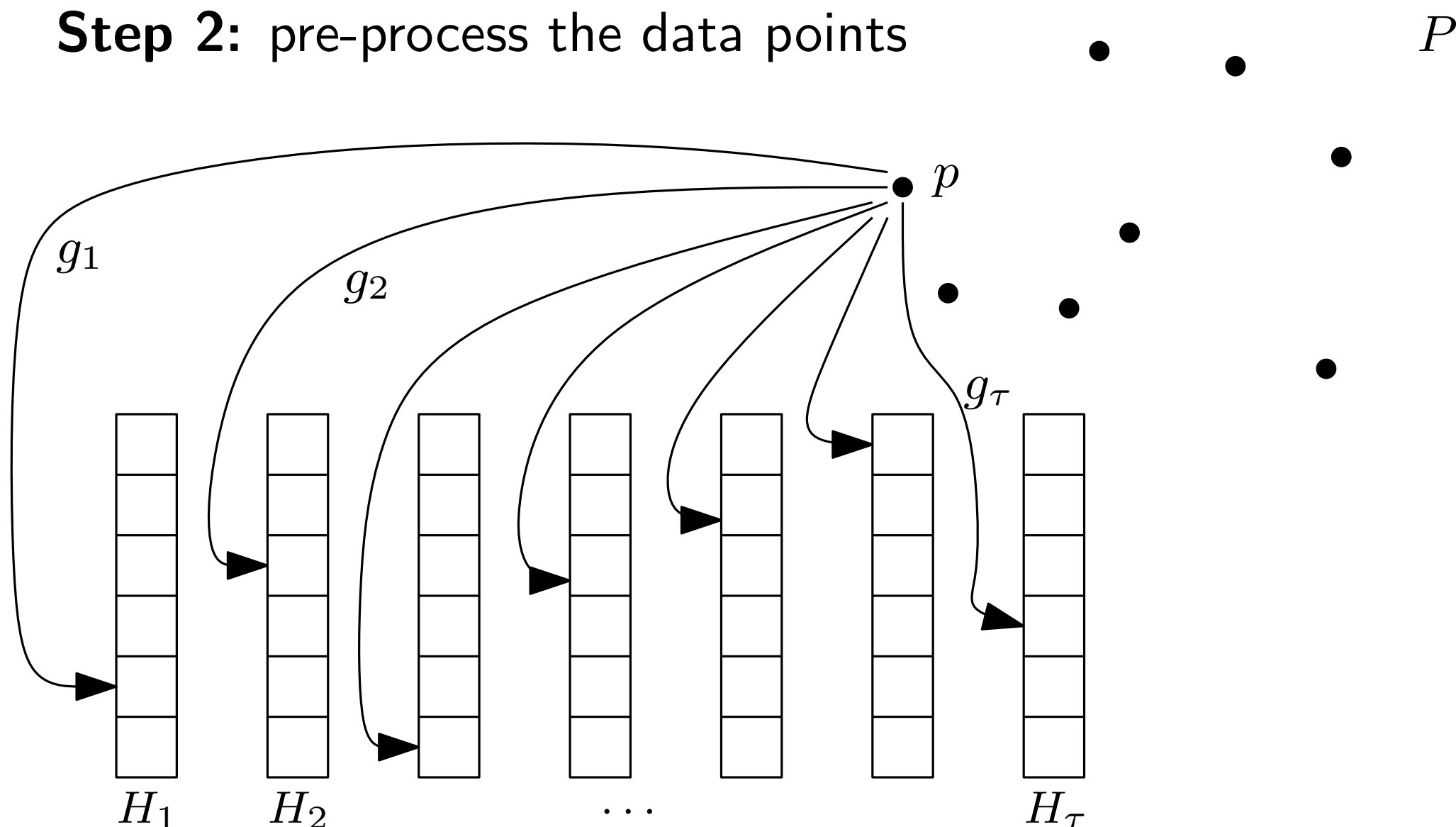
- choose  $\tau$  random functions  $g_1, \dots, g_\tau$  from boosted hash family  $\mathcal{G}$ ,
- initialise  $\tau$  hash tables  $H_1, \dots, H_\tau$
- $\forall i = 1, \dots, \tau$ , hash every point  $p \in P$  into  $H_i$  using  $g_i(p)$  as the key
- keep only one arbitrary point per non-empty entry

# The $(r, \varepsilon)$ -NN problem (PLEB)

**Goal:** pre-process  $P$  such that, for any query point  $q$ ,

- if  $d(q, P) \leq r$  then answer YES and return some  $p \in \text{NN}_P(q, r, \varepsilon)$ ,
- if  $d(q, P) > r(1 + \varepsilon)$  then answer NO,
- else  $(r < d(q, P) \leq r(1 + \varepsilon))$  give any of the above answers.

**Step 2:** pre-process the data points

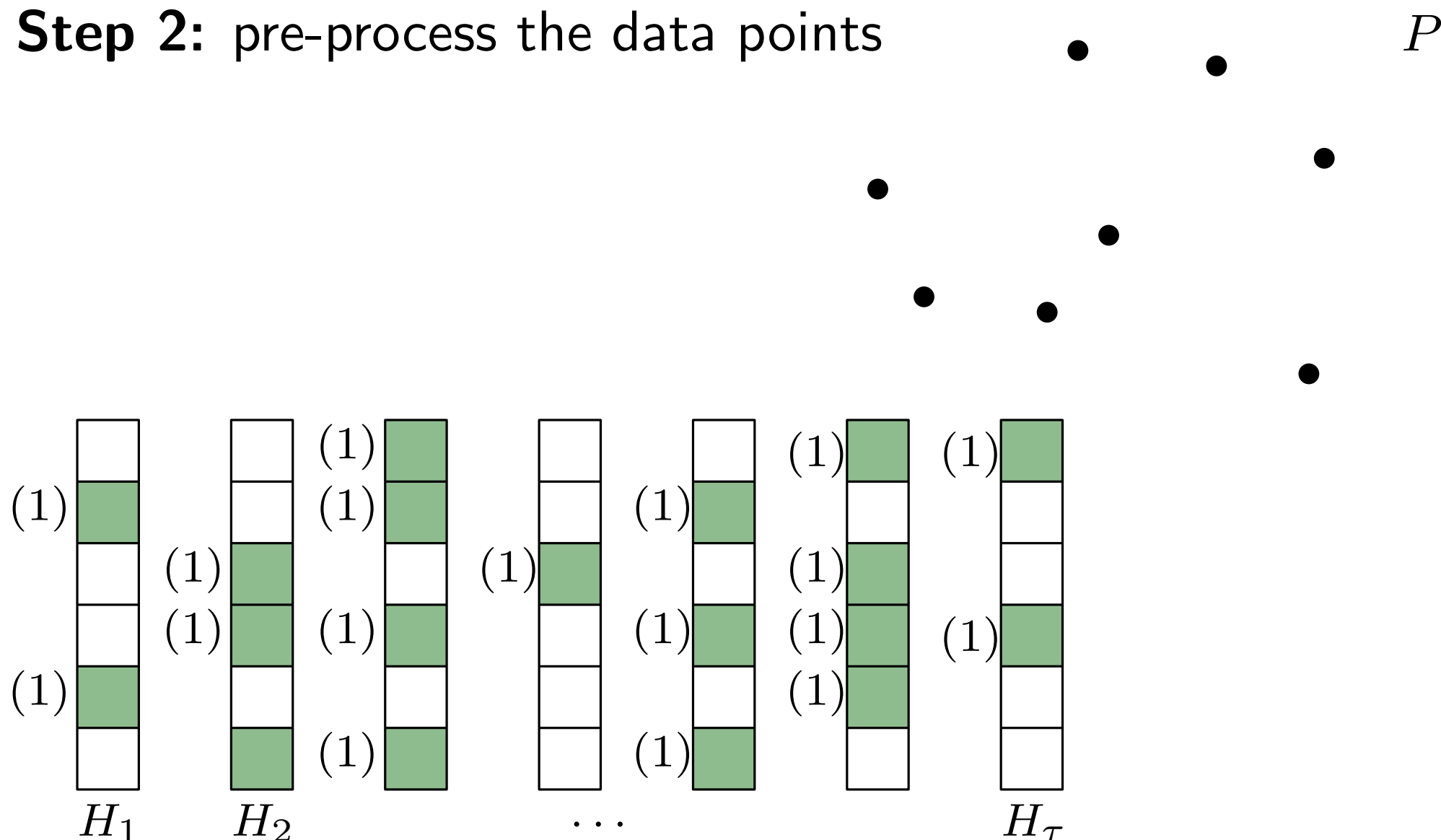


# The $(r, \varepsilon)$ -NN problem (PLEB)

**Goal:** pre-process  $P$  such that, for any query point  $q$ ,

- if  $d(q, P) \leq r$  then answer YES and return some  $p \in \text{NN}_P(q, r, \varepsilon)$ ,
- if  $d(q, P) > r(1 + \varepsilon)$  then answer NO,
- else  $(r < d(q, P) \leq r(1 + \varepsilon))$  give any of the above answers.

**Step 2:** pre-process the data points

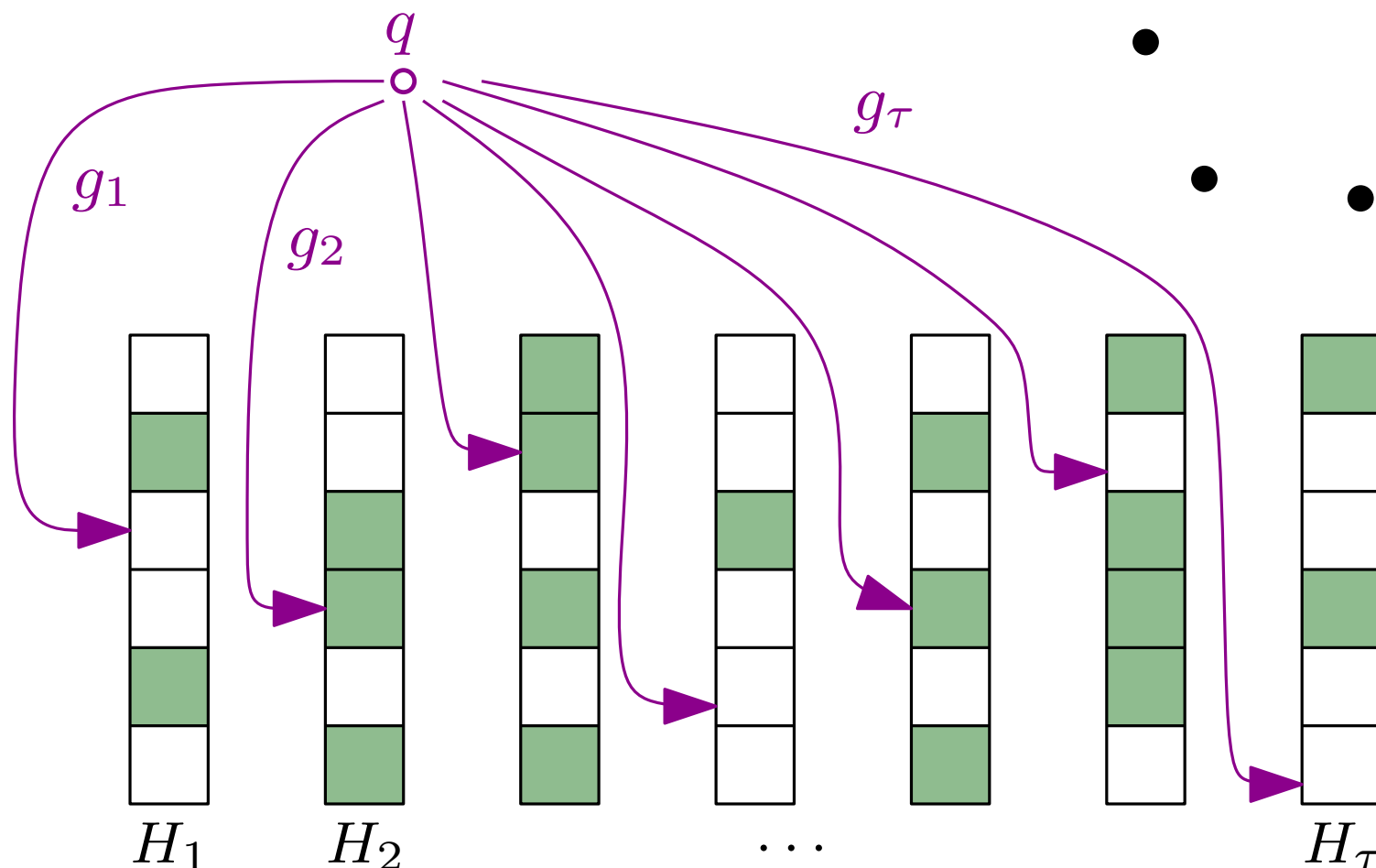


# The $(r, \varepsilon)$ -NN problem (PLEB)

**Goal:** pre-process  $P$  such that, for any query point  $q$ ,

- if  $d(q, P) \leq r$  then answer YES and return some  $p \in \text{NN}_P(q, r, \varepsilon)$ ,
- if  $d(q, P) > r(1 + \varepsilon)$  then answer NO,
- else  $(r < d(q, P) \leq r(1 + \varepsilon))$  give any of the above answers.

**Step 3:** hash the query point using the  $g_i$  • • •  $P$



# The $(r, \varepsilon)$ -NN problem (PLEB)

**Goal:** pre-process  $P$  such that, for any query point  $q$ ,

- if  $d(q, P) \leq r$  then answer YES and return some  $p \in \text{NN}_P(q, r, \varepsilon)$ ,
- if  $d(q, P) > r(1 + \varepsilon)$  then answer NO,
- else  $(r < d(q, P) \leq r(1 + \varepsilon))$  give any of the above answers.

**Step 3:** hash the query point using the  $g_i$

- let  $p_1, \dots, p_l$  be the collisions ( $l \leq \tau$ )
- if some  $p_j$  is such that  $d(p_j, q) \leq r(1 + \varepsilon)$  then return YES and  $p_j$
- else return NO

# The $(r, \varepsilon)$ -NN problem (PLEB)

**Goal:** pre-process  $P$  such that, for any query point  $q$ ,

- if  $d(q, P) \leq r$  then answer YES and return some  $p \in \text{NN}_P(q, r, \varepsilon)$ ,
- if  $d(q, P) > r(1 + \varepsilon)$  then answer NO,
- else  $(r < d(q, P) \leq r(1 + \varepsilon))$  give any of the above answers.

**Analysis in a nutshell:**

- test  $\Rightarrow$  return NO whenever  $d(q, P) > r(1 + \varepsilon)$

# The $(r, \varepsilon)$ -NN problem (PLEB)

**Goal:** pre-process  $P$  such that, for any query point  $q$ ,

- if  $d(q, P) \leq r$  then answer YES and return some  $p \in \text{NN}_P(q, r, \varepsilon)$ ,
- if  $d(q, P) > r(1 + \varepsilon)$  then answer NO,
- else ( $r < d(q, P) \leq r(1 + \varepsilon)$ ) give any of the above answers.

## Analysis in a nutshell:

- test  $\Rightarrow$  return NO whenever  $d(q, P) > r(1 + \varepsilon)$
- if  $\exists p \in P$  s.t.  $d(p, q) \leq r$ , then for a fixed  $i \in \{1, \dots, \tau\}$ ,

$$\mathbb{P}[p \text{ collides with } q \text{ in } H_i] \geq p_1^k$$

$$\forall p' \in P \setminus B(q, r(1 + \varepsilon)), \mathbb{P}[p' \text{ collides with } q \text{ in } H_i \mid p \text{ collides with } q \text{ in } H_i]$$

$$= \frac{\mathbb{P}[p \text{ and } p' \text{ collide with } q \text{ in } H_i]}{\mathbb{P}[p \text{ collides with } q \text{ in } H_i]}$$

$$\leq \frac{\mathbb{P}[p' \text{ collides with } q \text{ in } H_i]}{\mathbb{P}[p \text{ collides with } q \text{ in } H_i]} \leq \left(\frac{p_2}{p_1}\right)^k$$



# The $(r, \varepsilon)$ -NN problem (PLEB)

**Goal:** pre-process  $P$  such that, for any query point  $q$ ,

- if  $d(q, P) \leq r$  then answer YES and return some  $p \in \text{NN}_P(q, r, \varepsilon)$ ,
- if  $d(q, P) > r(1 + \varepsilon)$  then answer NO,
- else ( $r < d(q, P) \leq r(1 + \varepsilon)$ ) give any of the above answers.

## Analysis in a nutshell:

- test  $\Rightarrow$  return NO whenever  $d(q, P) > r(1 + \varepsilon)$
- if  $\exists p \in P$  s.t.  $d(p, q) \leq r$ , then for a fixed  $i \in \{1, \dots, \tau\}$ ,

$$\text{union bound} \Rightarrow \mathbb{P}[H_i \text{ succeeds}] \geq p_1^k \left( 1 - n \left( \frac{p_2}{p_1} \right)^k \right)$$

$\tau$  independent hash tables  $\Rightarrow$

$$\mathbb{P}[\text{some } H_i \text{ succeeds}] \geq 1 - \left( 1 - p_1^k \left( 1 - n \left( \frac{p_2}{p_1} \right)^k \right) \right)^\tau$$

# The $(r, \varepsilon)$ -NN problem (PLEB)

**Goal:** pre-process  $P$  such that, for any query point  $q$ ,

- if  $d(q, P) \leq r$  then answer YES and return some  $p \in \text{NN}_P(q, r, \varepsilon)$ ,
- if  $d(q, P) > r(1 + \varepsilon)$  then answer NO,
- else ( $r < d(q, P) \leq r(1 + \varepsilon)$ ) give any of the above answers.

## Analysis in a nutshell:

- test  $\Rightarrow$  return NO whenever  $d(q, P) > r(1 + \varepsilon)$
- if  $\exists p \in P$  s.t.  $d(p, q) \leq r$ , then for a fixed  $i \in \{1, \dots, \tau\}$ ,

$$\text{union bound} \Rightarrow \mathbb{P}[H_i \text{ succeeds}] \geq p_1^k \left( 1 - n \left( \frac{p_2}{p_1} \right)^k \right)$$

$\tau$  independent hash tables  $\Rightarrow$

$$\mathbb{P}[\text{some } H_i \text{ succeeds}] \geq 1 - \left( 1 - p_1^k \left( 1 - n \left( \frac{p_2}{p_1} \right)^k \right) \right)^\tau$$

Let  $k = c \log n$  and  $\tau = n^\varrho$  where  $\varrho = \frac{\ln p_1}{\ln p_2} \in (0, 1)$ .

$$\Rightarrow \boxed{\text{query time} = O(n^\varrho \log n), \Pr[\text{success}] \geq 1 - 1/n^{c\varrho}}$$

# From $(r, \varepsilon)$ -NN to $\varepsilon$ -NN

- Special case:  $(X, d) = (\{0, 1\}^d, d_{\mathcal{H}})$

**Observation:** inter-point distances lie within  $\{0, 1, 2, \dots, d\}$

# From $(r, \varepsilon)$ -NN to $\varepsilon$ -NN

- Special case:  $(X, d) = (\{0, 1\}^d, d_{\mathcal{H}})$

**Observation:** inter-point distances lie within  $\{0, 1, 2, \dots, d\}$

- solve case  $d_{\mathcal{H}}(q, P) = 0$  independently (use lexicographical sorting)
- take geometric sequence  $r_0 = 1, r_1 = 1 + \varepsilon, \dots, r_j = (1 + \varepsilon)^j, \dots$
- for  $j = 0$  to  $\lceil \log_{1+\varepsilon} d \rceil = O(\frac{1}{\varepsilon} \log d)$ , solve  $(r_j, r_{j+1})$ -NN query
- let  $j_l$  be the lowest  $j$  s.t. the answer to  $(r_j, r_{j+1})$ -MM query is YES
- return the output point of the  $(r_{j_l}, r_{j_l+1})$ -NN query
- if no YES answer, return output of case  $d_{\mathcal{H}}(q, P) = 0$

# From $(r, \varepsilon)$ -NN to $\varepsilon$ -NN

- Special case:  $(X, d) = (\{0, 1\}^d, d_{\mathcal{H}})$

**Observation:** inter-point distances lie within  $\{0, 1, 2, \dots, d\}$

→ query time =  $O(\frac{1}{\varepsilon} n^{\varrho} \log n \log d)$

(becomes  $O(d n^{\varrho} \log n)$  if arithmetic sequence is used)

# From $(r, \varepsilon)$ -NN to $\varepsilon$ -NN

- Special case:  $(X, d) = (\{0, 1\}^d, d_{\mathcal{H}})$

**Observation:** inter-point distances lie within  $\{0, 1, 2, \dots, d\}$

→ query time =  $O(\frac{1}{\varepsilon} n^{\varrho} \log n \log d)$

(becomes  $O(d n^{\varrho} \log n)$  if arithmetic sequence is used)

**Observation:** deterministically,  $r_{j_l} \geq d_{\mathcal{H}}(q, P)/(1 + \varepsilon)$

⇒ output  $\in NN_P(q, \varepsilon(2 + \varepsilon))$  iff LSH data structure works for  $j = j_l$

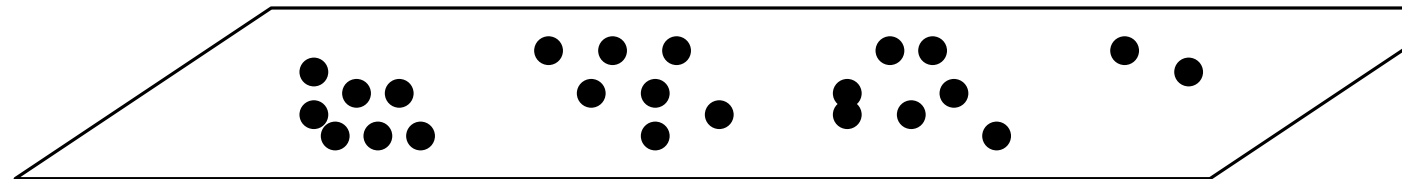
⇒  $\mathbb{P}[\text{success}] \geq 1 - 1/n^{c_{\varrho}}$

# From $(r, \varepsilon)$ -NN to $\varepsilon$ -NN

- General case: use **hierarchical clustering tree** [Har-Peled'01]
  - consider geometric sequences of scales as before
  - cluster data points in order to bound the lengths of the sequences

# From $(r, \varepsilon)$ -NN to $\varepsilon$ -NN

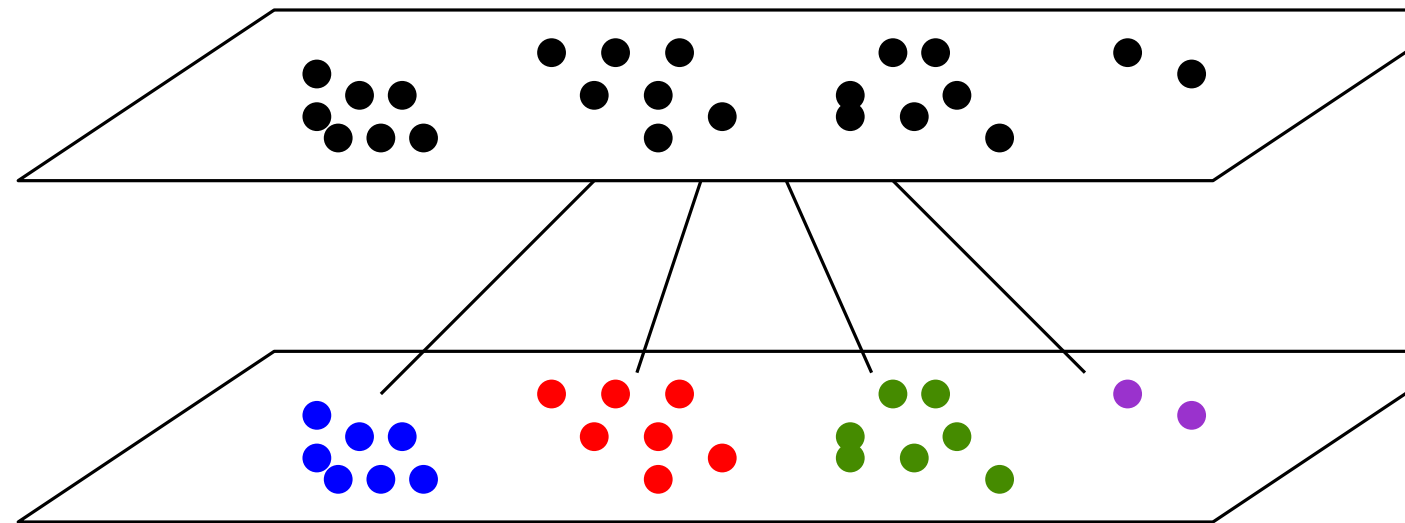
- General case: use **hierarchical clustering tree** [Har-Peled'01]
  - consider geometric sequences of scales as before
  - cluster data points in order to bound the lengths of the sequences





# From $(r, \varepsilon)$ -NN to $\varepsilon$ -NN

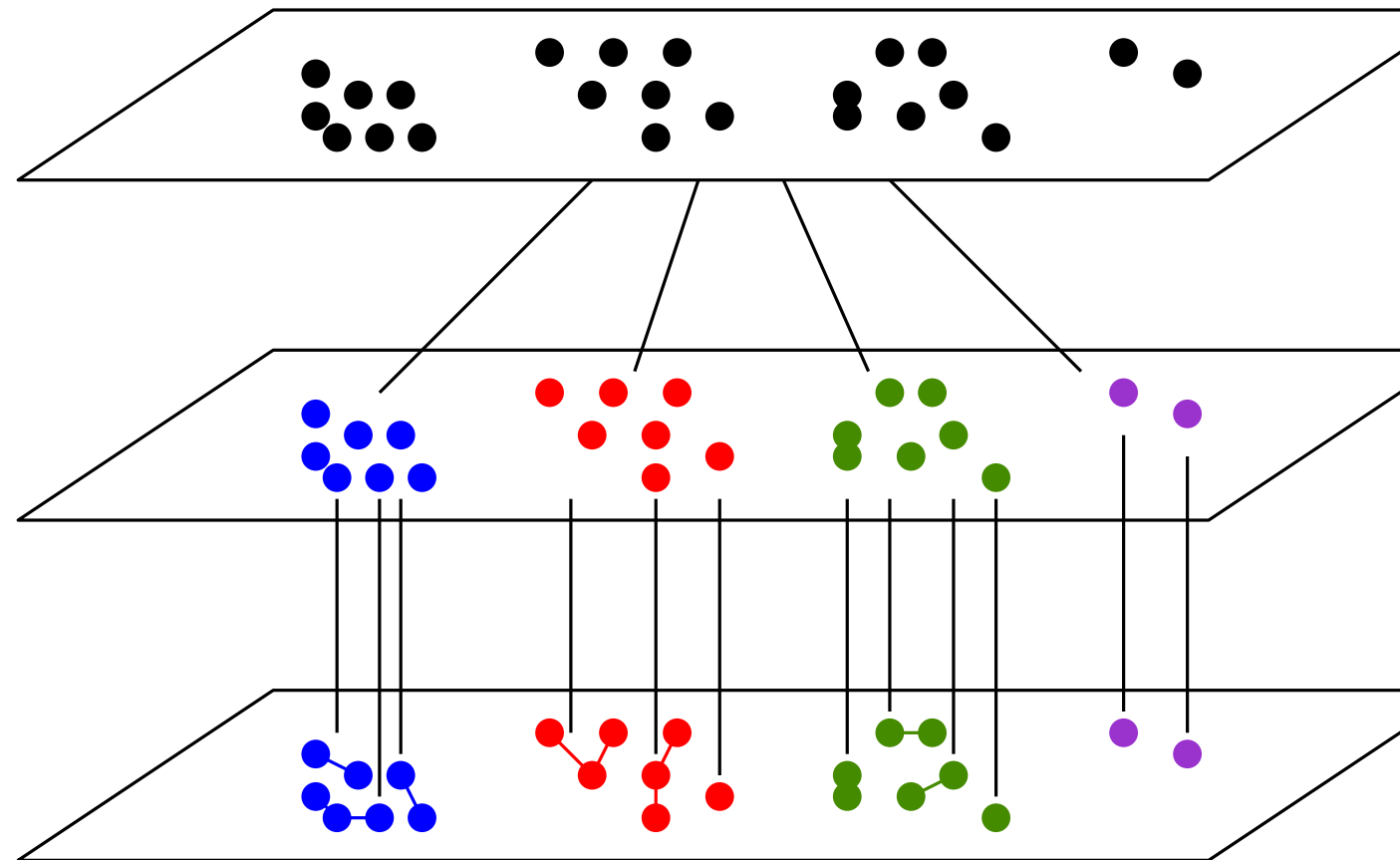
- General case: use **hierarchical clustering tree** [Har-Peled'01]
  - consider geometric sequences of scales as before
  - cluster data points in order to bound the lengths of the sequences



# From $(r, \varepsilon)$ -NN to $\varepsilon$ -NN

- General case: use **hierarchical clustering tree** [Har-Peled'01]
  - consider geometric sequences of scales as before
  - cluster data points in order to bound the lengths of the sequences

Assign  $P_v \subseteq P$  and  $[r_v, R_v]$  to each node  $v$



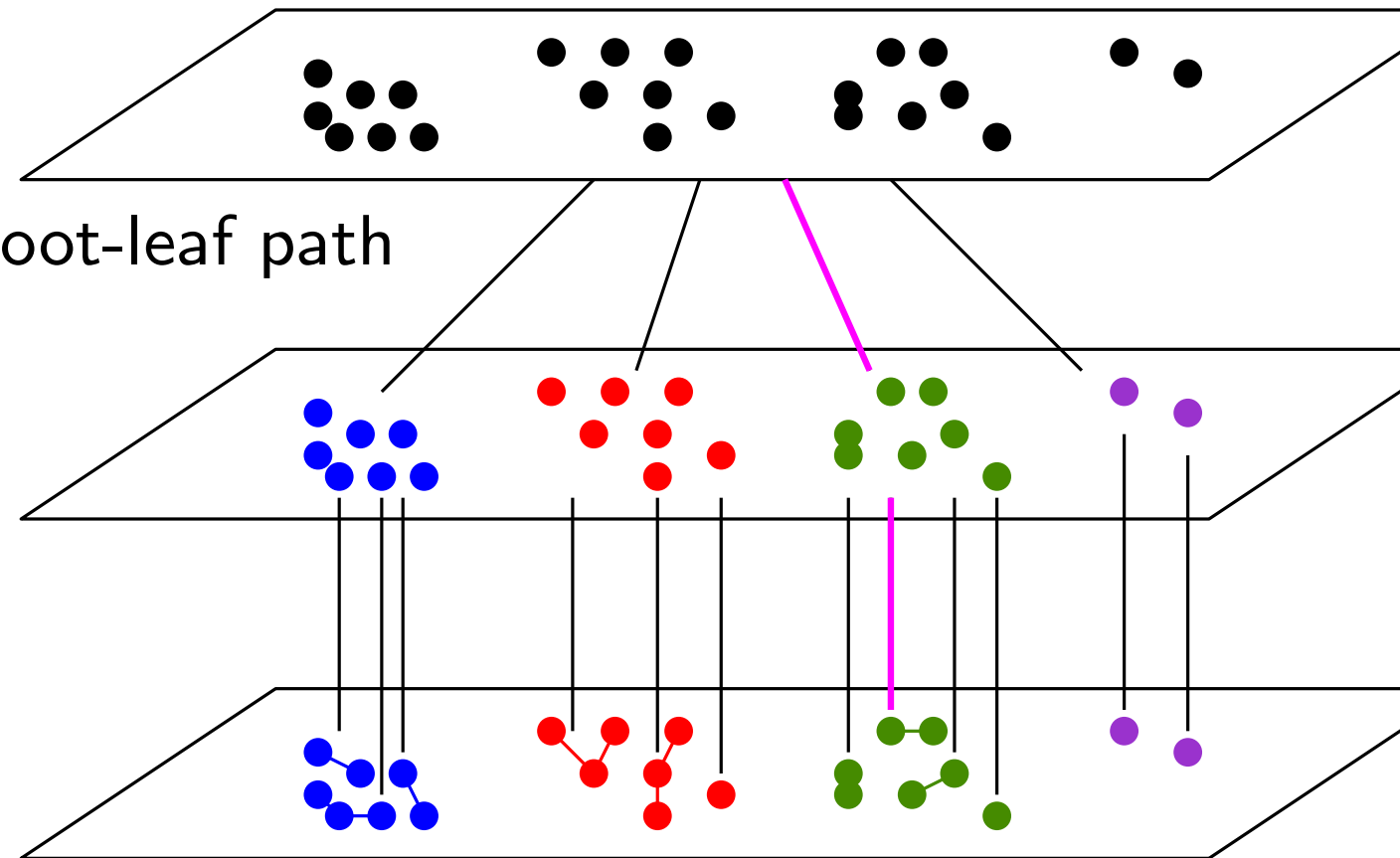
# From $(r, \varepsilon)$ -NN to $\varepsilon$ -NN

- General case: use **hierarchical clustering tree** [Har-Peled'01]
  - consider geometric sequences of scales as before
  - cluster data points in order to bound the lengths of the sequences

Assign  $P_v \subseteq P$  and  $[r_v, R_v]$  to each node  $v$

$\varepsilon$ -NN query:

- traverse down the tree along one root-leaf path



# From $(r, \varepsilon)$ -NN to $\varepsilon$ -NN

- General case: use **hierarchical clustering tree** [Har-Peled'01]
  - consider geometric sequences of scales as before
  - cluster data points in order to bound the lengths of the sequences

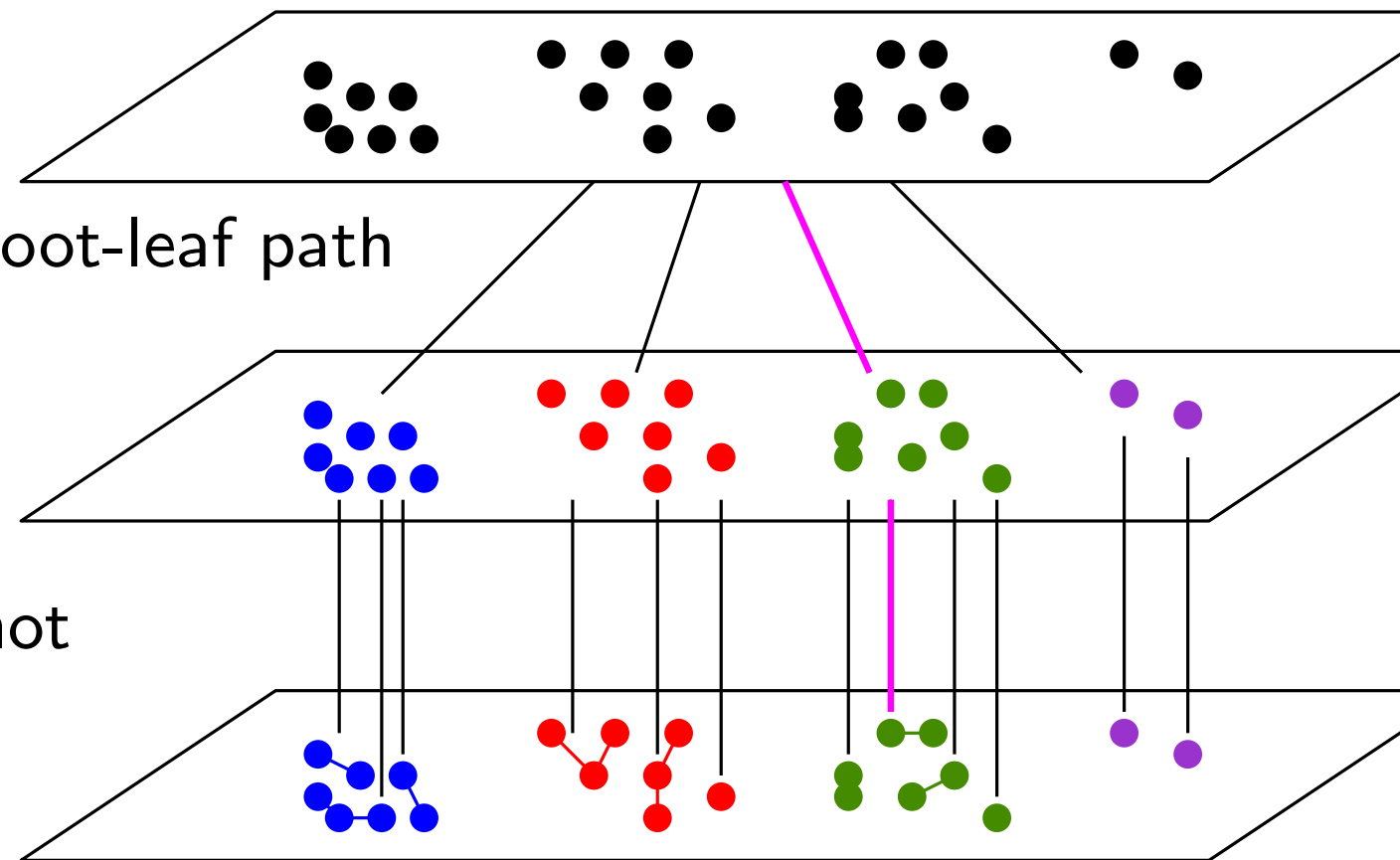
Assign  $P_v \subseteq P$  and  $[r_v, R_v]$  to each node  $v$

$\varepsilon$ -NN query:

- traverse down the tree along one root-leaf path
- at each visited node  $v$ , perform

$(r_v, \varepsilon)$ -NN and  $(R_v, \varepsilon)$ -NN queries

→ decide if  $d(q, P) \in [r_v, R_v]$  or not



# From $(r, \varepsilon)$ -NN to $\varepsilon$ -NN

- General case: use **hierarchical clustering tree** [Har-Peled'01]
  - consider geometric sequences of scales as before
  - cluster data points in order to bound the lengths of the sequences

Assign  $P_v \subseteq P$  and  $[r_v, R_v]$  to each node  $v$

$\varepsilon$ -NN query:

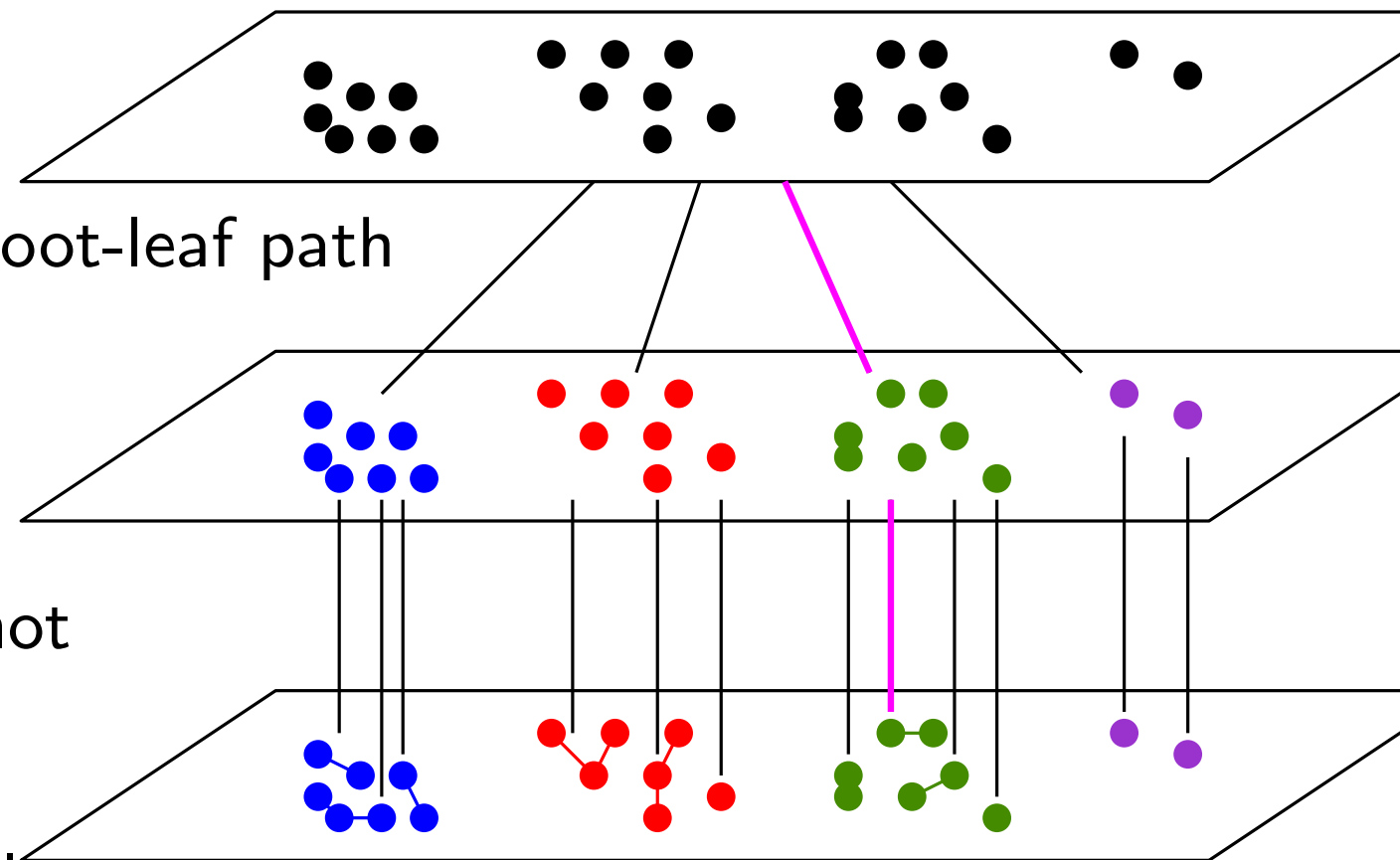
- traverse down the tree along one root-leaf path
- at each visited node  $v$ , perform

$(r_v, \varepsilon)$ -NN and  $(R_v, \varepsilon)$ -NN queries

→ decide if  $d(q, P) \in [r_v, R_v]$  or not

→ if so, locate  $d(q, P)$  in  $[r_v, R_v]$

→ if not, recurse into one child only



# From $(r, \varepsilon)$ -NN to $\varepsilon$ -NN

- General case: use **hierarchical clustering tree** [Har-Peled'01]
  - consider geometric sequences of scales as before
  - cluster data points in order to bound the lengths of the sequences

Assign  $P_v \subseteq P$  and  $[r_v, R_v]$  to each node  $v$

$\varepsilon$ -NN query:

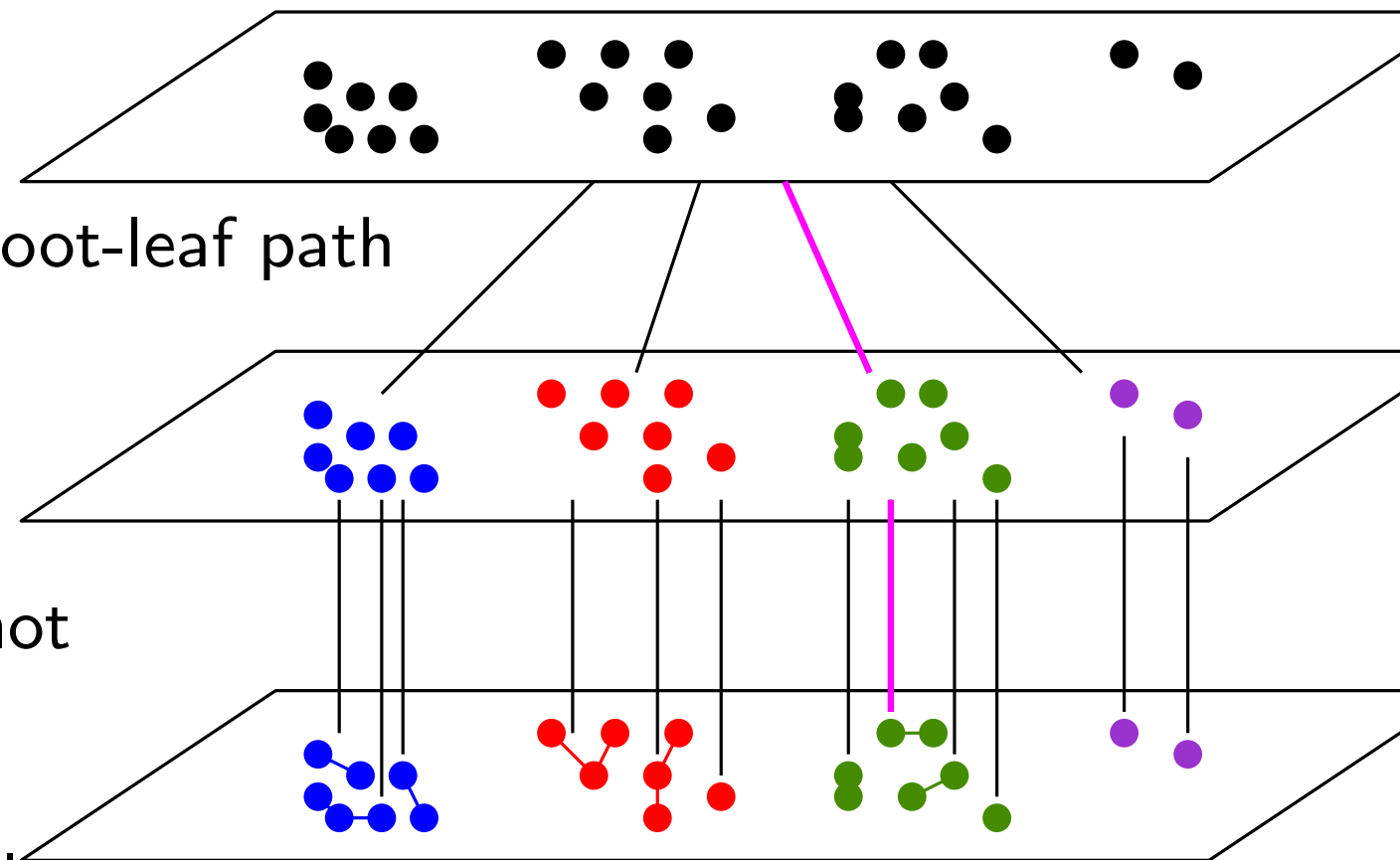
- traverse down the tree along one root-leaf path
- at each visited node  $v$ , perform

$(r_v, \varepsilon)$ -NN and  $(R_v, \varepsilon)$ -NN queries

→ decide if  $d(q, P) \in [r_v, R_v]$  or not

→ if so, locate  $d(q, P)$  in  $[r_v, R_v]$

→ if not, recurse into one child only



$O(\frac{1}{\varepsilon} \log \frac{n}{\varepsilon})$   $(r, \varepsilon)$ -NN queries per  $\varepsilon$ -NN query  $\Rightarrow O(\frac{1}{\varepsilon} n^\ell \log \frac{n}{\varepsilon})$  query time

# Take-Home Messages

- (Approximate) NN search requires an exponential amount of resources (space/time) in the algebraic comparison tree model [Arya et al. 98].
- Using random hashing allows to beat the *curse of dimensionality*.
- The price to pay is that algorithms become almost linear  
→ in practice, a trade-off must be found.
- The complexity of the exact NN search problem is not fully understood.  
→ what about *reverse* NN search? [Cheong et al. 09], [Arthur, O. 10], ...