

5800

*data structures*

apr8/apr11 2022  
shelat

# Dictionary

data structure

Key-value mapping.

Insert (key, value) : creates the association b/w key-value.

Lookup (key) : returns (key, value) if key was previously inserted into the data structure.

Delete (key) : removes key.

Find next (key) : finds the "next" lexicographic <sup>(pair)</sup> key in the data structure that is greater than (key).

# DICTIONARY

insert(key, value)

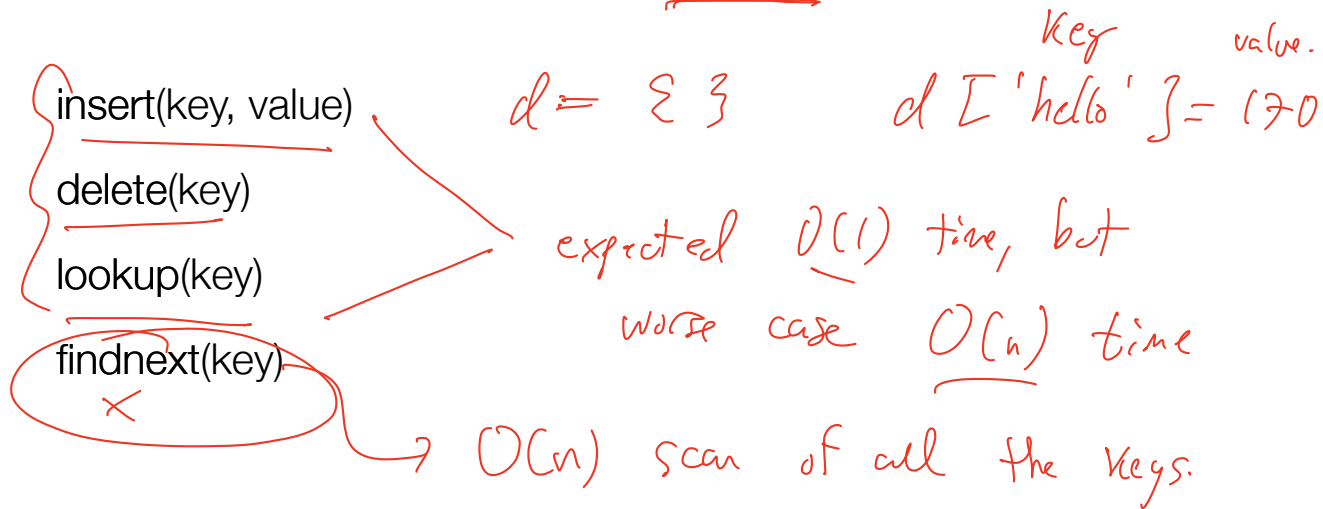
delete(key)

lookup(key)

findnext(key)

# DICTIONARY

standard solution: hashtable



# Python DICTIONARY $N = 2^k$

hash function  $h$

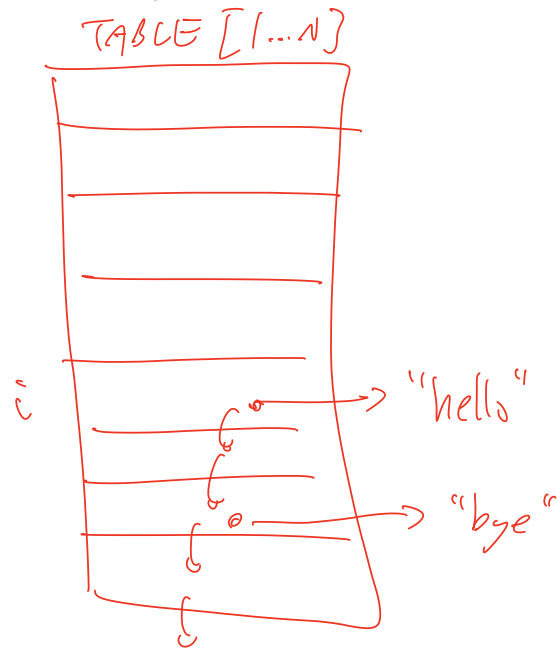
$$h(\text{key}) \rightarrow [0, N-1]$$

$$h(\text{"hello"}) = i$$

$$h(\text{"bye"}) = i$$

might be

$$(5 * i + p + 1) \% N.$$



# Hashtables are tricky

```
1
2 import time
3 import sys
4 import d
5
6 dd = {}
7
8 # make a dictionary with elements from the list
9 for l in d.list:
10     dd[l] = l
11
12 def lookup(v):
13     start = time.time()
14     t = 0
15     for j in range(10000):
16         if v in dd:
17             t = t + 1
18     end = time.time()
19     print(end - start)
20     return t
21
```

*empty dictionary*

*insert (k, v) into the dictionary*

*lookup(v)*

# Hashtables are tricky

```
1
2 import time
3 import sys
4 import d
5
6 dd = {}
7
8 # make a dictionary with elements from the list
9 for l in d.list:
10     dd[l] = l
11
12 def lookup(v):
13     start = time.time()
14     t = 0
15     for j in range(10000):
16         if v in dd:
17             t = t + 1
18     end = time.time()
19     print(end - start)
20     return t
21
```

This is a trivial lookup experiment.  
Looking up 1 key takes 2000x longer.

```
MacBook-Pro-2:hashing abhi$ python3 bad.py
size of dictionary: 43689
Starting experiment to lookup 1000:
0.0005161762237548828
Starting experiment to lookup 100000:
1.0303189754486084
MacBook-Pro-2:hashing abhi$
```

# Hashtables are tricky

2003

```
1
2 import time
3 import sys
4 import d
5
6 dd = {}
7
8 # make a dictionary with elements from the list
9 for l in d.list:
10     dd[l] = l
11
12 def lookup(v):
13     start = time.time()
14     t = 0
15     for j in range(10000):
16         if v in dd:
17             t = t + 1
18     end = time.time()
19     print(end - start)
20     return t
21
```

This is a trivial lookup experiment.  
Looking up 1 key takes 2000x longer.

```
MacBook-Pro-2:hashing abhi$ python3 bad.py
size of dictionary: 43689
Starting experiment to lookup 1000:
0.0005161762237548828
Starting experiment to lookup 100000:
1.0303189754486084
MacBook-Pro-2:hashing abhi$ █
```

Worst case performance:  $O(n)$





# DICTIONARY

new constraint: keys belong to limited range:

$\{1, \dots, n\}$

insert(key, value) :  $O(\log \log n)$

delete(key)

"

worst case

lookup(key)

"

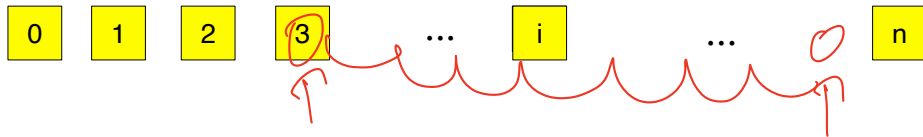
performance !!

findnext(key)

"

# A simple solution: bit vector

Maintain an array of bits



insert(key, value) : *easy*  $O(1)$

delete(key)  $O(1)$

lookup(key)  $O(1)$

findnext(key)  $O(n)$

*min(-)*

*max(-)*

CAN WE DO BETTER THAN  $O(N)$  FINDNEXT?

# van emde Boas Q

## THE BIG IDEA:

(1) recursive data structure

the data structure that handles the universe  $[1 \dots n]$  consists of

many smaller data structures that handle the universe  $[1 \dots \sqrt{n}]$

# van emde Boas Q

## THE BIG IDEA:

Use recursion for a data structure.

A data structure that handles  $1..n$  can be designing using several smaller versions of the same structure.

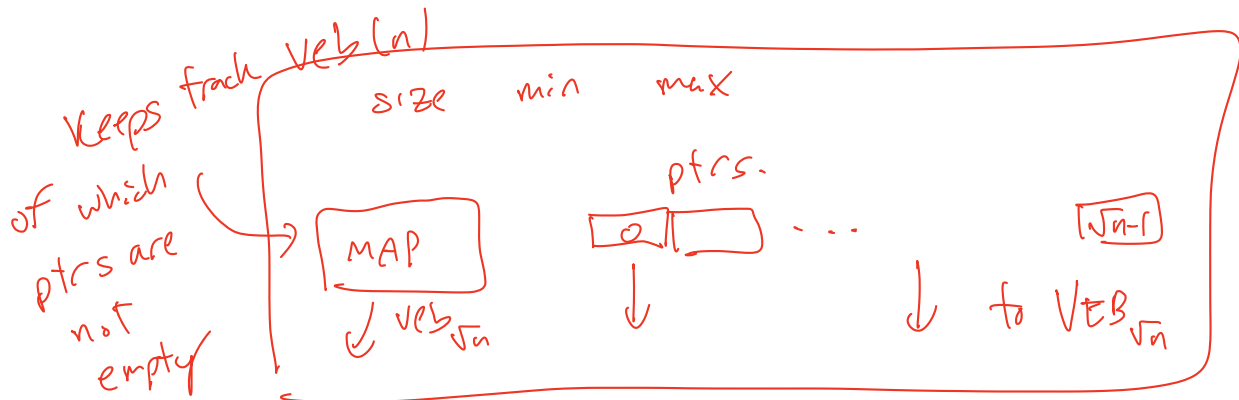
# VEB queue

$$N - 2^k$$

VEB<sub>(N)</sub>      size      min      max

Base case: 2 bit vector for 1...2

recursive case:



# VEB queue

$\text{VEB}_{(N)}$

SZ, MIN, MAX

# VEB queue

$\text{VEB}_{(N)}$

SZ, MIN, MAX

BASE CASE: 1 BIT VECTOR.



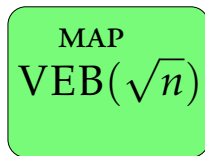
# VEB queue

$VEB_{(N)}$

SZ, MIN, MAX

BASE CASE: 4 BIT VECTOR.

NORMAL CASE:



Pointers to recursive, smaller instances of VEB.



Keeps track of which ptrs  
are not empty.

EXAMPLE:

$n = 256$

$VEB_{(n)}$   
SZ, MIN, MAX

map  
 $VEB(\sqrt{n})$

0

1

2

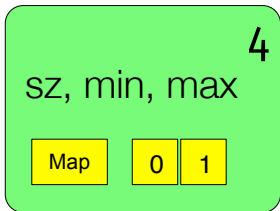
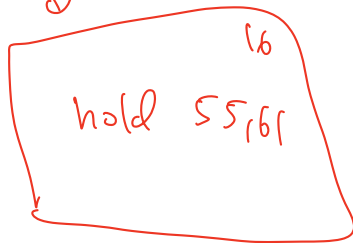
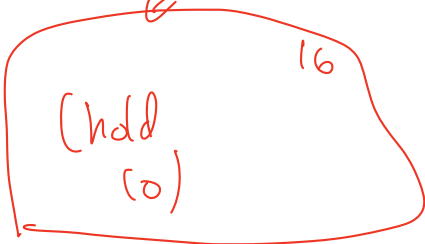
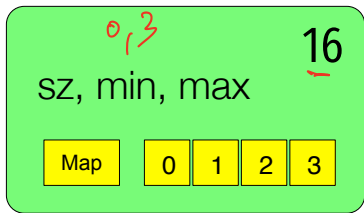
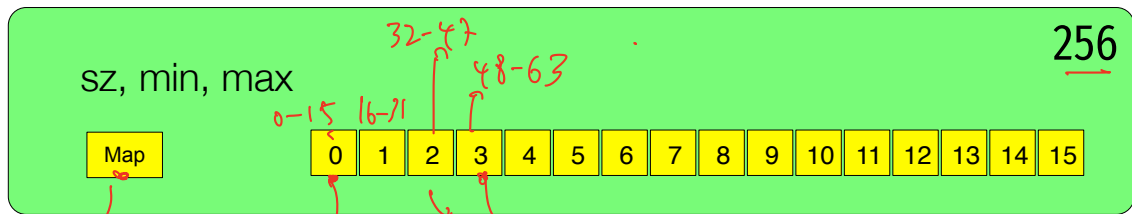
3

...

$\sqrt{n}$

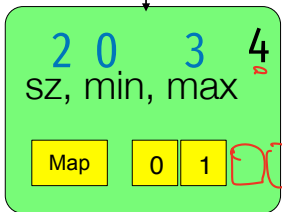
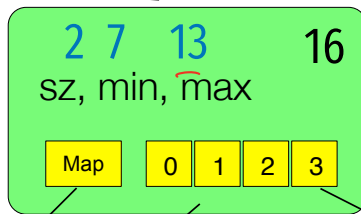
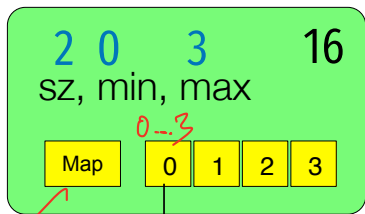
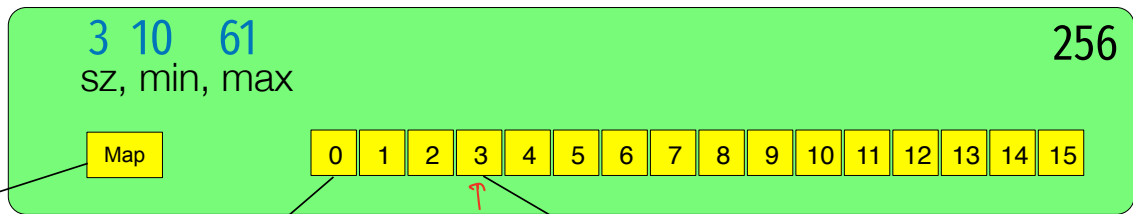
16

# Example $n=256$ , keys=10,55,61

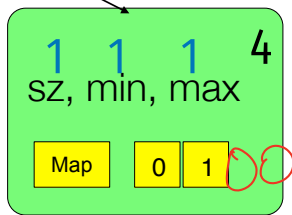
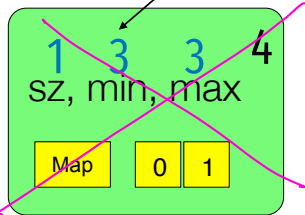


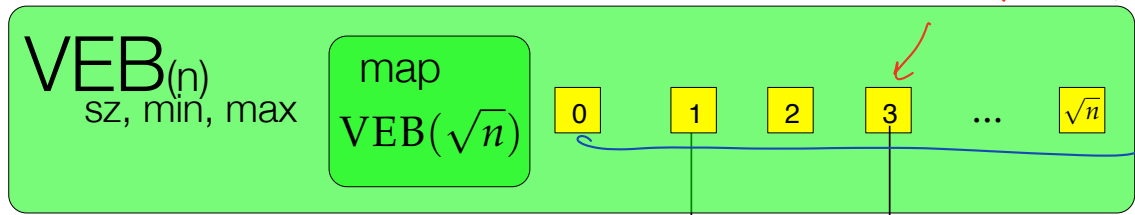
# Example $n=256$ , keys={10,55,61}

48



(Roughly correct)





LOOKUP(i)

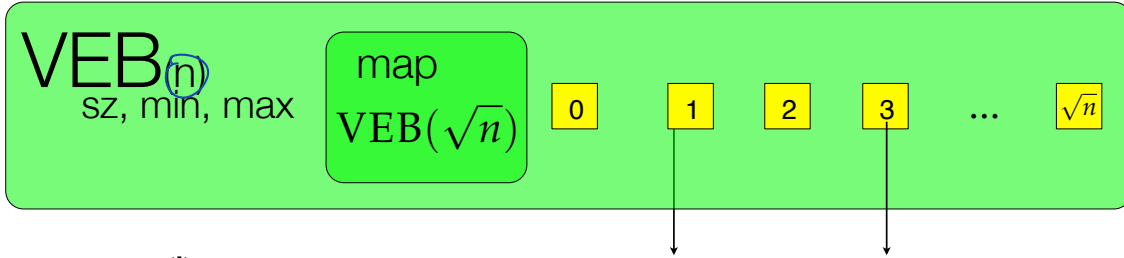
o(1) write  $i = a \cdot \sqrt{n} + b$  such that  $a, b \in [0.. \sqrt{n}-1]$

o(1) if BASE case then check the bit vector.

if size = 0 or a.size = 0 return false.

else return a.lookup(b)

$$T(n) = T(\sqrt{n}) + \Theta(1) = \Theta(\log \log n)$$



LOOKUP(i)

WRITE  $i = a\sqrt{n} + b$

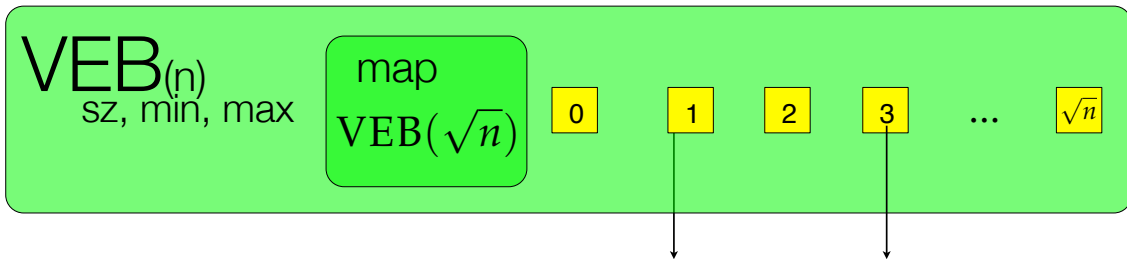
IF <BASE CASE>: CHECK BIT VECTOR

IF SIZE = 0 OR **a**.SIZE = 0 THEN RETURN FALSE

ELSE RETURN **a**.LOOKUP(b)

*N<sup>-</sup>  
size  
of  
the  
universe.*

(Almost right, we will have to slightly change this later.)



LOOKUP(*i*)

WRITE  $i = a\sqrt{n} + b$

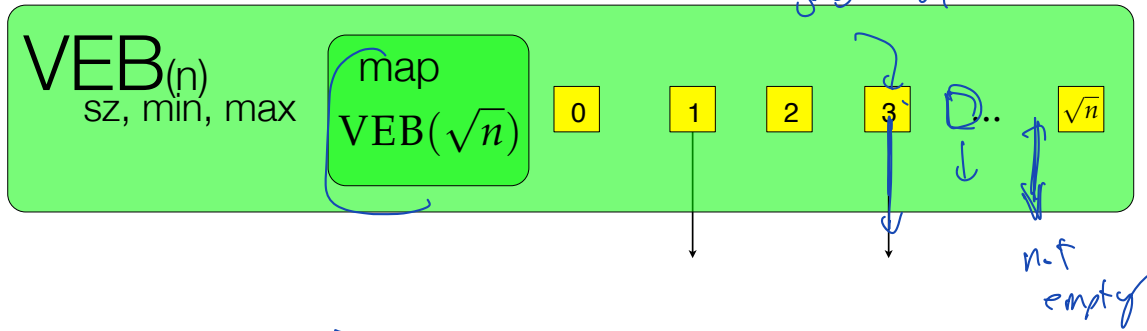
IF <BASE CASE>: CHECK BIT VECTOR

IF SIZE = 0 OR **a**.SIZE = 0 THEN RETURN FALSE

ELSE RETURN **a**.LOOKUP(*b*)

Running time:  $T(n) = T(\sqrt{n}) + \Theta(1) = \Theta(\log \log n)$

(Almost right, we will have to slightly change this later.)



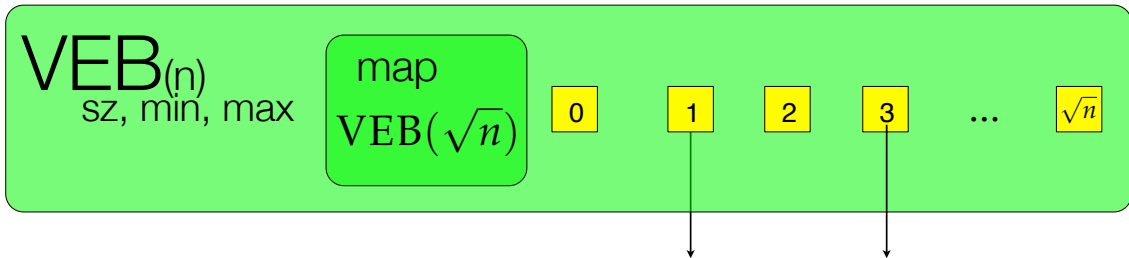
FINDNEXT(i)

IDEA:

write  $e = a\sqrt{n} + b$ .

- 2 cases -
- ① either bucket  $a$  contains the next key
  - ② or the next bucket at this level contains the next key





FINDNEXT(*i*)

IDEA:

Write  $i = a\sqrt{n} + b$  as usual.

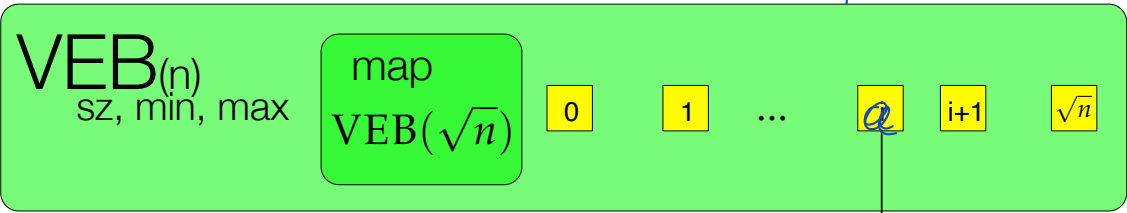
Case 1: Bucket *a* has the next value.

Recursively use findnext<sub>*a*</sub>(*b*)

Case 2: Bucket *a* does not have the next value.

Use  $x = \text{findnext}_{\text{map}}(a)$ , return  $x.\text{min}$ .

55,61



FINDNEXT(i)

ec1) write  $i = a \cdot \sqrt{n} + b$

< handle base case with scanning >

IF  $a.max \geq b$  then

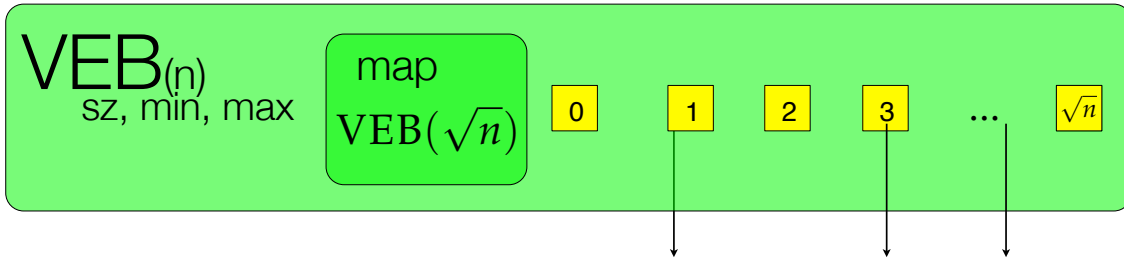
return  $a.findnext(b)$

else

return  $map.findnext(a)$ .min

if no such entry, return  $\infty$ .

$$T(a) = T(\sqrt{n}) + O(1)$$



FINDNEXT(i)

WRITE  $i = a\sqrt{n} + b$

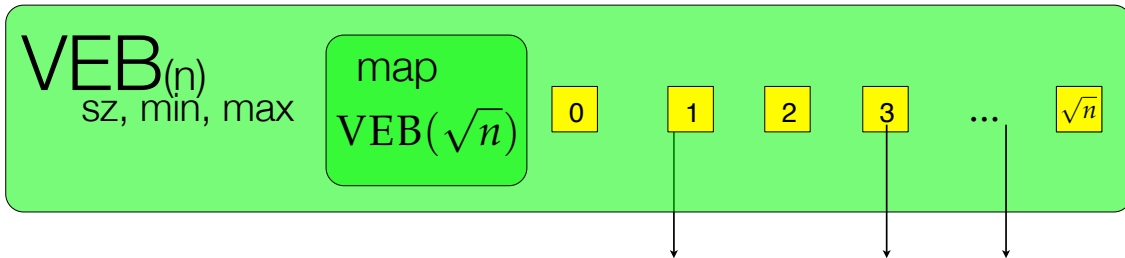
<BASE CASE IF SIZE IS ZERO>

IF  $a$ .MAX  $> b$  THEN

RETURN  $a$ .FINDNEXT(b)

ELSE

RETURN MAP.FINDNEXT(a).MIN



FINDNEXT(i)

WRITE  $i = a\sqrt{n} + b$

<BASE CASE IF SIZE IS ZERO>

IF  $a$ .MAX  $> b$  THEN

RETURN  $a$ .FINDNEXT(b)

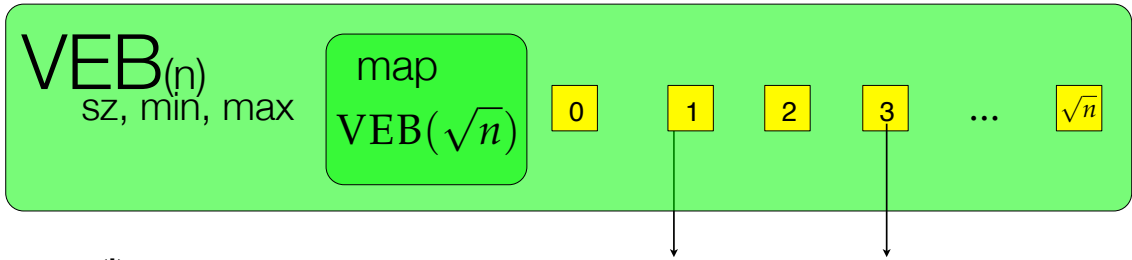
ELSE

RETURN  $MAP$ .FINDNEXT(a).MIN

Running time:

$$T(n) = T(\sqrt{n}) + \Theta(1)$$

$$\Theta(\log \log n)$$

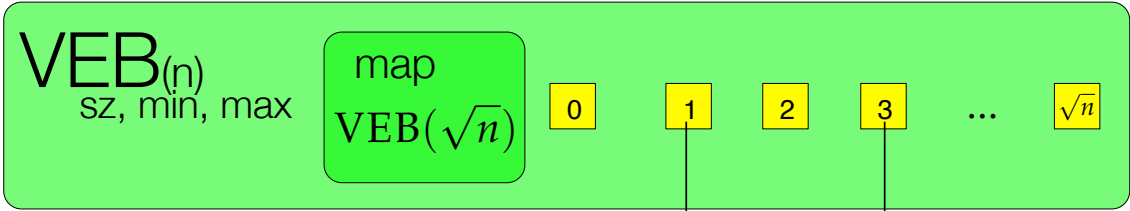


INSERT(*i*)

WRITE  $i = a\sqrt{n} + b$

*a.insert(b)*

*map.insert(a)*



INSERT(i)

WRITE  $i = a\sqrt{n} + b$

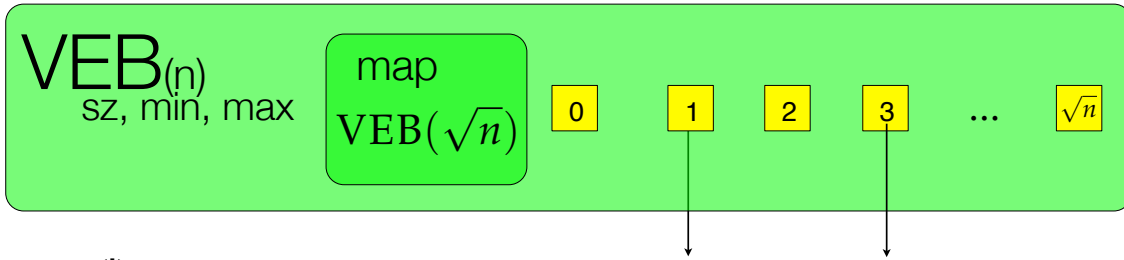
A.INSERT(B)

MAP.INSERT(A)

$$T(n) = 2T(\sqrt{n}) + \Theta(1)$$

$$= \Theta(\log n)$$

It is too much work to perform 2 inserts here!!



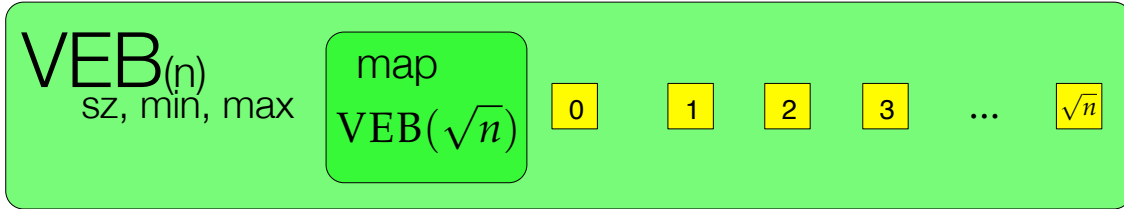
INSERT(*i*)

WRITE  $i = a\sqrt{n} + b$

A.INSERT(*B*)

MAP.INSERT(*A*)

WHAT IS THE PROBLEM WITH THIS?



INSERT(i)

WHAT IS THE PROBLEM WITH THIS?

WRITE  $i = a\sqrt{n} + b$

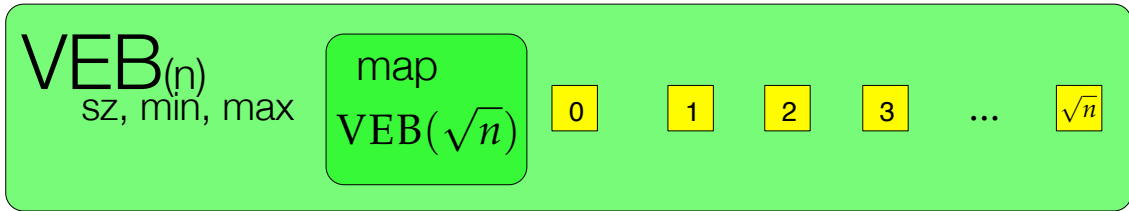
A.INSERT(B)

MAP.INSERT(A)

HOW CAN WE GET AROUND THE PROBLEM OF  
INSERTING TWICE?

ANSWER: LAZY INSERTS. HOW MANY TIMES DO WE NEED  
TO INSERT INTO MAP?





INSERT(i)

WRITE  $i = a\sqrt{n} + b$

IF SZ==0 THEN update  $SZ = 1$   $min = max = i$

ELSE

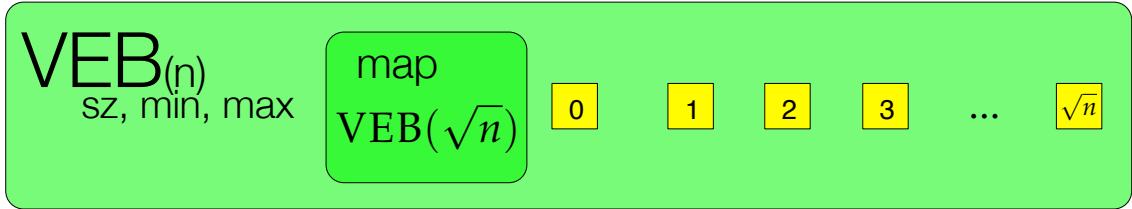
if  $min > i$  swap( $i, min$ )

write  $i = a\sqrt{n} + b$

if  $a.SZ == 0$  then map.insert( $a$ )

$a.insert(b)$ .

update  $max, sz$ .



INSERT(i)

IF  $SZ == 0$  THEN UPDATE  $SZ = 1, MIN = MAX = i$

ELSE

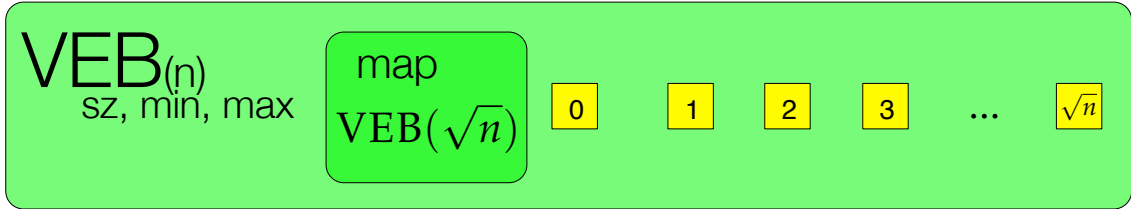
IF  $MIN > i$  SWAP(i, MIN)

WRITE  $i = a\sqrt{n} + b$

IF  $a$ .SZ == 0 THEN MAP.INSERT(a).

$a$ .INSERT(b)

UPDATE SZ, MAX



INSERT(i)

IF SZ==0 THEN UPDATE SZ=1, MIN=MAX=i

ELSE

IF MIN>i SWAP(i, MIN)

WRITE  $i = a\sqrt{n} + b$

IF **a**.SZ==0 THEN **MAP**.INSERT(a).

**a**.INSERT(b)

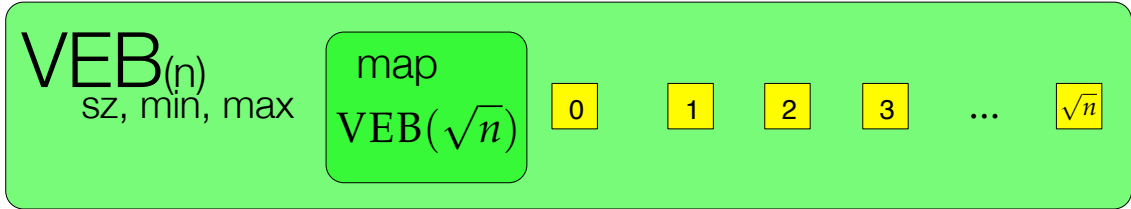
UPDATE SZ, MAX

If a is empty:  
then 1 full recursive call + 1 base case

*If a is not empty=*

*← we just run this case*

$$T(u) = T(\sqrt{u}) + \Theta(1) = \Theta(\log \log u)$$



INSERT(i)

IF SZ==0 THEN UPDATE SZ=1, MIN=MAX=i

ELSE

IF MIN>i SWAP(i, MIN)

WRITE  $i = a\sqrt{n} + b$

IF  $a$ .SZ==0 THEN MAP.INSERT(a).

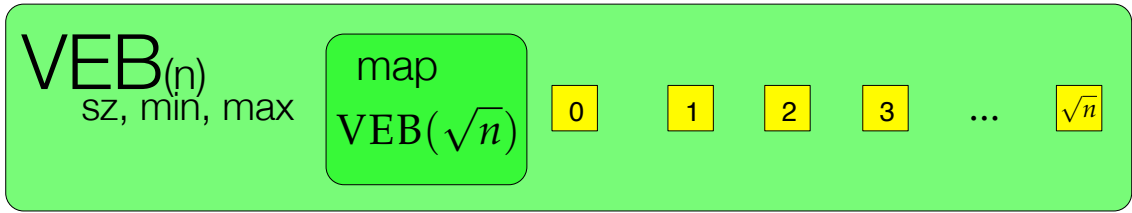
$a$ .INSERT(b)

UPDATE SZ, MAX

If a is empty:  
then 1 full recursive call + 1 base case

If a is not empty:  
Then this line does not run  
but 1 full recursive call is made





INSERT(i)

IF SZ==0 THEN UPDATE SZ=1, MIN=MAX=i

ELSE

IF MIN>i SWAP(i,MIN)

WRITE  $i = a\sqrt{n} + b$

IF **a**.SZ==0 THEN **MAP**.INSERT(a).

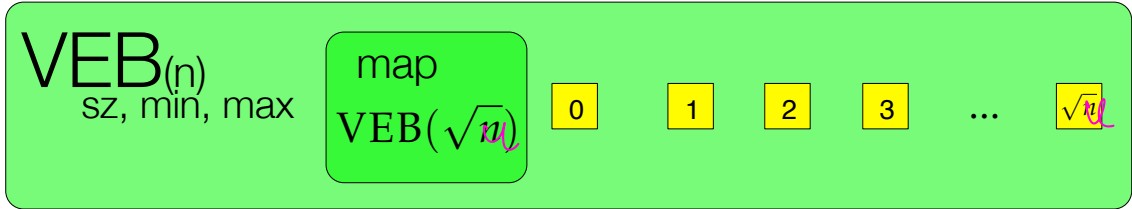
**a**.INSERT(b)

UPDATE SZ, MAX

If a is empty:  
then 1 full recursive call + 1 base case

If a is not empty:  
Then this line does not run  
but 1 full recursive call is made

$$T(n) = T(\sqrt{n}) + \Theta(1)$$



LOOKUP(i)

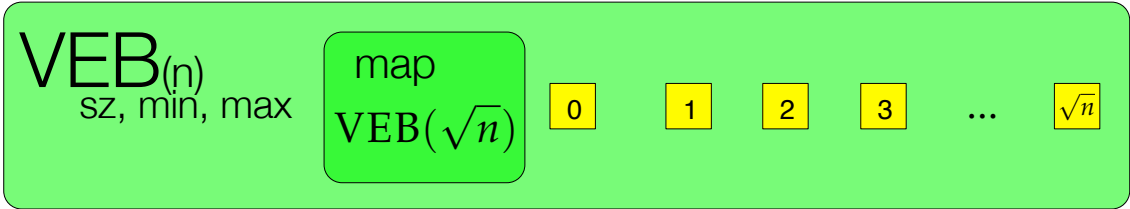
WRITE  $i = a\sqrt{n} + b$

If  $sz == 0$  return false

If  $i == min$  return true

else return  $a \cdot lookup(b)$

We need to fix the Lookup to work with Lazy inserts.



LOOKUP(i)

WRITE  $i = a\sqrt{n} + b$

IF SIZE==0 RETURN FALSE

IF  $i$ ==MIN RETURN TRUE

ELSE RETURN  $a$ .LOOKUP(b)

$u = \lceil 2^u \rceil \quad u \sim 2^{\frac{1}{2}}$

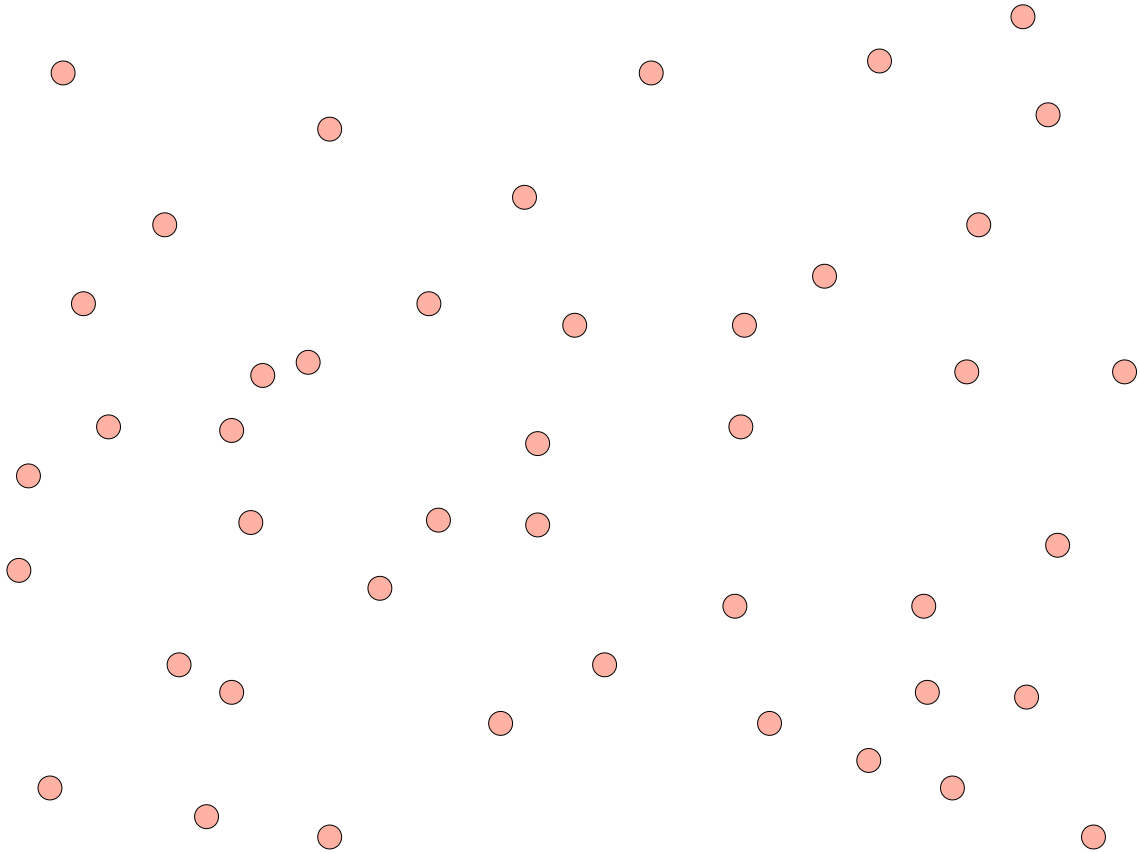
$\log \log(u) = 5$

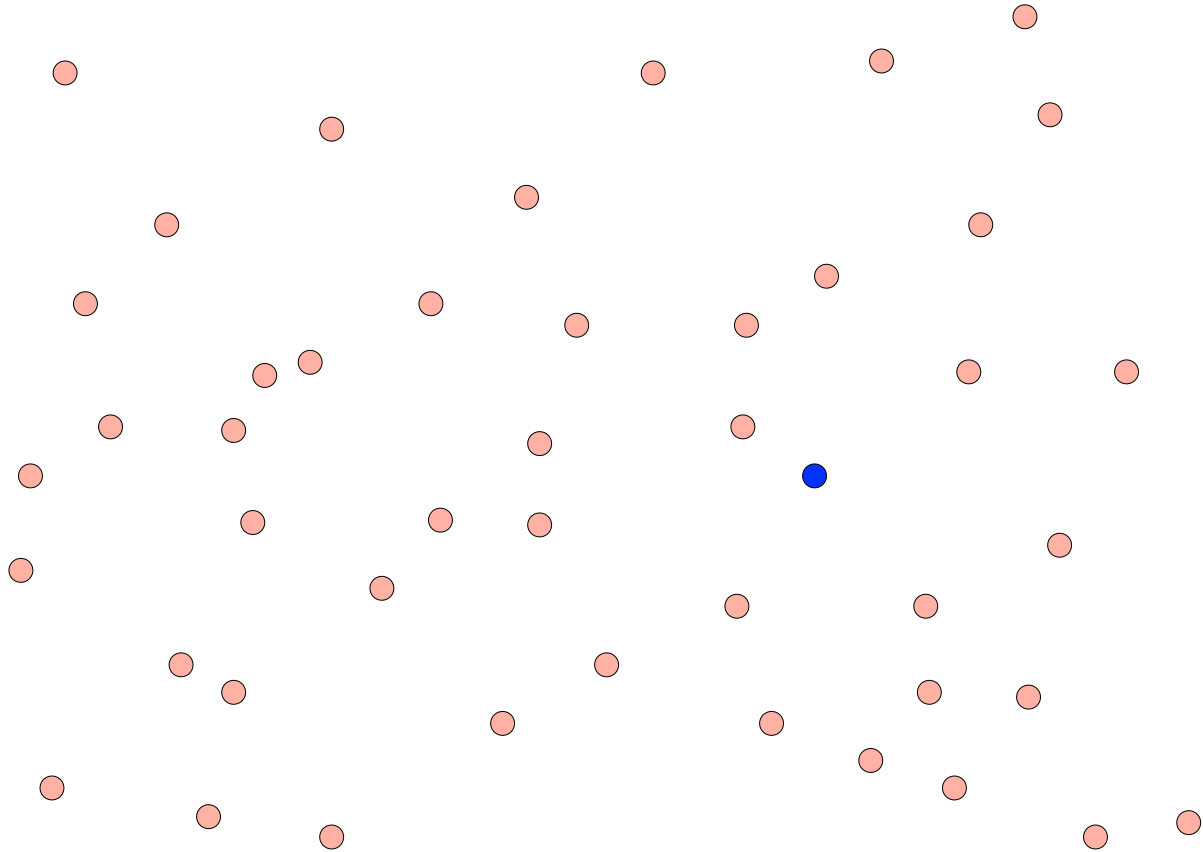
$\log(b)$   
 $\lceil \log \log(u) \rceil = 5$

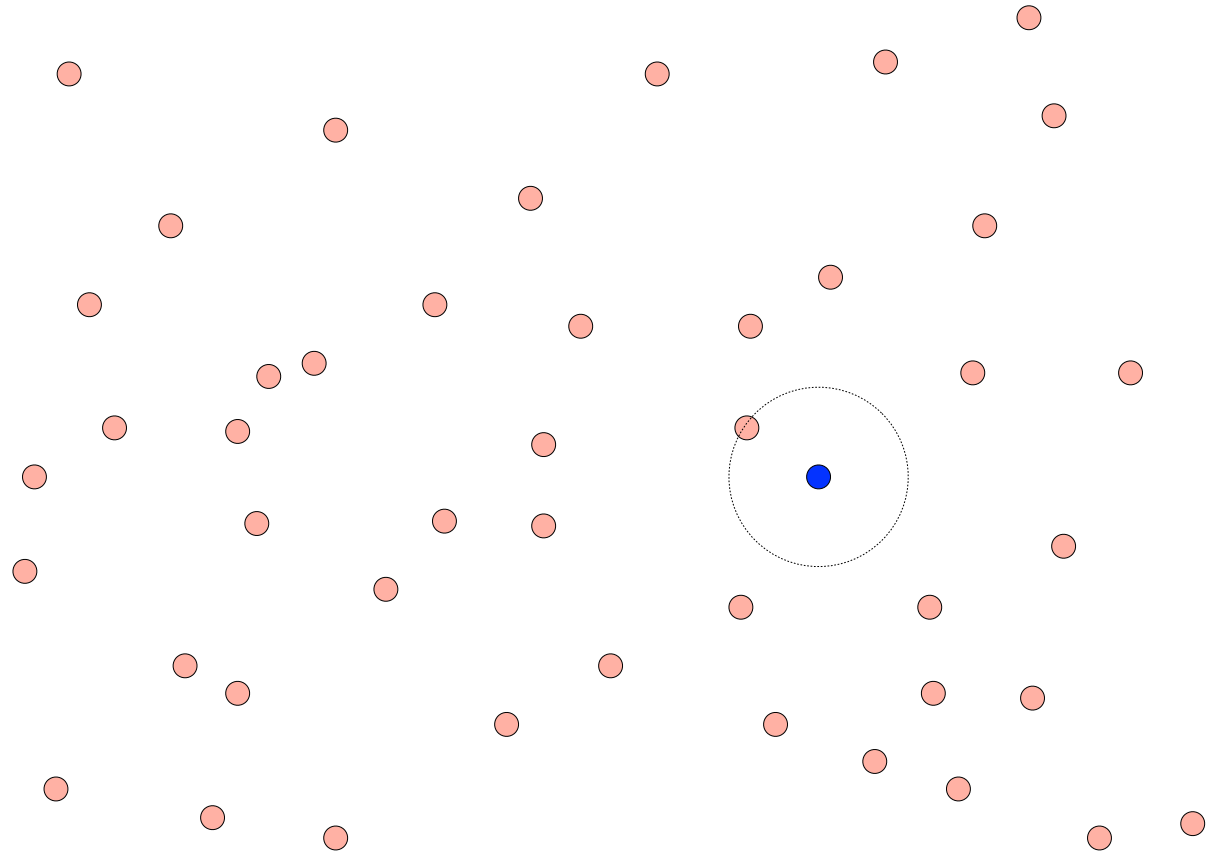
We need to fix the Lookup to work with Lazy inserts.

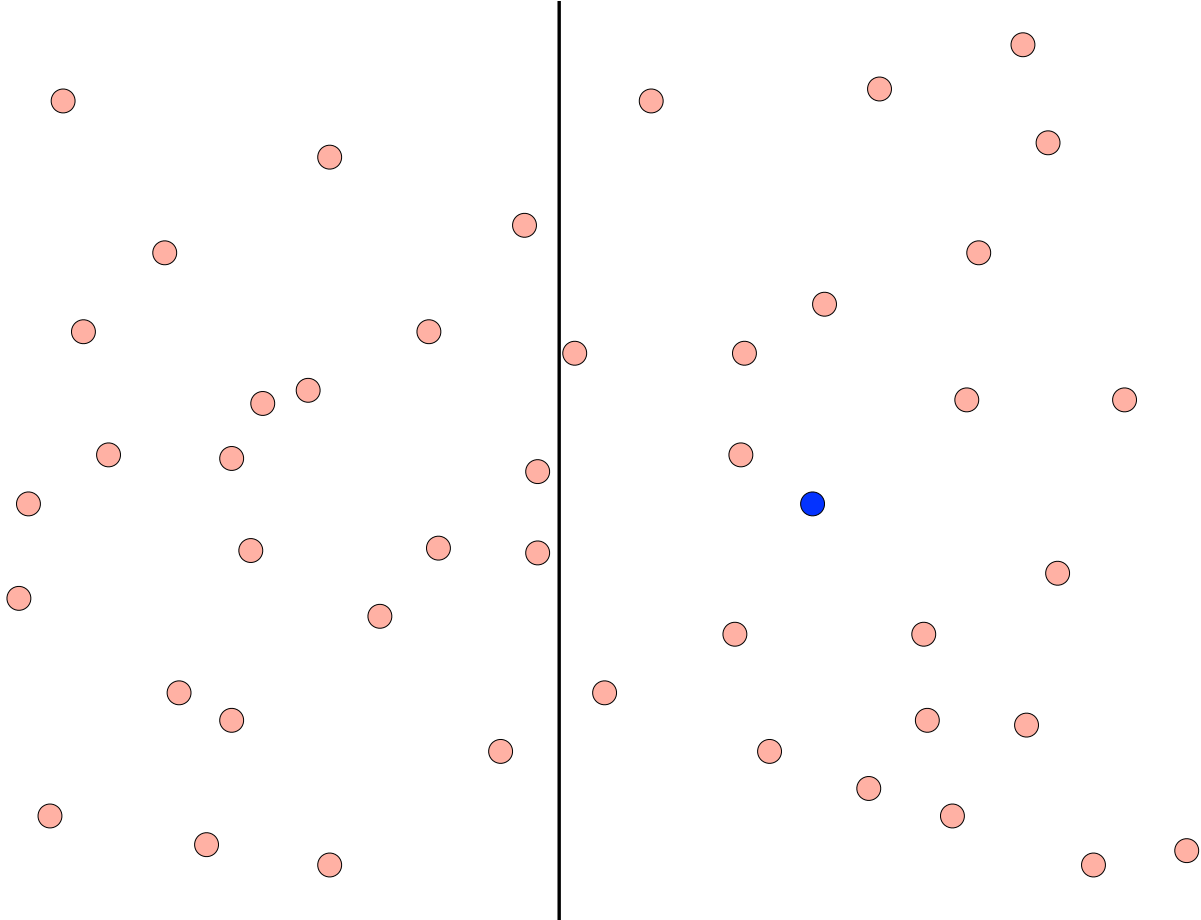
Nearest  
neighbor  
queries

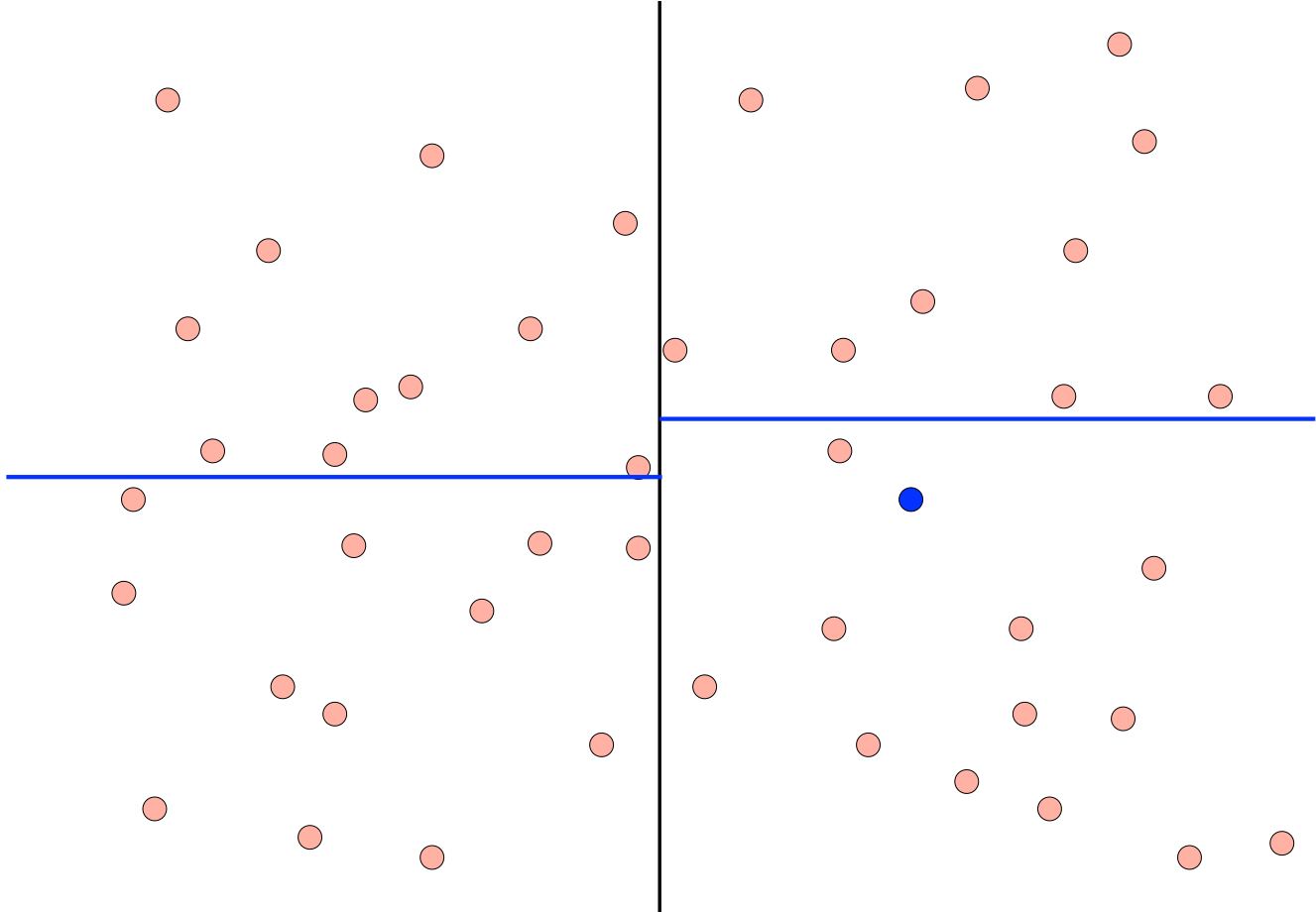


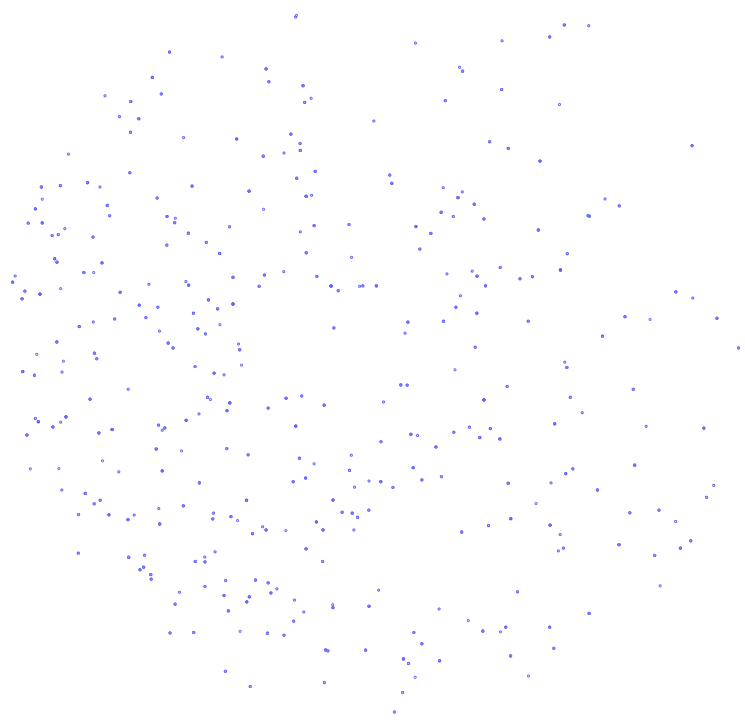


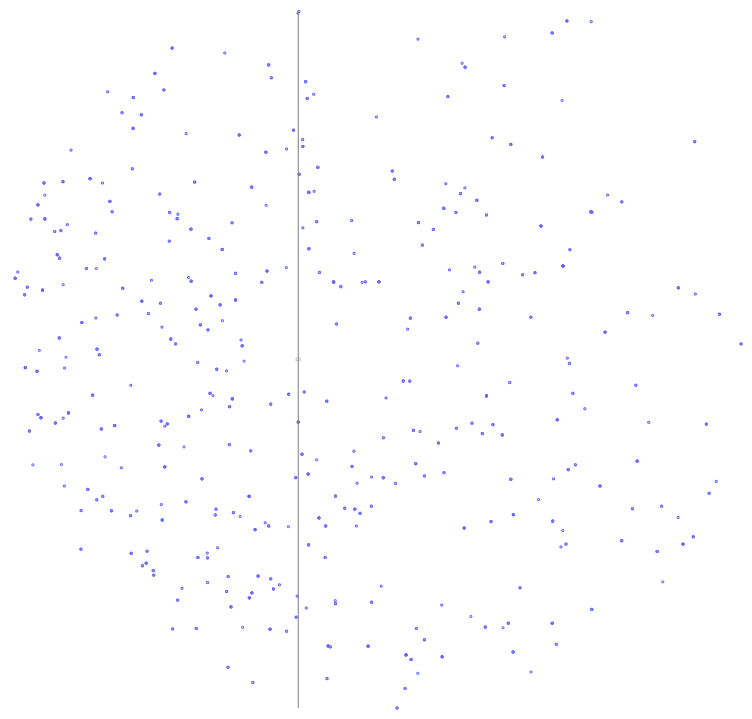


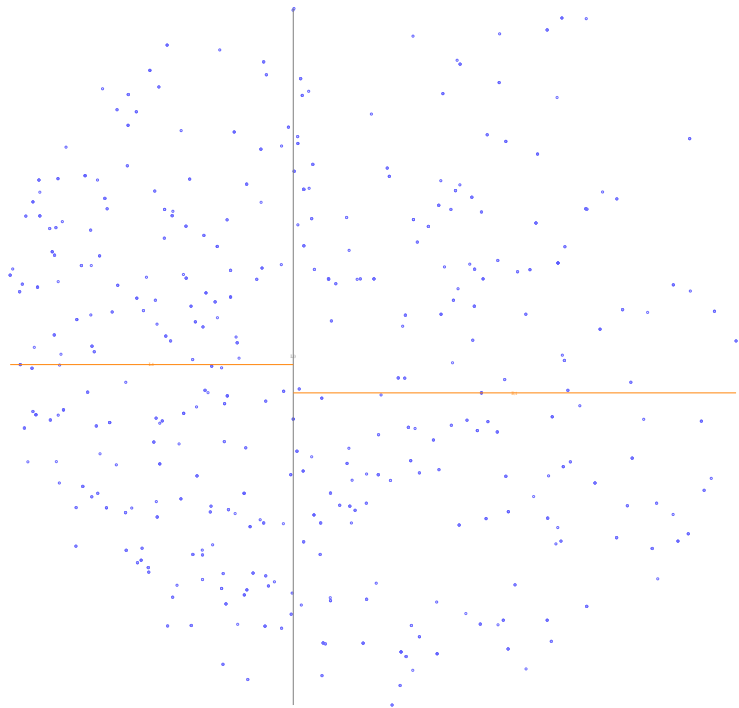




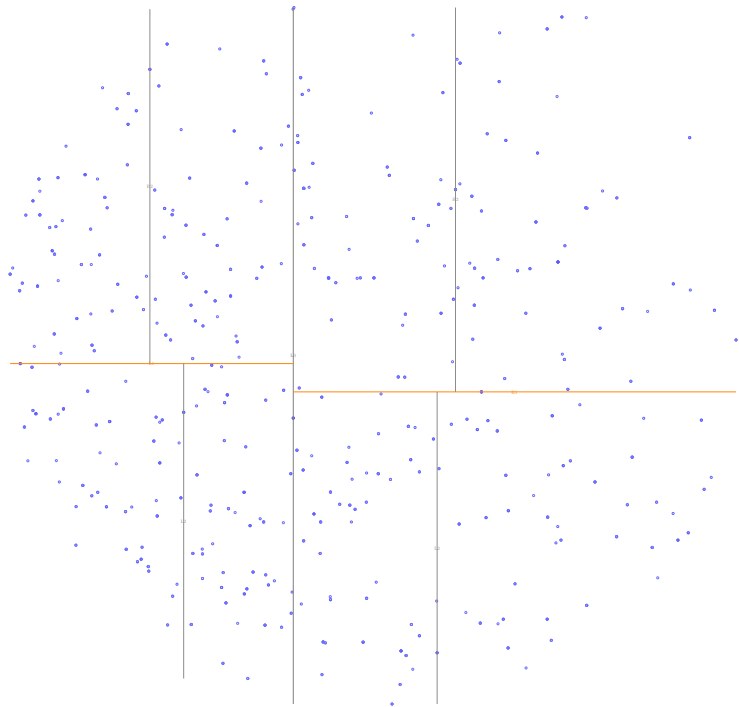


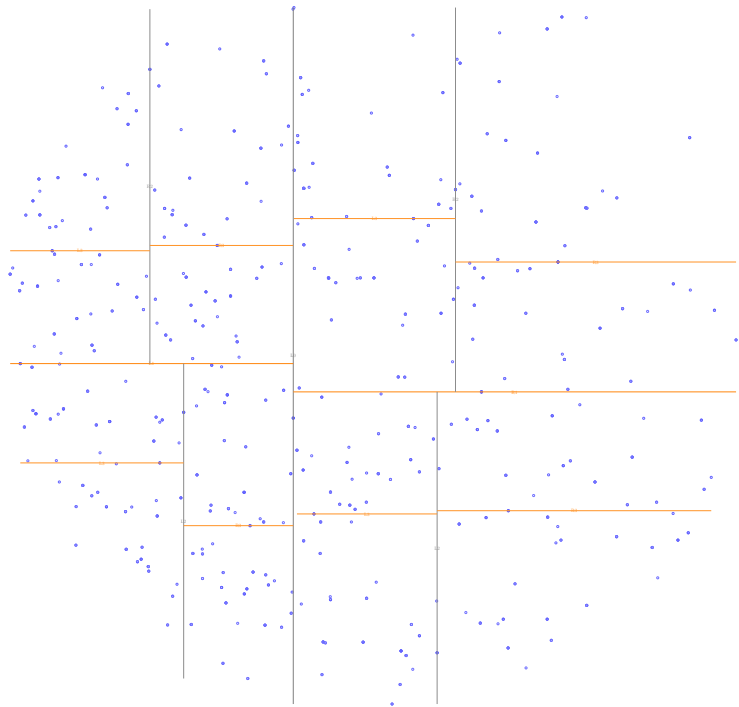


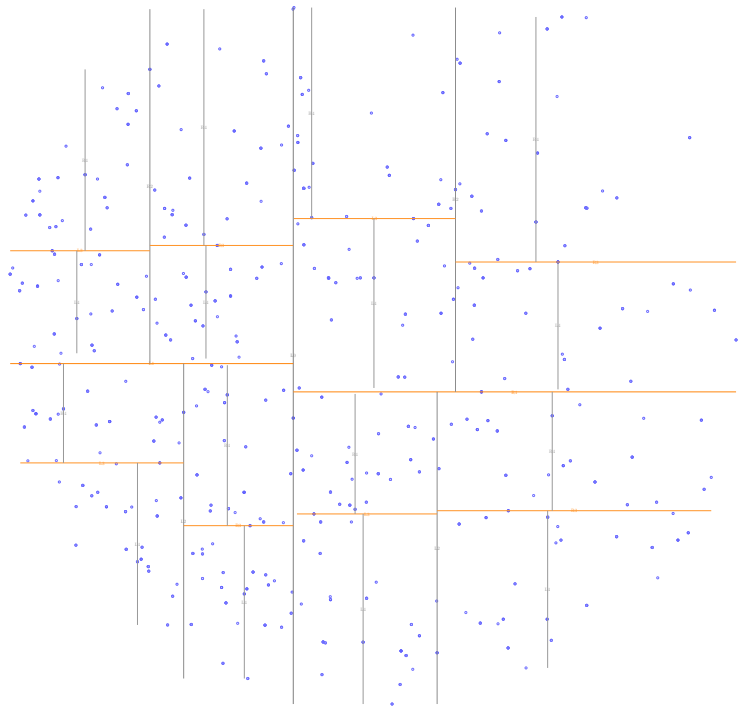


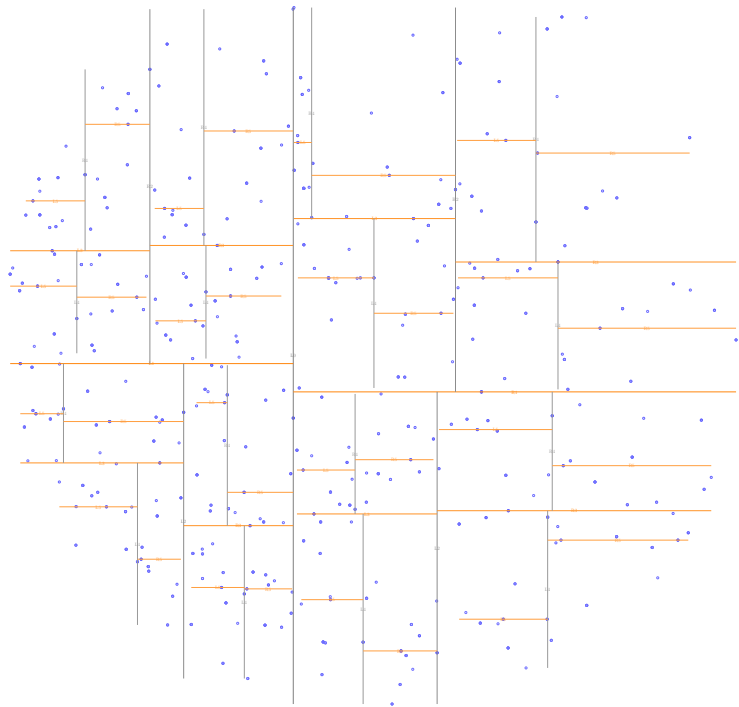








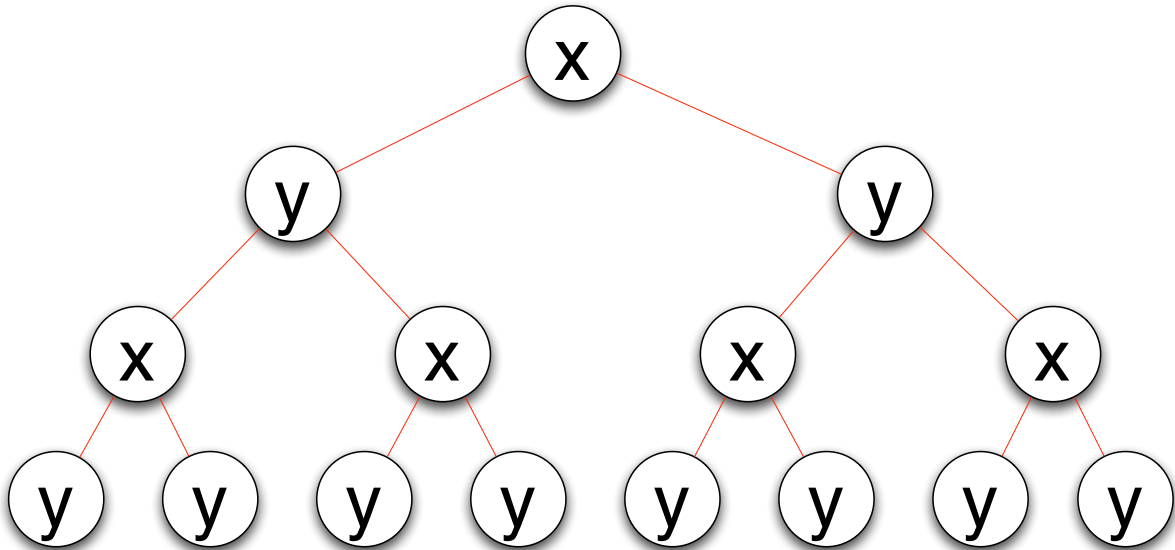


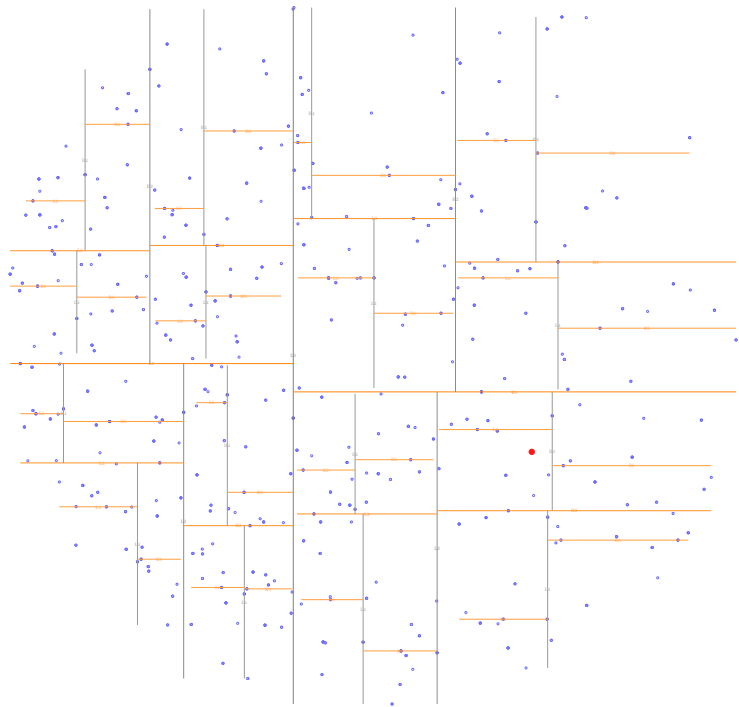


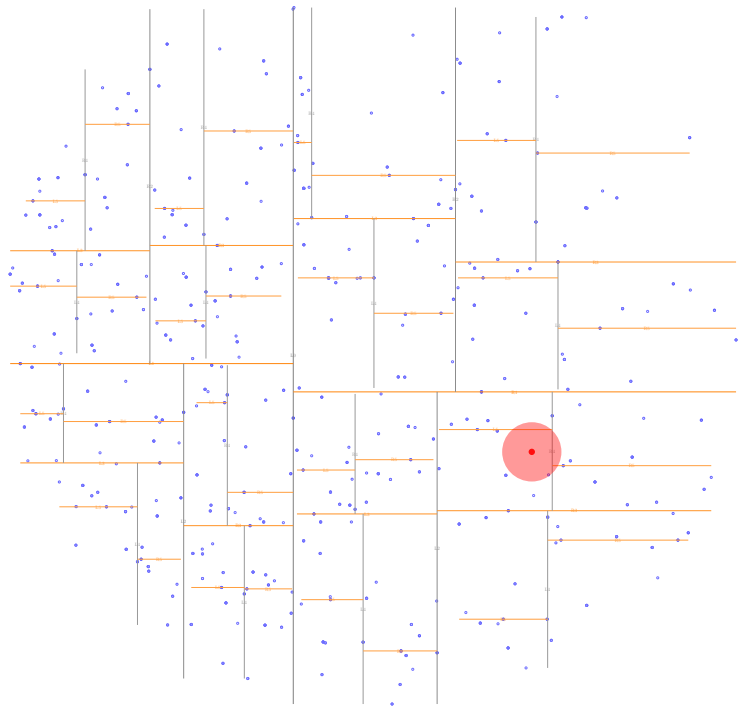
# KD-Tree

Each node in tree maintains variable “box”

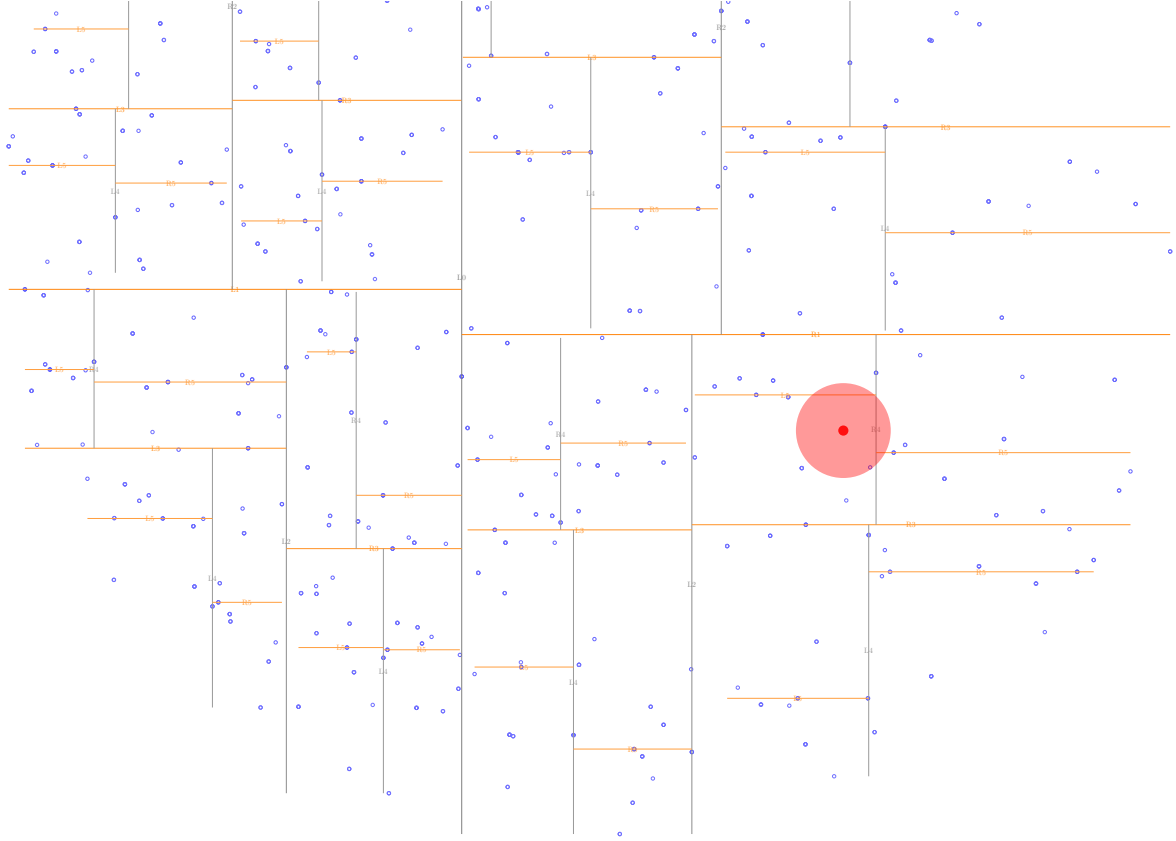
```
node {  
  rect box  
  point split  
  node* left  
  node* right  
}
```

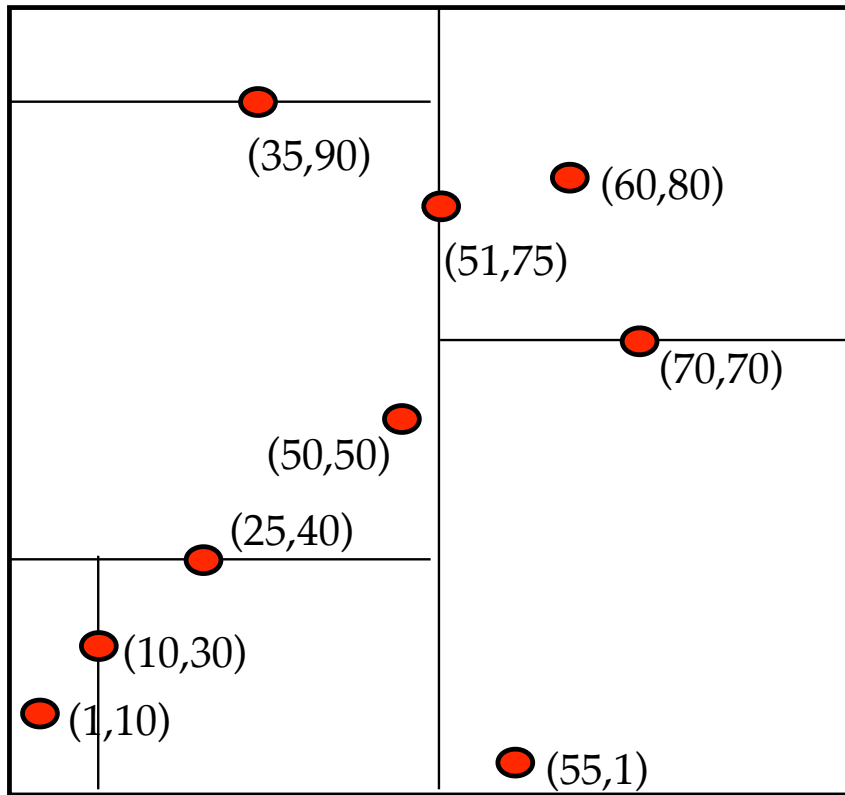












NN(q, tree, dir, closest-so-far)

if empty(tree) or dist(q, tree.box) > closest return

if dist(q, tree.root) < closest { update closest }

if q.dir < tree.dir {

    NN(q, tree.left, nextdir, closest)

    NN(q, tree.right, nextdir, closest)

} else {

    NN(q, tree.left, nextdir, closest)

    NN(q, tree.right, nextdir, closest)

}