

UNIVERSITY OF ILLINOIS  
AT URBANA-CHAMPAIGN

# CS411 - Extreme Indexing



[illinois.edu](http://illinois.edu)

# Announcements

- Get MP2 finished
- Enjoy your break!



# Extreme Indexing

- Indexing tricks we've used so far are great for relational data.
- What if we wanted to query other kinds of data quickly?



# Examples

- Where is the nearest grocery store?
- What coffee shop am I in now?
- What web pages contain the words “Belle” and “Sebastian”?
- Where in the human genome (a string of 3 billion characters) does the string “GATTACA” appear?

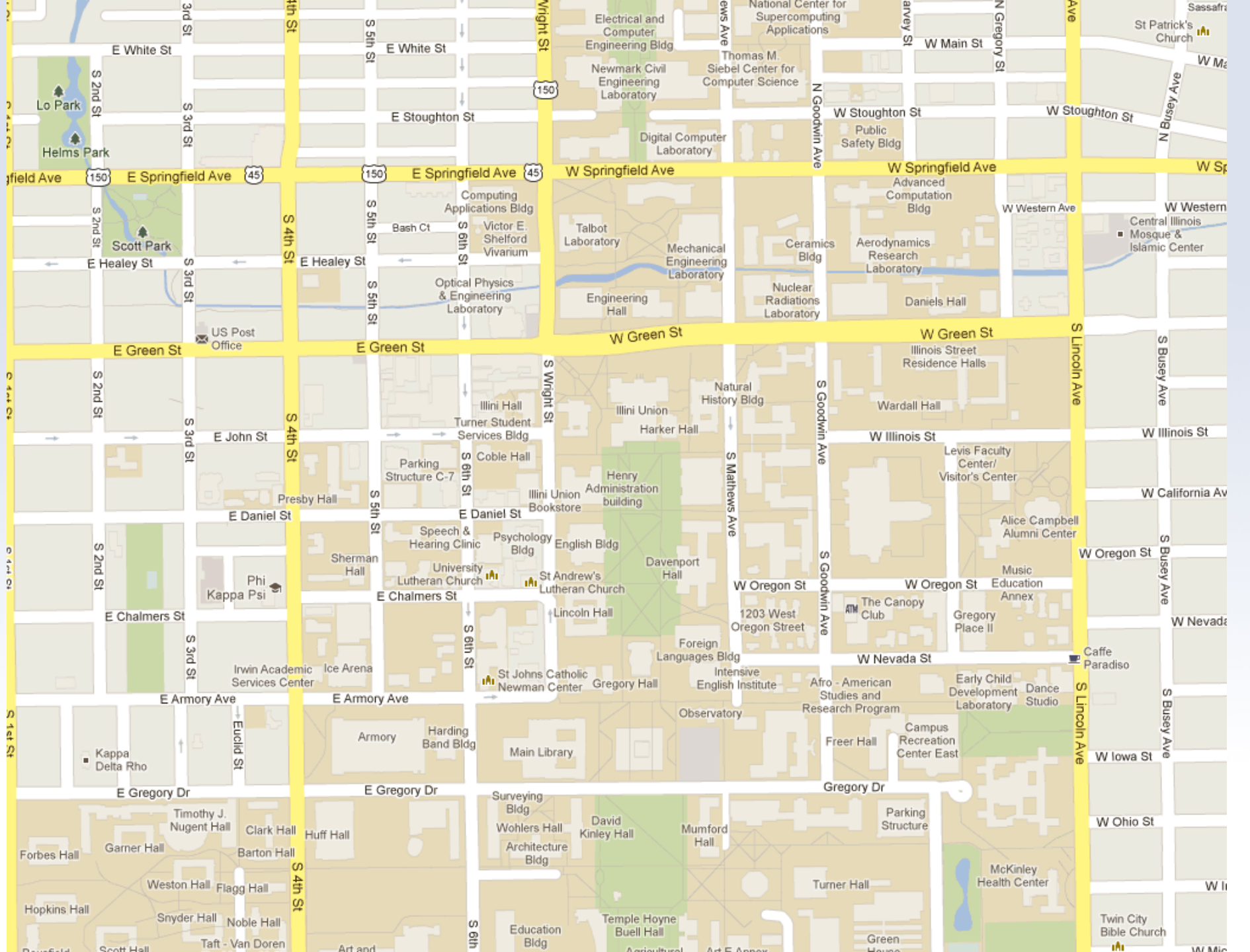


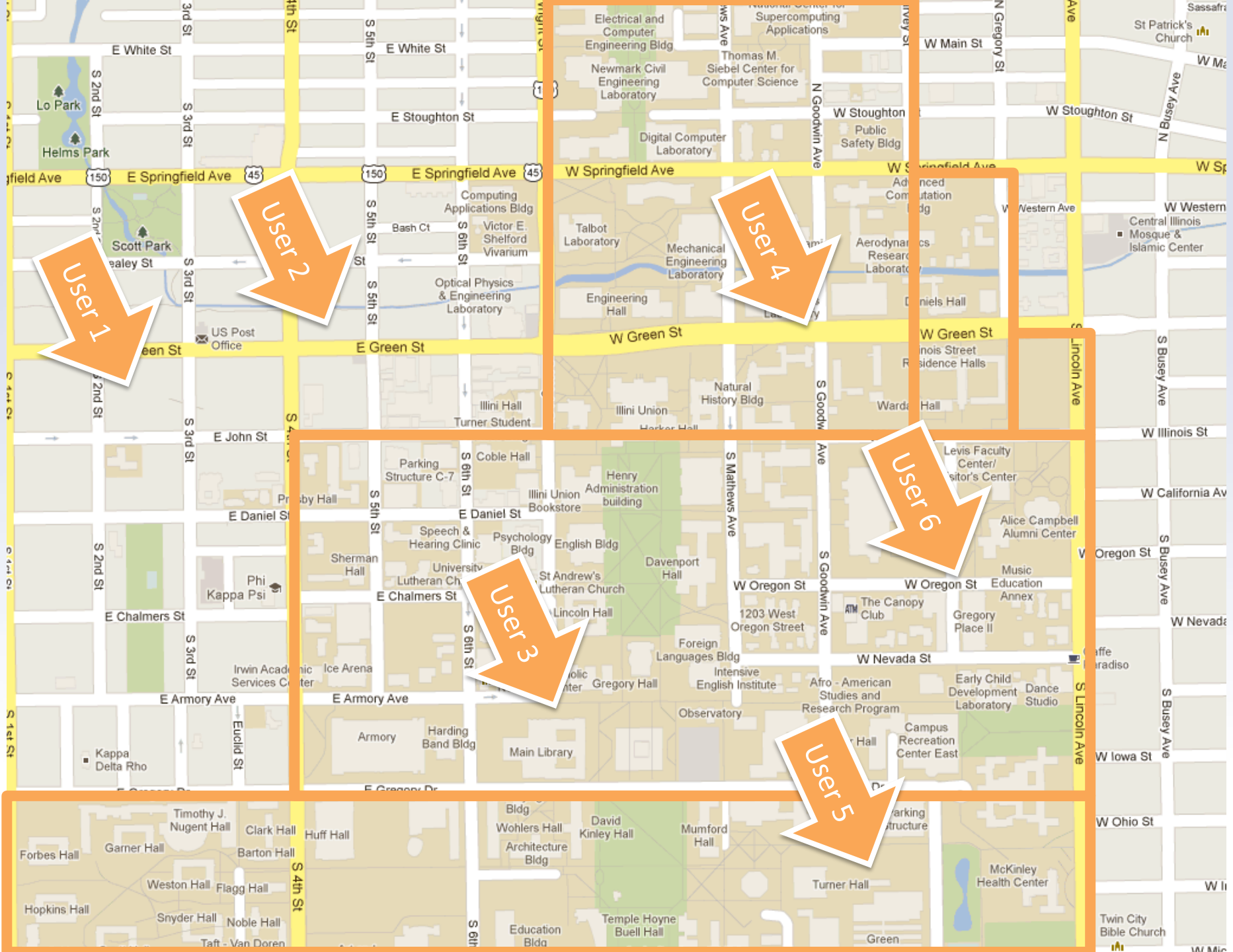
# Multidimensional Data

- We're focus on spatial data
- Queries involve points, lines, and polygons
  - Range queries: find points/shapes in a region
  - Nearest-neighbor: find the nearest point
  - Location queries: find shapes containing a point
- Our examples will be two dimensions:  $x, y$









User 1

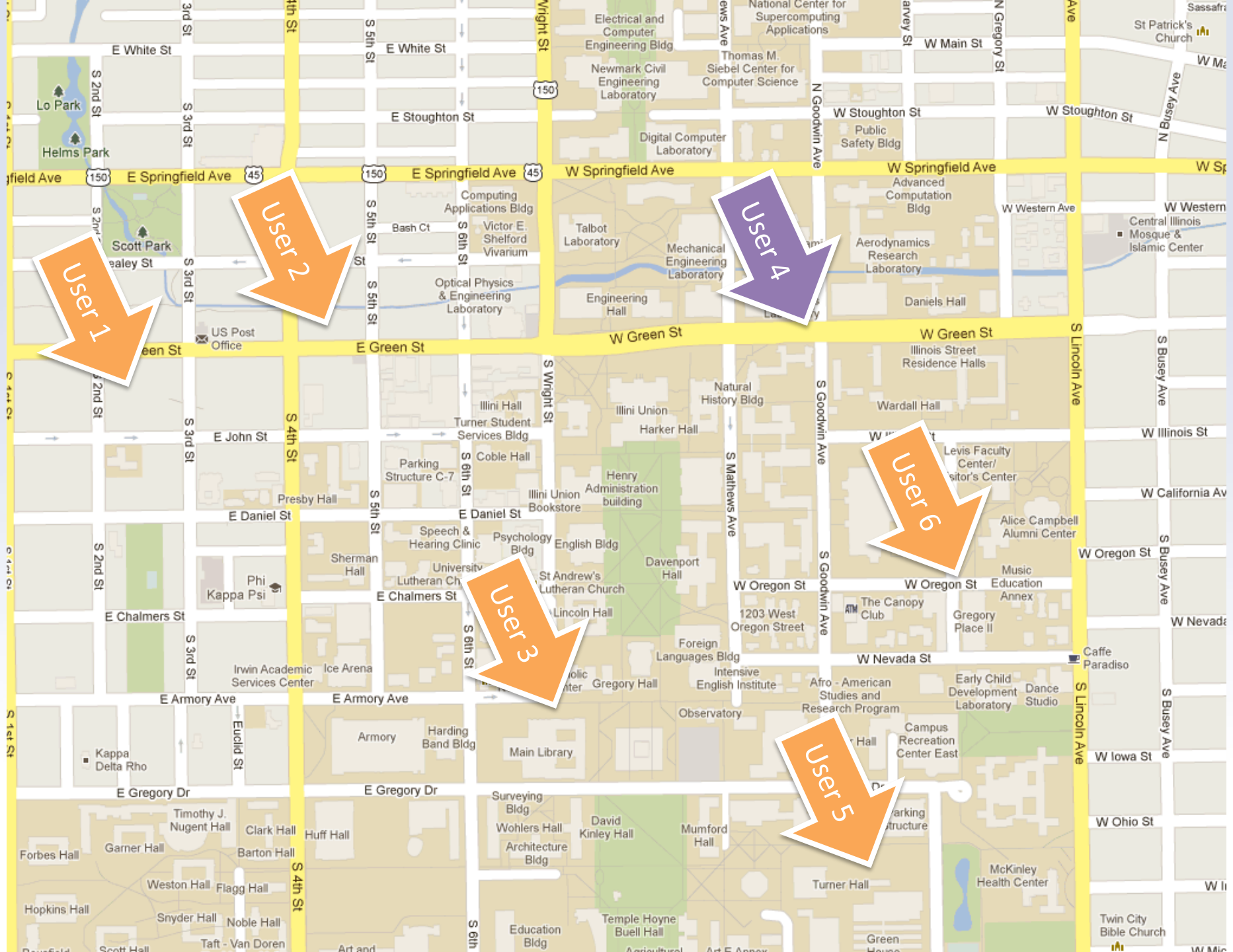
User 2

User 4

User 6

User 3

User 5



User 1

User 2

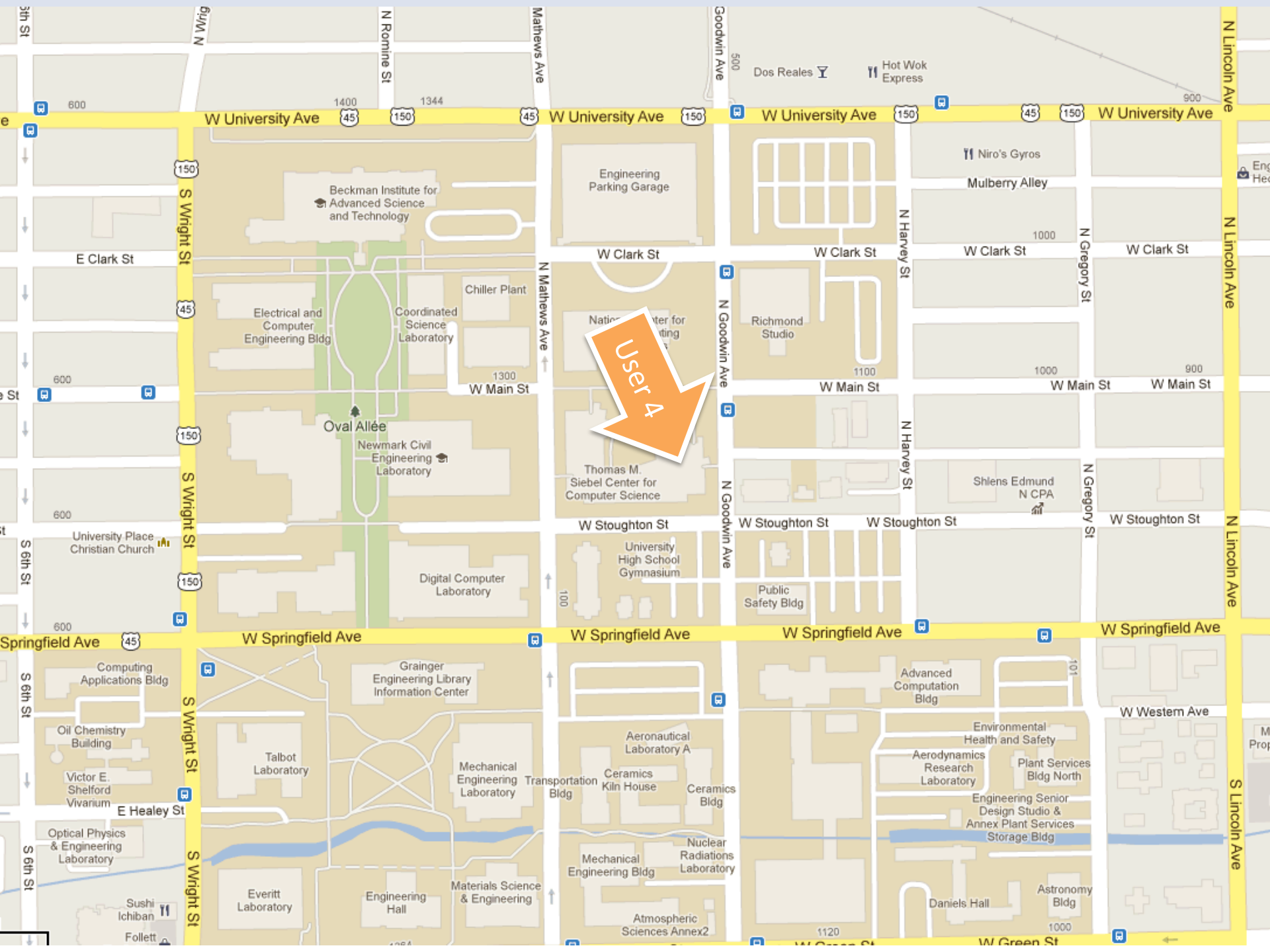
User 3

User 4

User 5

User 6





# MySQL

```
create table Points (  
  name VARCHAR(20) PRIMARY KEY,  
  location Point NOT NULL,  
  description VARCHAR(200),  
  SPATIAL INDEX(location)  
);
```



# MySQL

```
SELECT name, AsText(location)
FROM Points
WHERE Intersects(location, GeomFromText(@bbox) );
```

name	AsText(location)
name_243.0	POINT(6.1908 3.1375)
name_429.0	POINT(6.5194 0.4023)
name_533.0	POINT(7.7479 3.0894)
name_808.0	POINT(4.6818 7.618)



# Multidimensional Indexing

- The book covers a lot of structures for this:
  - Grid Files
  - Partitioned Hash Tables
  - Multiple-Key Indexes
  - Quad trees
- We'll only cover two





# k-d trees

- Useful for range and nearest neighbor queries
- Similar to binary search tree
- Internal nodes alternate dimensions
  - Root splits  $x$
  - Child splits  $y$
  - Grandchild splits  $x$
  - Great grandchild splits  $y...$

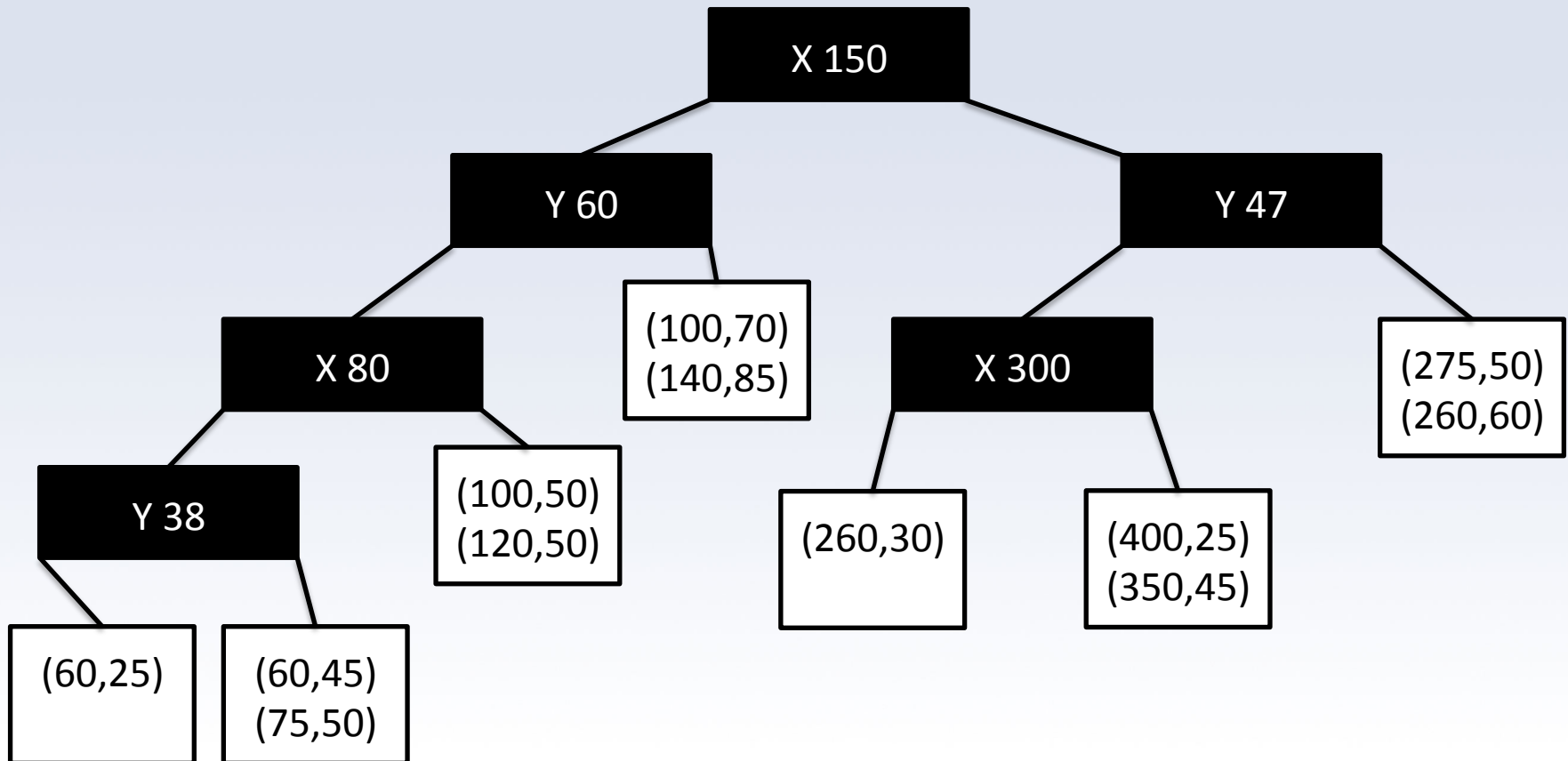


# k-d trees

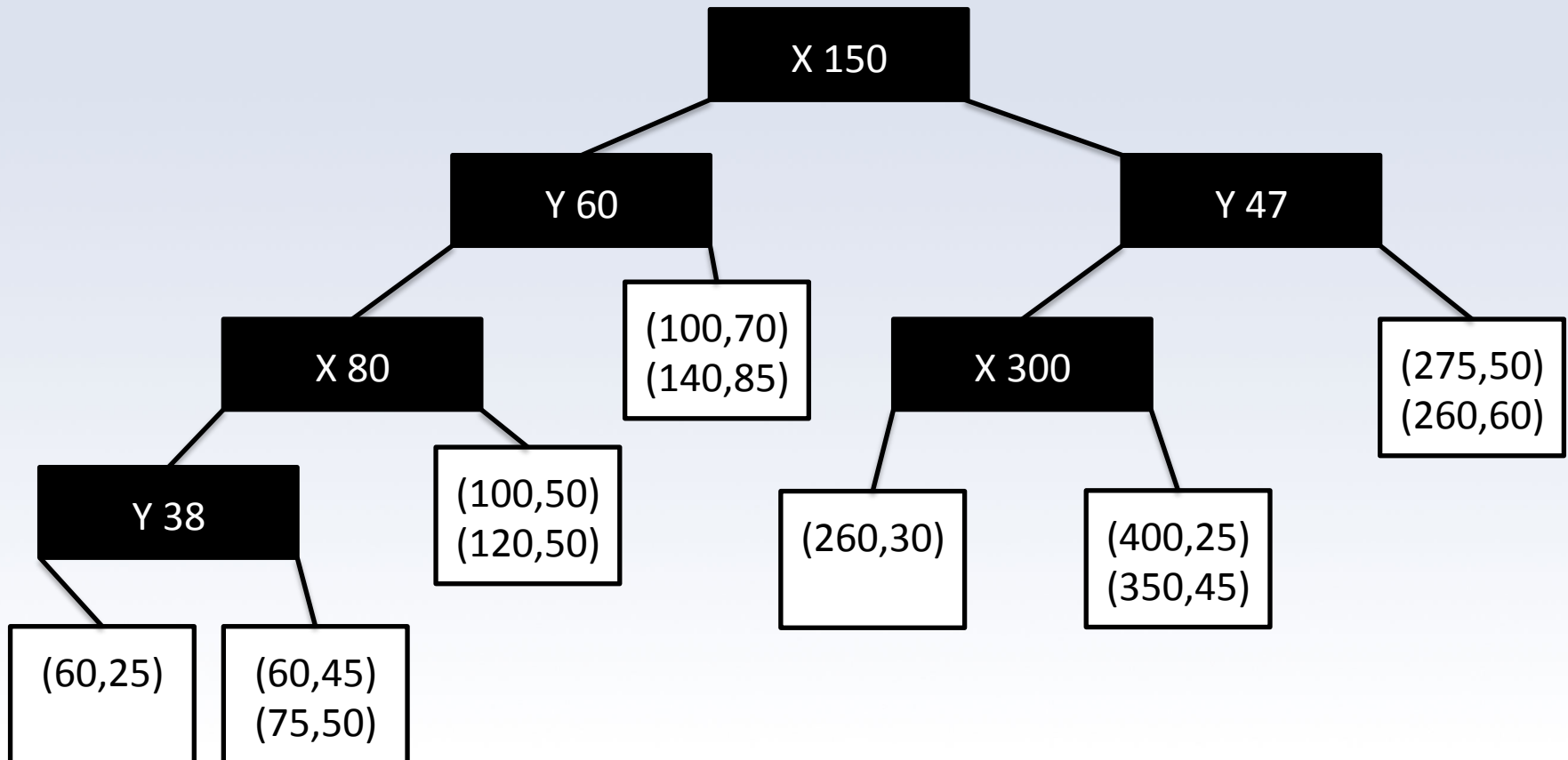
- Each internal node has a dimension  $D$  and value  $V$  associated with it
  - If  $R.D < V$ , point in left subtree
  - If  $R.D \geq V$ , point in right subtree
- Leaves contain blocks
  - not typical for memory-based k-d trees



# Example

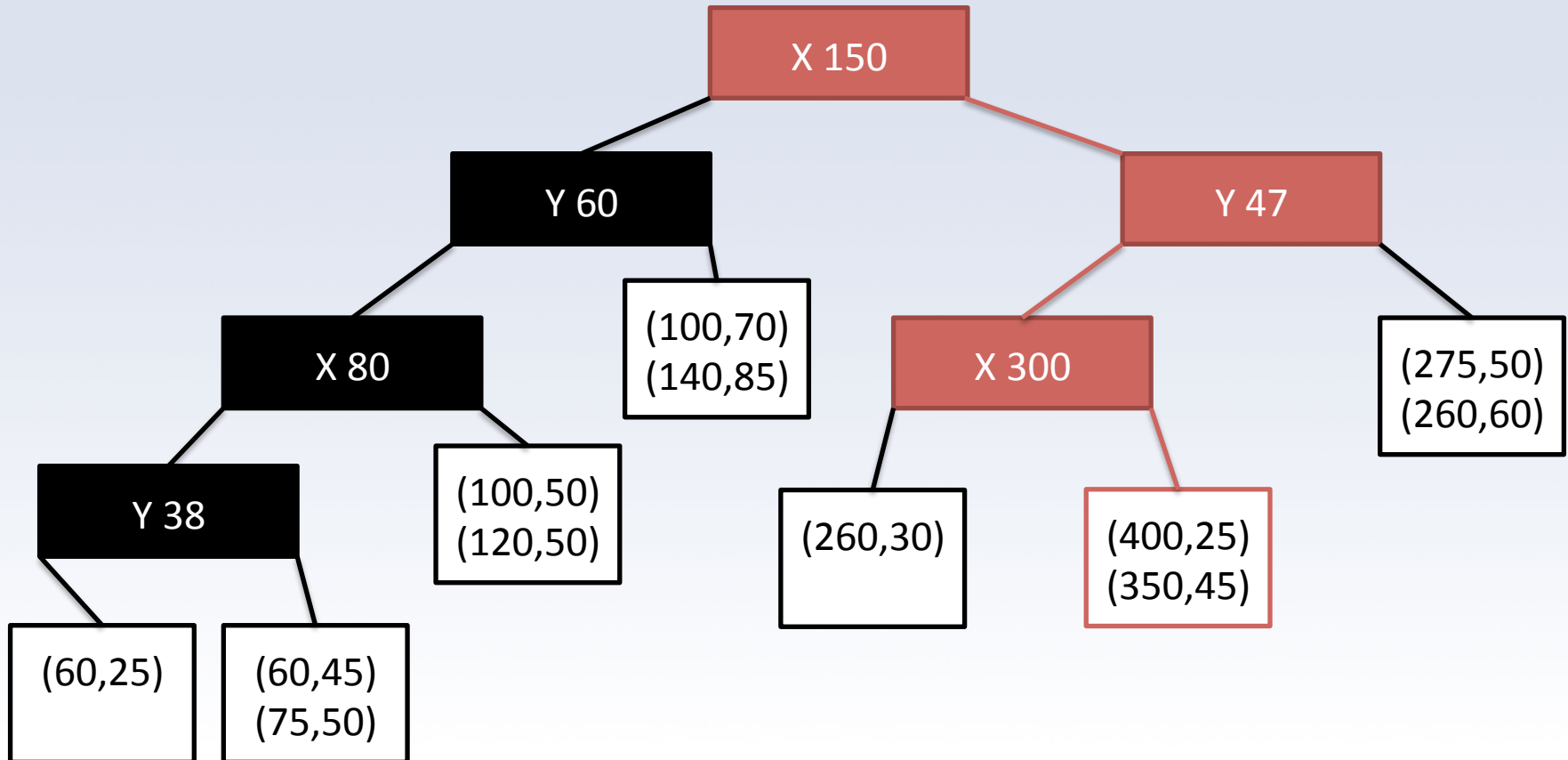


# Insert (500,35)

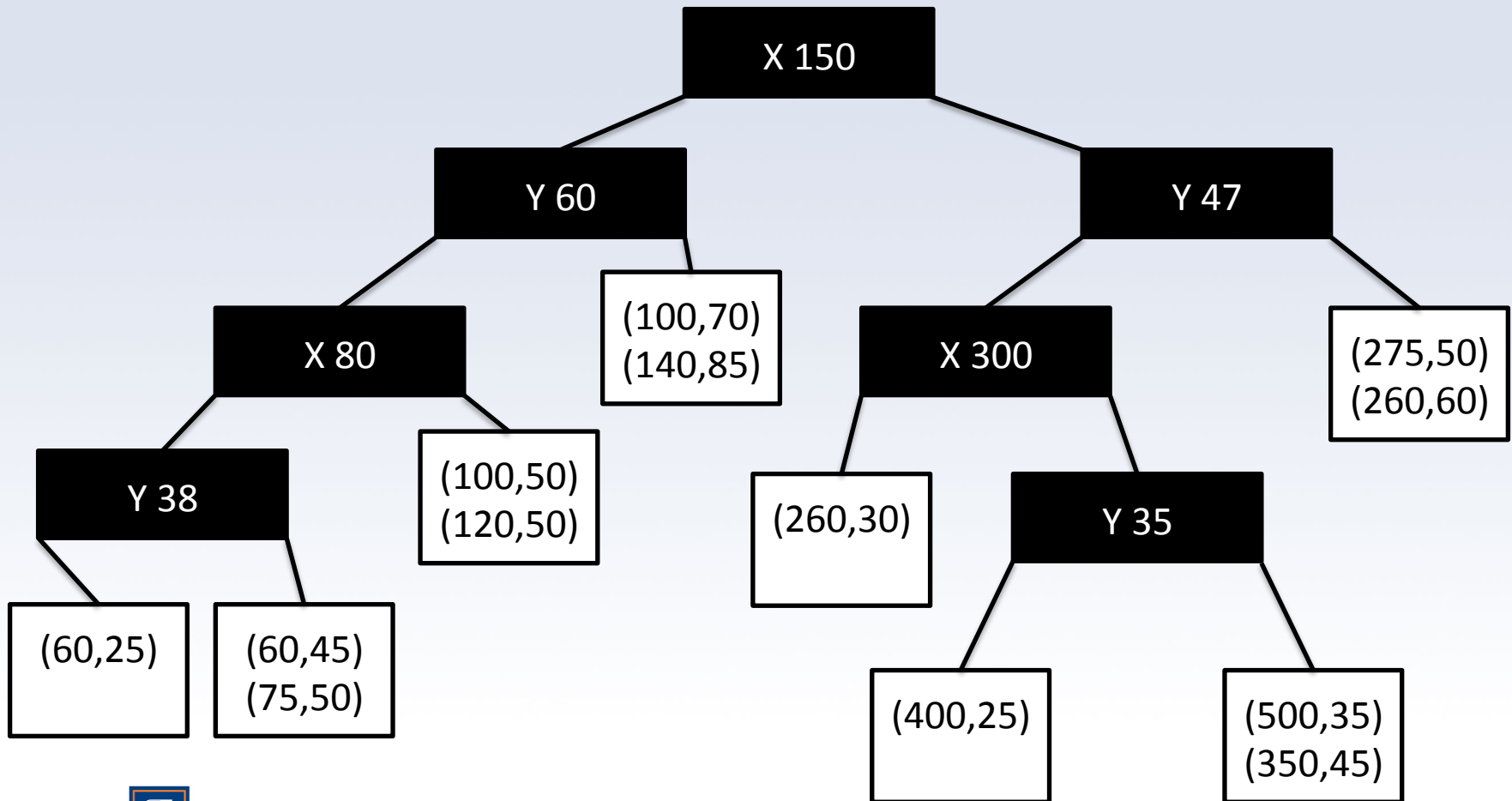




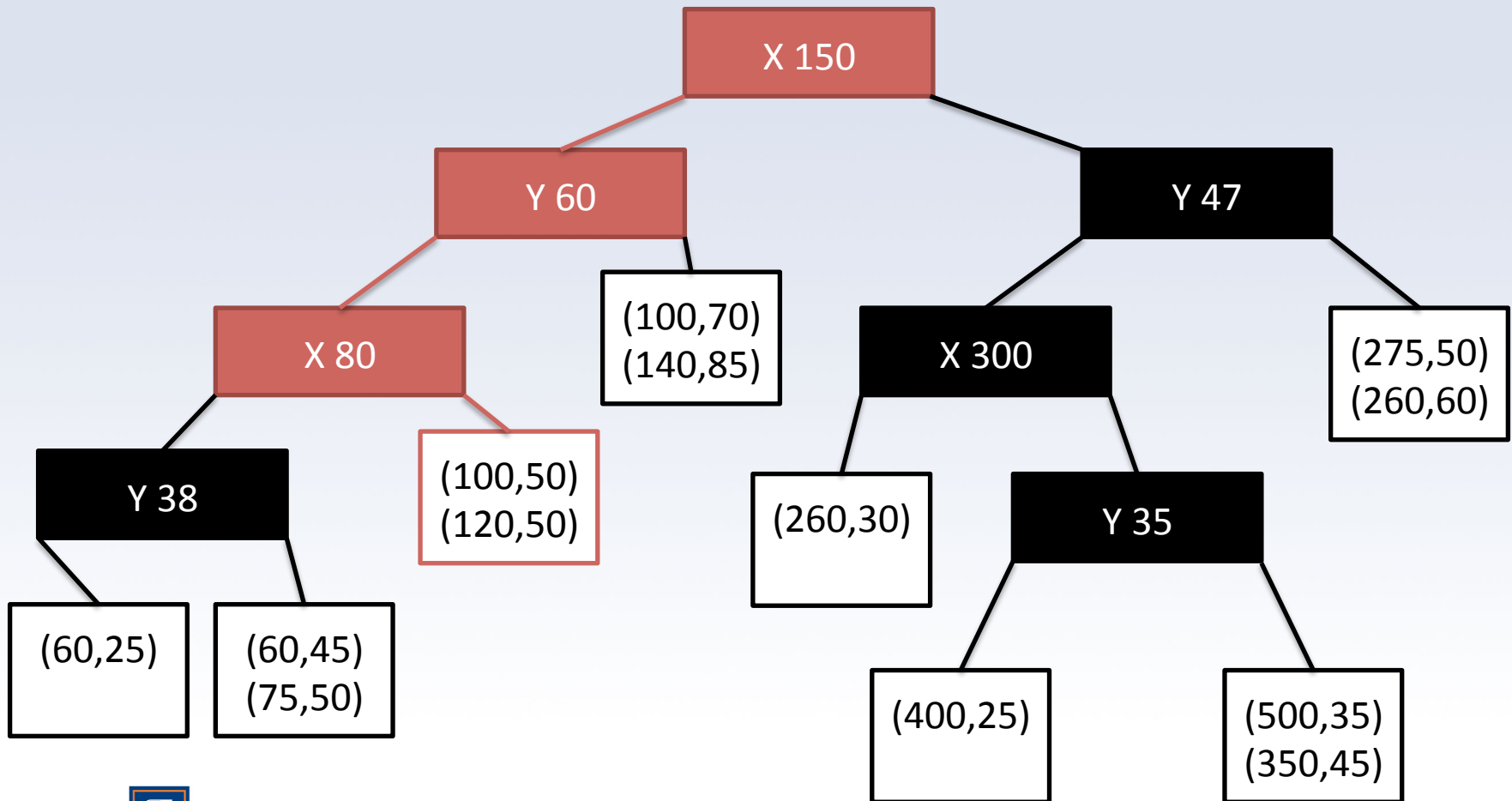
# Insert (500,35)



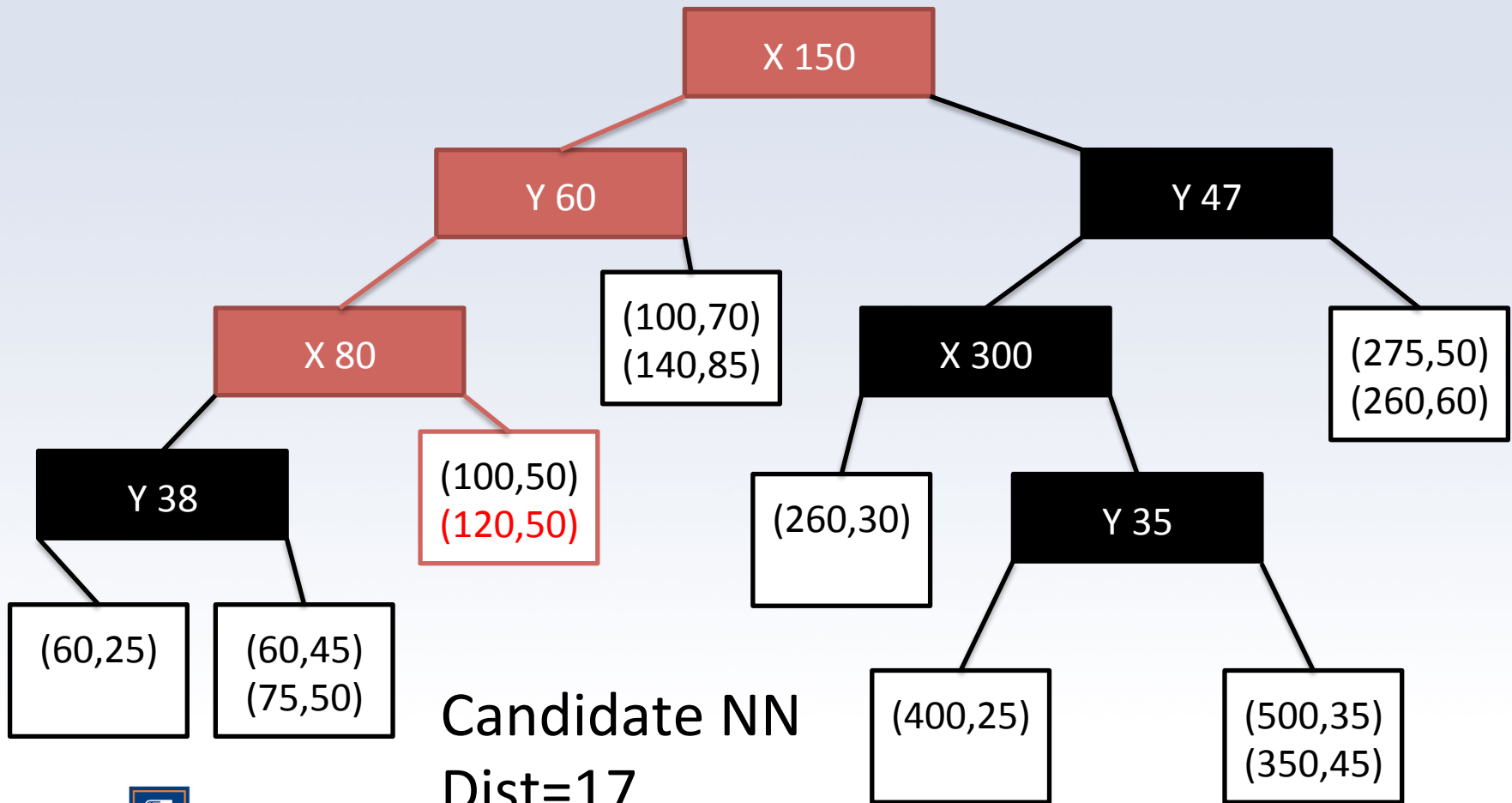
# Insert (500,35)



# Find NN (120,33)



# Find NN (120,33)





# Find NN (120,33)

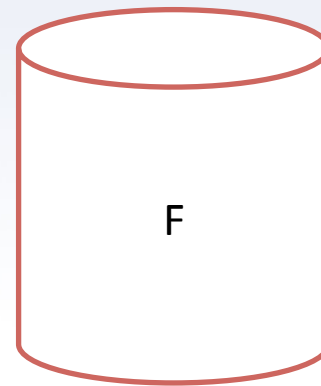
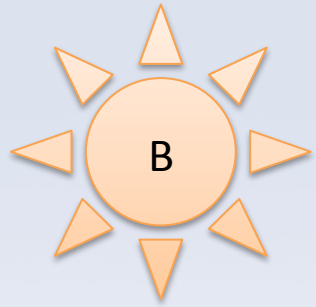
- After finding a candidate, we have to “unwind” back up the search tree to check if any points are closer
- Requires more data at internal nodes
- Not going to cover this

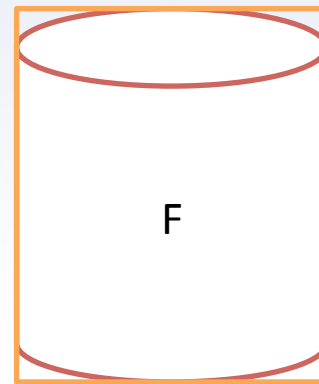
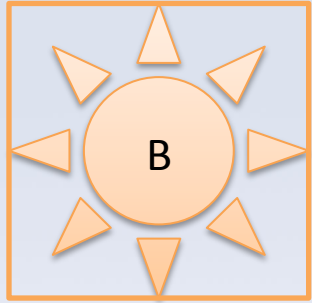
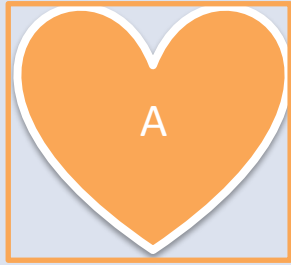


# R-trees

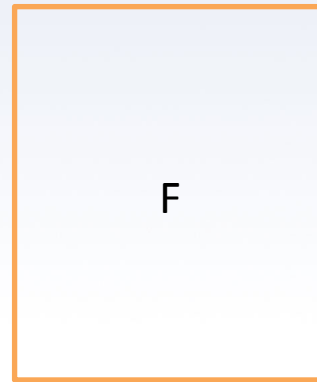
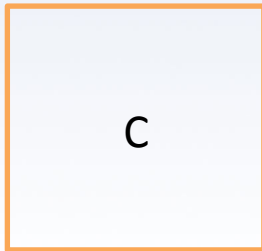
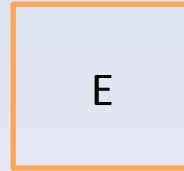
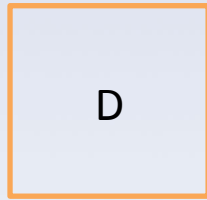
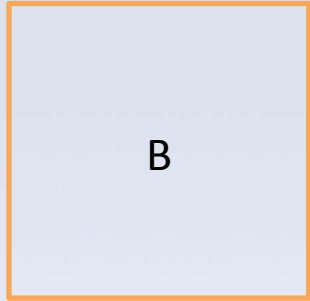
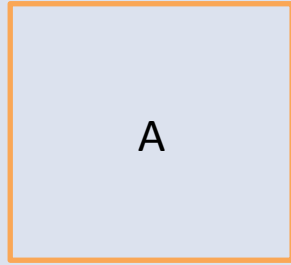
- Groups nearby objects by *minimum bounding rectangles*
  - Smallest rectangle that contains them
- Groups the groups with more bounding rectangles
- Creates a tree of bounding regions

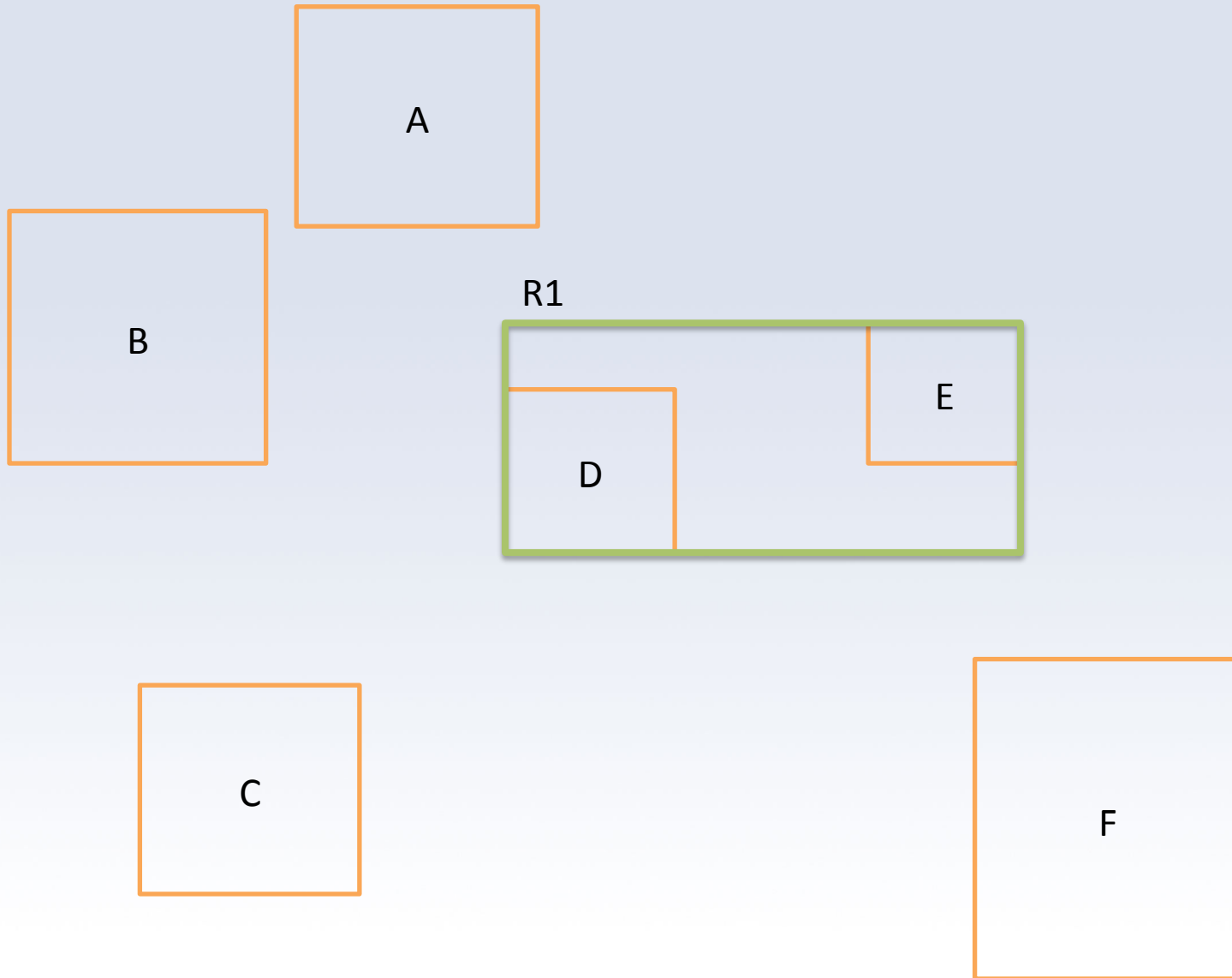


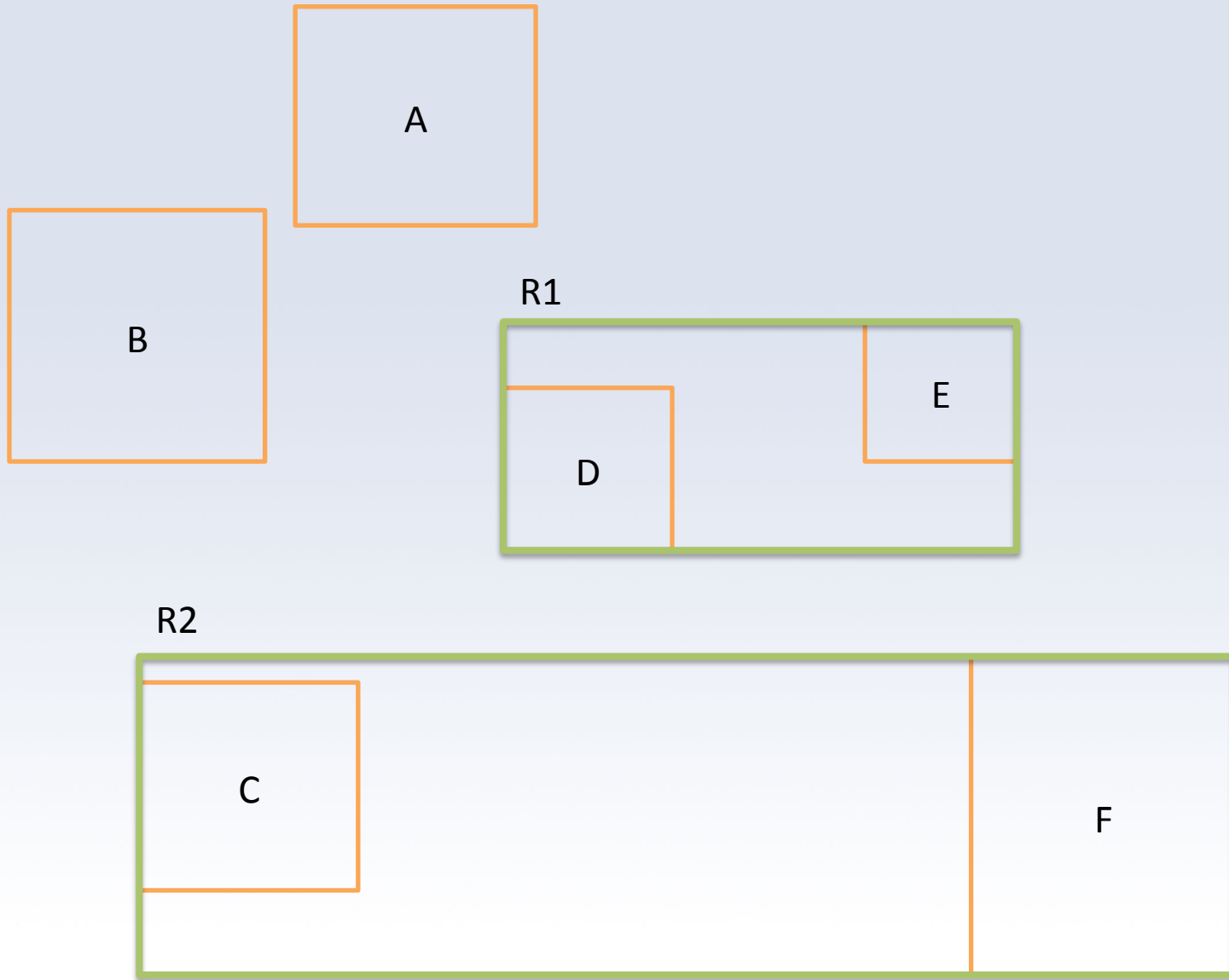


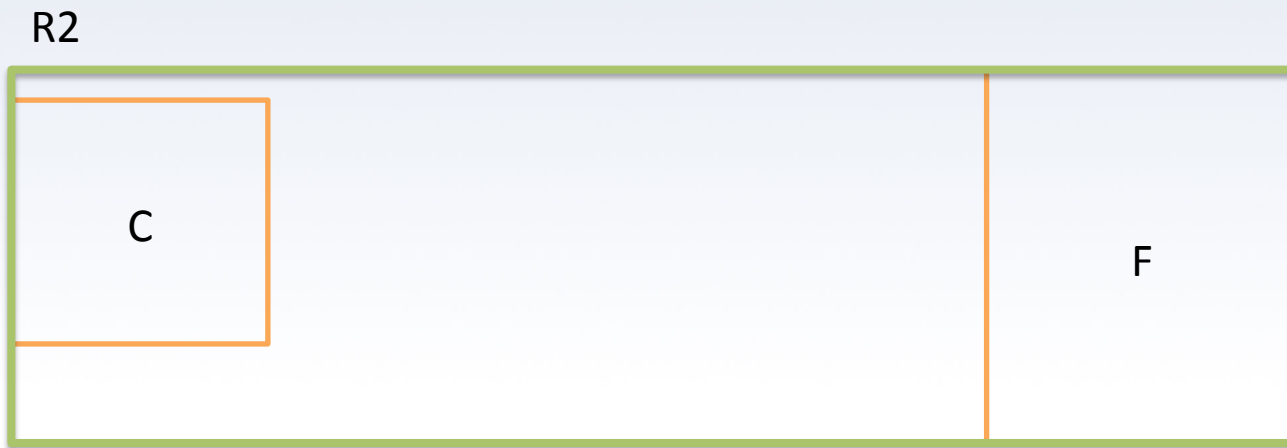
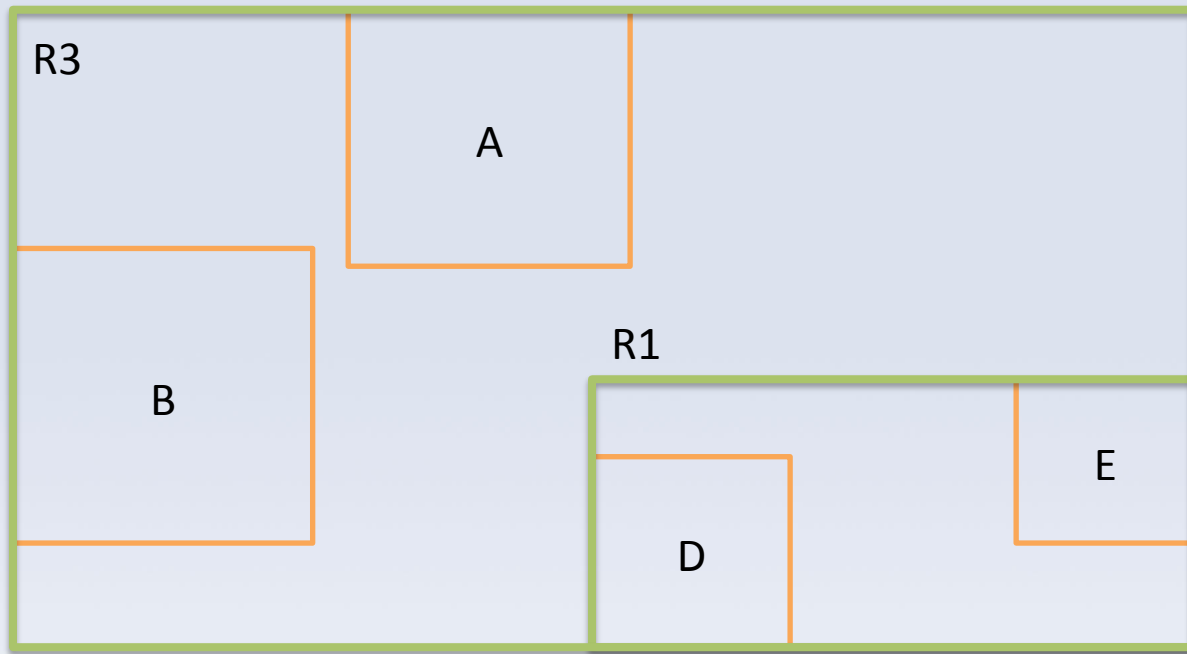












R4

R3

A

B

R1

D

E

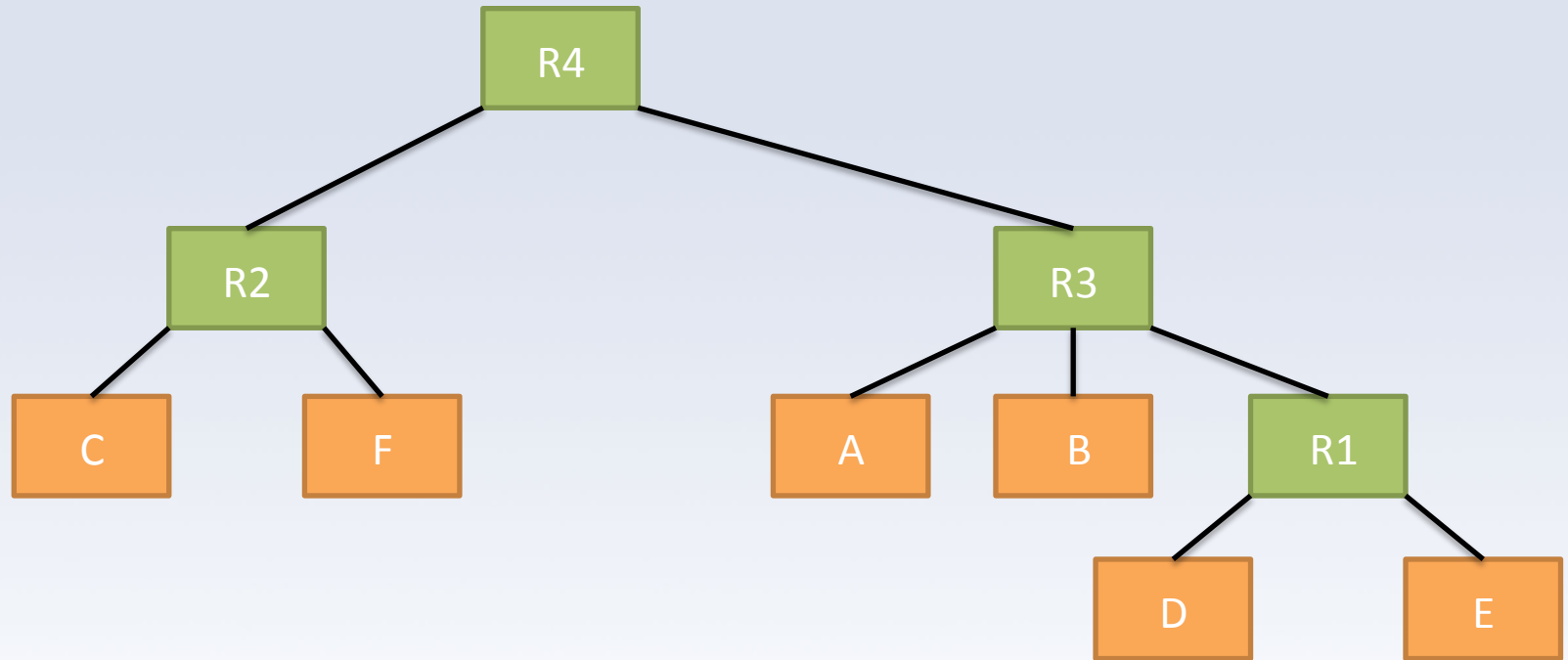
R2

C

F



# R-tree





# Location query

- Recursively check if point is in each rectangle



R4

R3

A

B

R1

D

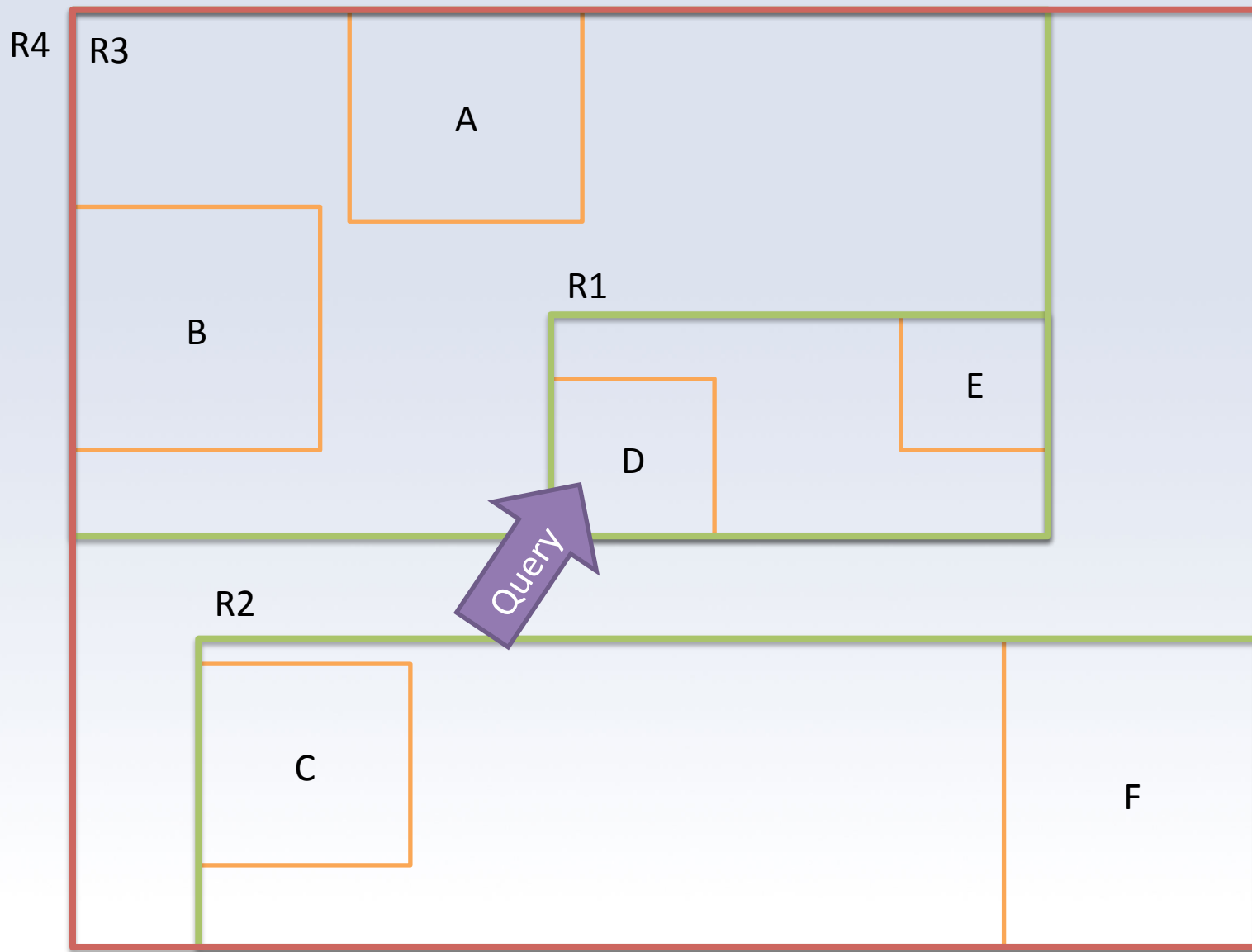
E

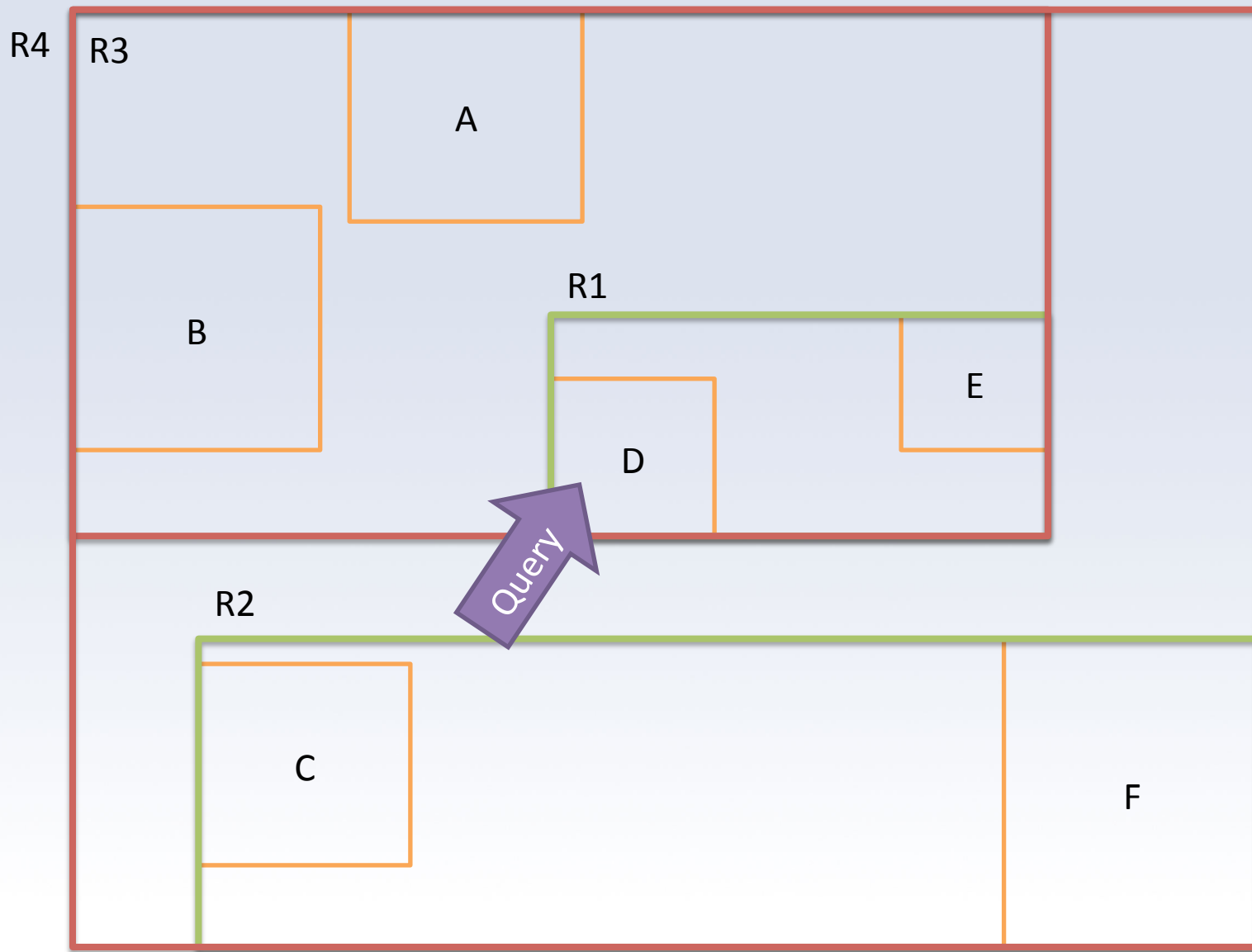
R2

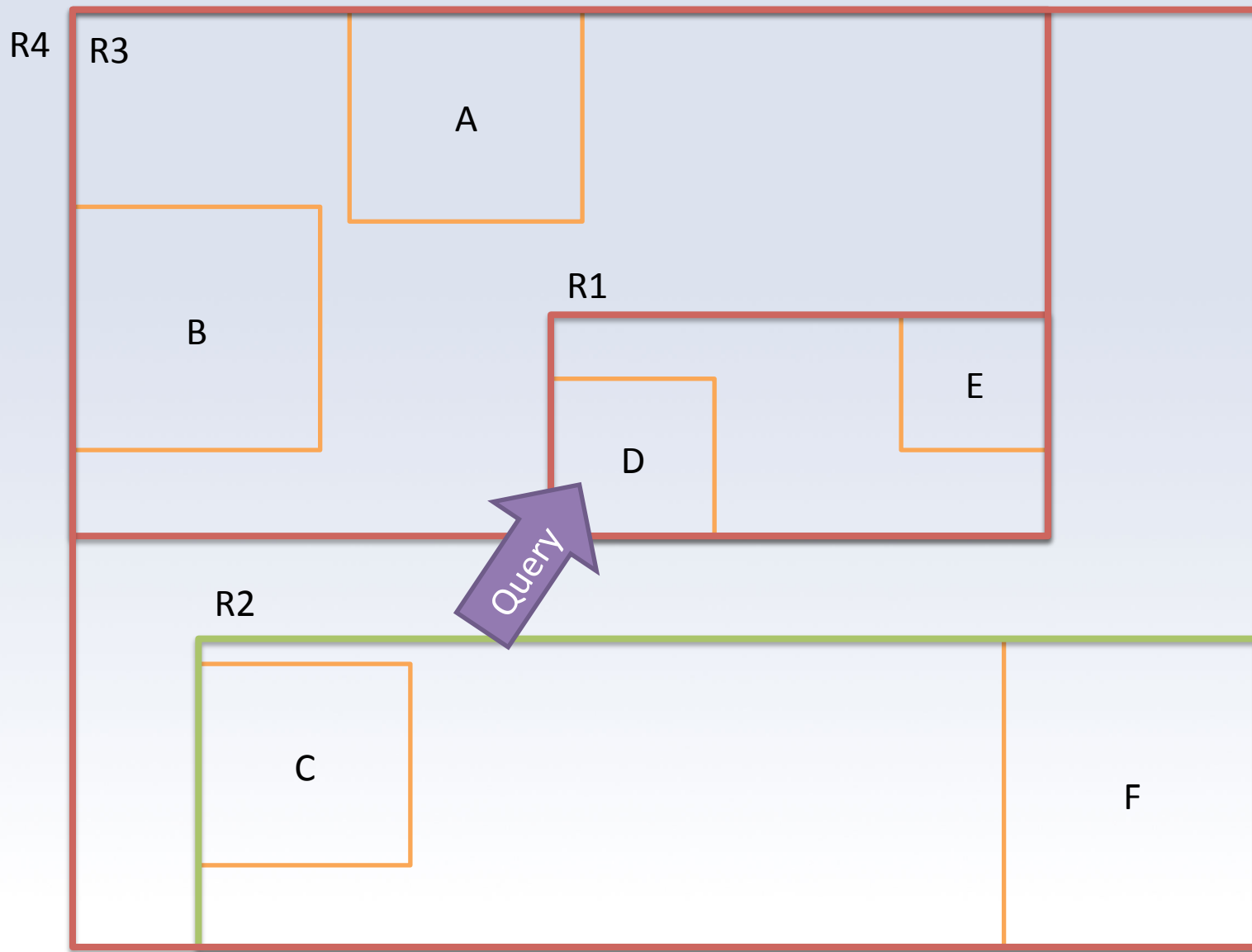
C

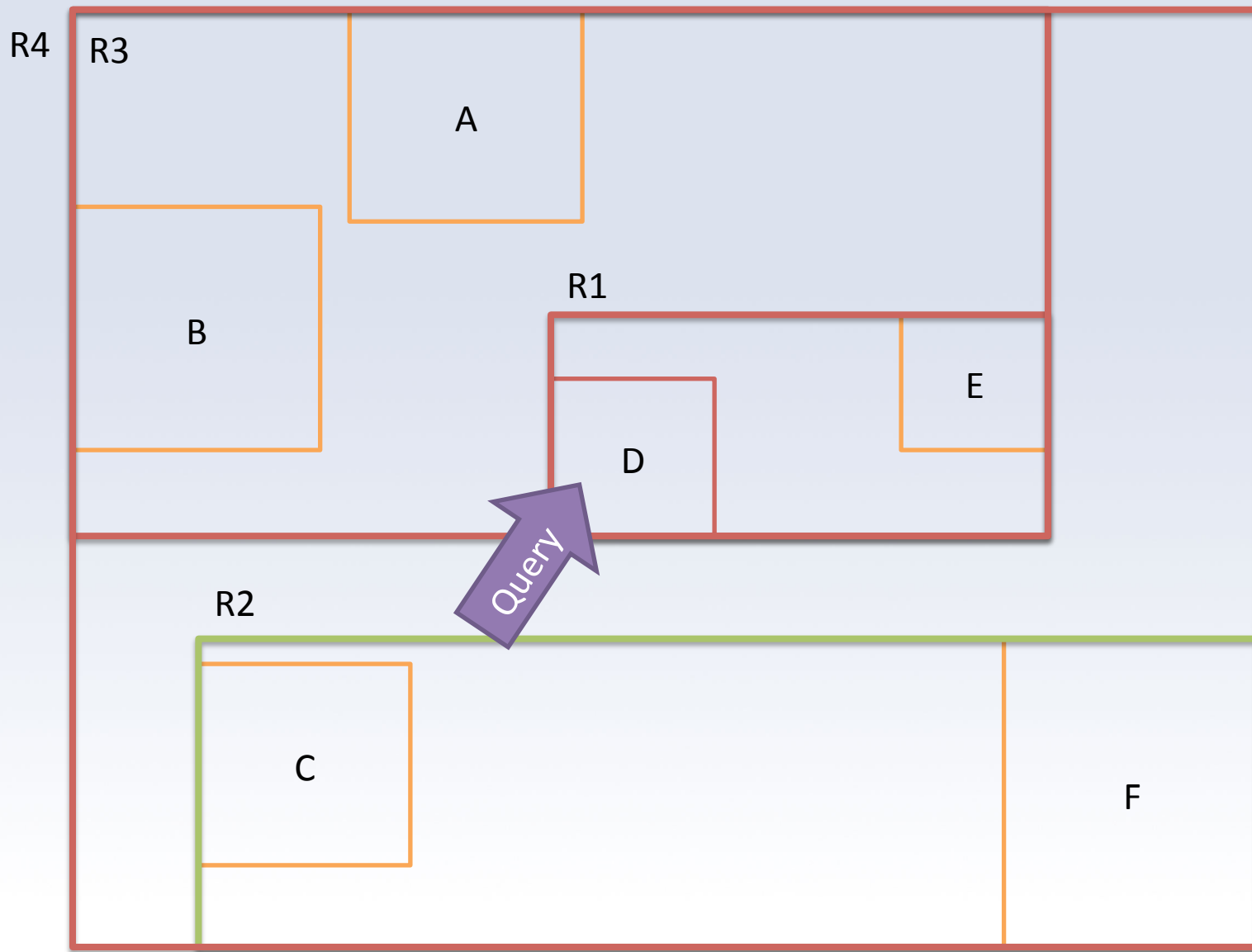
F



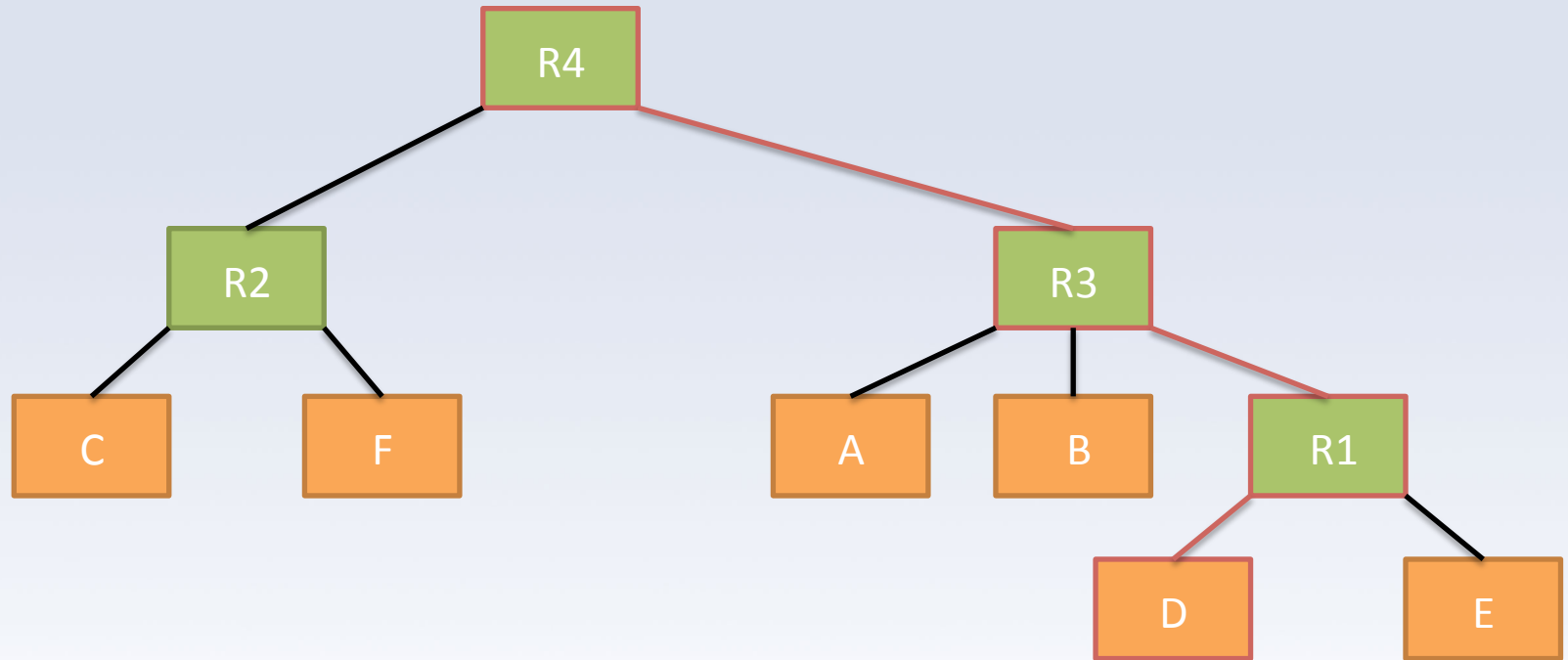








# Location Query





# R-trees

- Actually implemented in MySQL

```
CREATE SPATIAL INDEX spatial_ind  
ON Points(location);
```

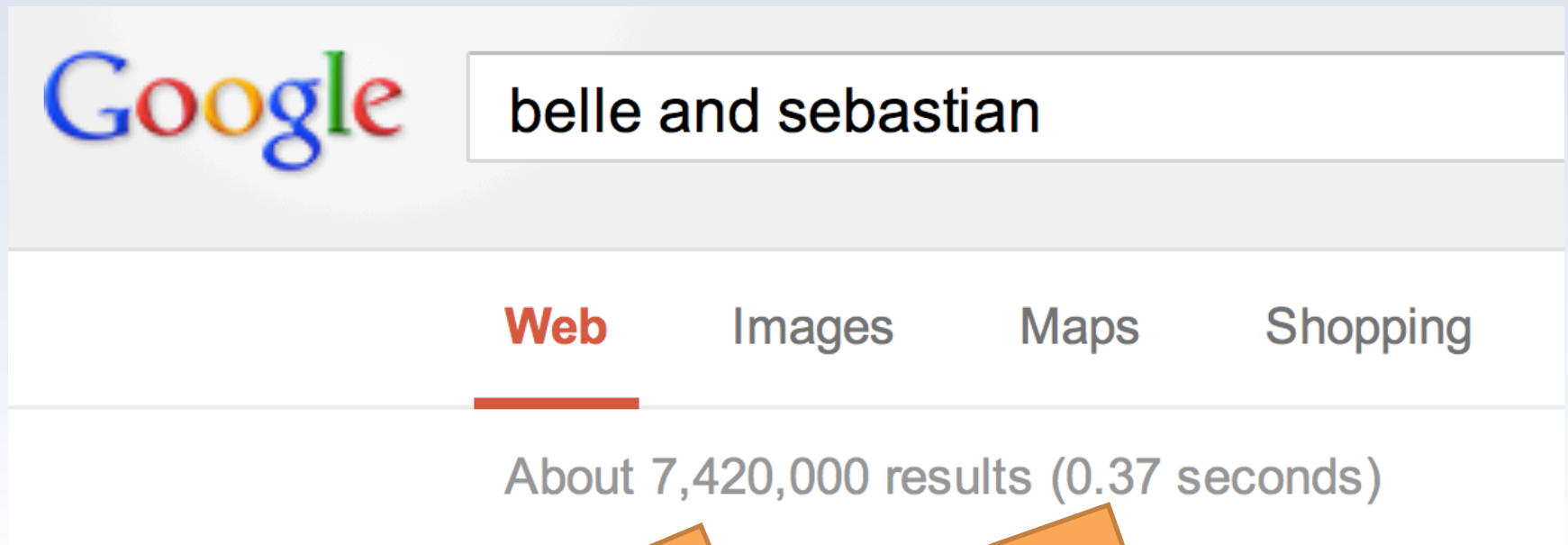


# Document Indexing

- Given a collection of text documents, how can we quickly find all the documents containing a word



# Document search



# Inverted Index

## DOCUMENT 1

Don't you love  
her madly?

## DOCUMENT 2

Say you love me.

## DOCUMENT 3

Come as you are.

## DOCUMENT 4

All you need is  
love.

don't	{1}
you	{1,2,3,4}
love	{1,2,4}
her	{1}
madly	{1}
say	{1}

me	{2}
come	{3}
as	{3}
are	{3}
all	{4}
need	{4}

is	{4}
----	-----

# Query: you love

## DOCUMENT 1

Don't you love  
her madly?

## DOCUMENT 2

Say you love me.

## DOCUMENT 3

Come as you are.

## DOCUMENT 4

All you need is  
love.

don't	{1}
you	{1,2,3,4}
love	{1,2,4}
her	{1}
madly	{1}
say	{1}

me	{2}
come	{3}
as	{3}
are	{3}
all	{4}
need	{4}

is	{4}
----	-----

# Query: you love

## DOCUMENT 1

Don't you love  
her madly?

## DOCUMENT 2

Say you love me.

## DOCUMENT 3

Come as you are.

## DOCUMENT 4

All you need is  
love.

don't	{1}
you	{1,2,3,4}
love	{1,2,4}
her	{1}
madly	{1}
say	{1}

me	{2}
come	{3}
as	{3}
are	{3}
all	{4}
need	{4}

is	{4}
----	-----

$$\{1,2,3,4\} \cap \{1,2,4\} = \{1,2,4\}$$

# Query: you love

## DOCUMENT 1

Don't you love  
her madly?

## DOCUMENT 2

Say you love me.

## DOCUMENT 3

Come as you are.

## DOCUMENT 4

All you need is  
love.

don't	{1}
you	{1,2,3,4}
love	{1,2,4}
her	{1}
madly	{1}
say	{1}

me	{2}
come	{3}
as	{3}
are	{3}
all	{4}
need	{4}

is	{4}
----	-----

$$\{1,2,3,4\} \cap \{1,2,4\} = \{1,2,4\}$$



# Inverted Indexing

- Real inverted indexes also keep track of the location of the word in each page
- This allows exact phrases to be looked up



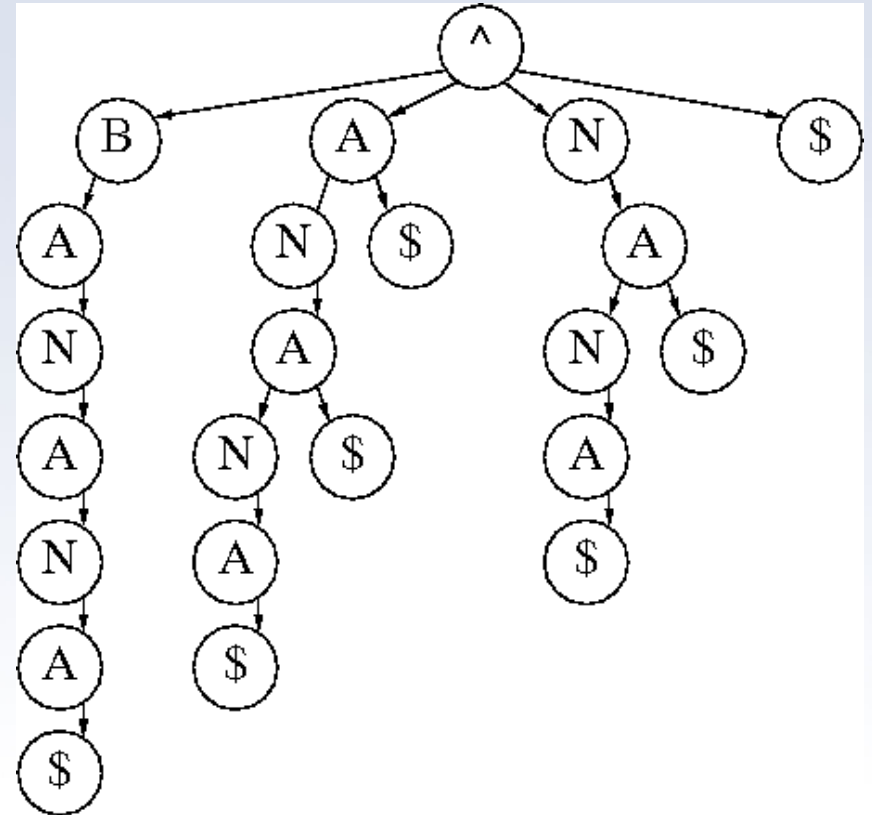
# String indexing

- How can we index an extremely long string for search?



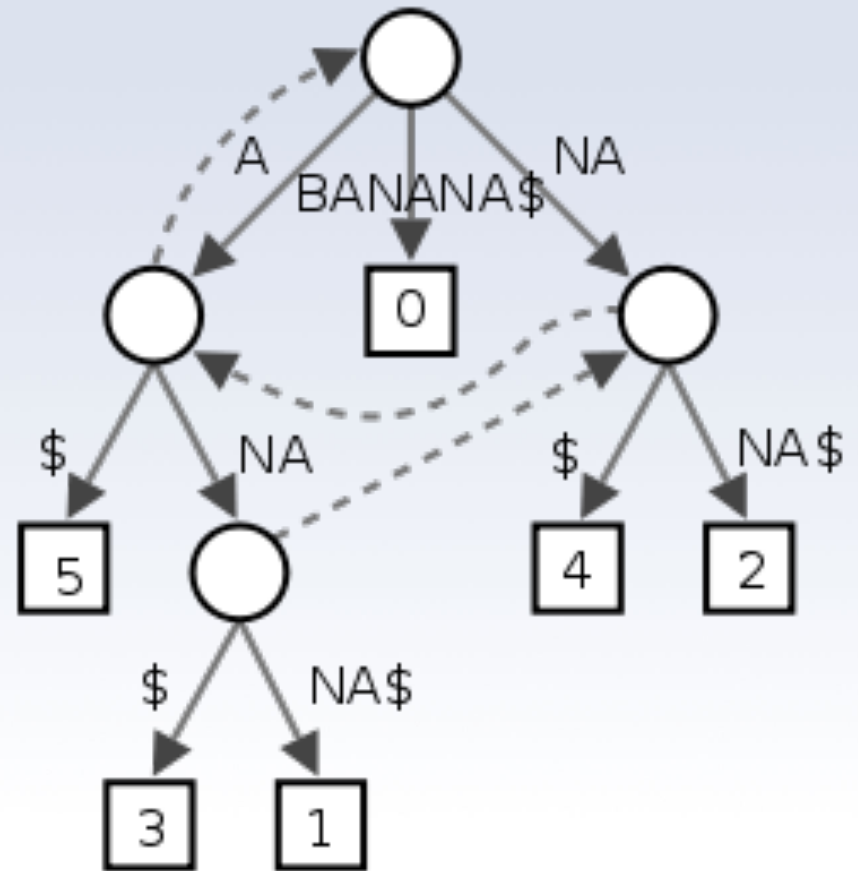
# Suffix Trie

- A tree where each path from root to leaf corresponds to a suffix
- $O(n)$  to search
- $O(n^2)$  memory



# Suffix Tree

- Compressed version of the suffix trie
- Can be searched in  $O(n)$
- Takes  $O(n)$  memory
- Can be constructed in  $O(n)$  !?!?!?
- Ukkonen's algorithm



# Suffix Array

- SA:
  - Given a text  $T$ , add a string terminator character '\$'
  - Compute all  $n$  rotations (cyclic rotations)
  - Sort the rotations
  - Suffix array is the index of the rotation
- Can be computed in  $O(n)$  from a suffix tree



# Suffix Array

Suffix	i
banana\$	1
anana\$	2
nana\$	3
ana\$	4
na\$	5
a\$	6
\$	7



# Suffix Array

Suffix	i
\$	7
a\$	6
ana\$	4
anana\$	2
banana\$	1
na\$	5
nana\$	3



# Burrows-Wheeler Transform

- BWT:
  - Given a text  $T$ , add a string terminator character '\$'
  - Compute all  $n$  rotations (cyclic rotations)
  - Sort the rotations
  - BWT is last column
- Can be computed in  $O(n)$  from a suffix array





A	B	R	A	C	A	D	A	B	R	A
---	---	---	---	---	---	---	---	---	---	---



A	B	R	A	C	A	D	A	B	R	A	\$
---	---	---	---	---	---	---	---	---	---	---	----



0	A	B	R	A	C	A	D	A	B	R	A	\$
---	---	---	---	---	---	---	---	---	---	---	---	----



0	A	B	R	A	C	A	D	A	B	R	A	\$
1	B	R	A	C	A	D	A	B	R	A	\$	A



0	A	B	R	A	C	A	D	A	B	R	A	\$
1	B	R	A	C	A	D	A	B	R	A	\$	A
2	R	A	C	A	D	A	B	R	A	\$	A	B



0	A	B	R	A	C	A	D	A	B	R	A	\$
1	B	R	A	C	A	D	A	B	R	A	\$	A
2	R	A	C	A	D	A	B	R	A	\$	A	B
3	A	C	A	D	A	B	R	A	\$	A	B	R



0	A	B	R	A	C	A	D	A	B	R	A	\$
1	B	R	A	C	A	D	A	B	R	A	\$	A
2	R	A	C	A	D	A	B	R	A	\$	A	B
3	A	C	A	D	A	B	R	A	\$	A	B	R
4	C	A	D	A	B	R	A	\$	A	B	R	A



0	A	B	R	A	C	A	D	A	B	R	A	\$
1	B	R	A	C	A	D	A	B	R	A	\$	A
2	R	A	C	A	D	A	B	R	A	\$	A	B
3	A	C	A	D	A	B	R	A	\$	A	B	R
4	C	A	D	A	B	R	A	\$	A	B	R	A
5	A	D	A	B	R	A	\$	A	B	R	A	C
6	D	A	B	R	A	\$	A	B	R	A	C	A
7	A	B	R	A	\$	A	B	R	A	C	A	D
8	B	R	A	\$	A	B	R	A	C	A	D	A
9	R	A	\$	A	B	R	A	C	A	D	A	B
10	A	\$	A	B	R	A	C	A	D	A	B	R
11	\$	A	B	R	A	C	A	D	A	B	R	A





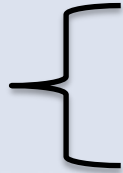
11	\$	A	B	R	A	C	A	D	A	B	R	A
10	A	\$	A	B	R	A	C	A	D	A	B	R
7	A	B	R	A	\$	A	B	R	A	C	A	D
0	A	B	R	A	C	A	D	A	B	R	A	\$
3	A	C	A	D	A	B	R	A	\$	A	B	R
5	A	D	A	B	R	A	\$	A	B	R	A	C
8	B	R	A	\$	A	B	R	A	C	A	D	A
1	B	R	A	C	A	D	A	B	R	A	\$	A
4	C	A	D	A	B	R	A	\$	A	B	R	A
6	D	A	B	R	A	\$	A	B	R	A	C	A
9	R	A	\$	A	B	R	A	C	A	D	A	B
2	R	A	C	A	D	A	B	R	A	\$	A	B

Suffix  
Array

BWT  
Transform



ABRA



11	\$	A	B	R	A	C	A	D	A	B	R	A
10	A	\$	A	B	R	A	C	A	D	A	B	R
7	A	B	R	A	\$	A	B	R	A	C	A	D
0	A	B	R	A	C	A	D	A	B	R	A	\$
3	A	C	A	D	A	B	R	A	\$	A	B	R
5	A	D	A	B	R	A	\$	A	B	R	A	C
8	B	R	A	\$	A	B	R	A	C	A	D	A
1	B	R	A	C	A	D	A	B	R	A	\$	A
4	C	A	D	A	B	R	A	\$	A	B	R	A
6	D	A	B	R	A	\$	A	B	R	A	C	A
9	R	A	\$	A	B	R	A	C	A	D	A	B
2	R	A	C	A	D	A	B	R	A	\$	A	B



Suffix  
Array

11	\$	A
10	A	R
7	A	D
0	A	\$
3	A	R
5	A	C
8	B	A
1	B	A
4	C	A
6	D	A
9	R	B
2	R	B

BWT  
Transform



11	\$	A	B	R	A	C	A	D	A	B	R	A
10	A	\$	A	B	R	A	C	A	D	A	B	R
7	A	B	R	A	\$	A	B	R	A	C	A	D
0	A	B	R	A	C	A	D	A	B	R	A	\$
3	A	C	A	D	A	B	R	A	\$	A	B	R
5	A	D	A	B	R	A	\$	A	B	R	A	C
8	B	R	A	\$	A	B	R	A	C	A	D	A
1	B	R	A	C	A	D	A	B	R	A	\$	A
4	C	A	D	A	B	R	A	\$	A	B	R	A
6	D	A	B	R	A	\$	A	B	R	A	C	A
9	R	A	\$	A	B	R	A	C	A	D	A	B
2	R	A	C	A	D	A	B	R	A	\$	A	B

Suffix  
Array

BWT  
Transform



	\$	A	B	R	A	C	A	D	A	B	R	A	1
1	A	\$	A	B	R	A	C	A	D	A	B	R	
2	A	B	R	A	\$	A	B	R	A	C	A	D	
3	A	B	R	A	C	A	D	A	B	R	A	\$	
4	A	C	A	D	A	B	R	A	\$	A	B	R	
5	A	D	A	B	R	A	\$	A	B	R	A	C	
	B	R	A	\$	A	B	R	A	C	A	D	A	2
	B	R	A	C	A	D	A	B	R	A	\$	A	3
	C	A	D	A	B	R	A	\$	A	B	R	A	4
	D	A	B	R	A	\$	A	B	R	A	C	A	5
	R	A	\$	A	B	R	A	C	A	D	A	B	
	R	A	C	A	D	A	B	R	A	\$	A	B	



Search  
Interval:  
(0,11)

Search  
String:  
ABRA

11	\$	A
10	A	R
7	A	D
0	A	\$
3	A	R
5	A	C
8	B	A
1	B	A
4	C	A
6	D	A
9	R	B
2	R	B



Search  
Interval:  
(0,11)

Search  
String:  
ABR**A**

11	\$	A
10	<b>A</b>	R
7	<b>A</b>	D
0	<b>A</b>	\$
3	<b>A</b>	R
5	<b>A</b>	C
8	B	A
1	B	A
4	C	A
6	D	A
9	R	B
2	R	B



Search  
Interval:  
(0,11)

Search  
String:  
ABRA

11	\$	A
10	A	R
7	A	D
0	A	\$
3	A	R
5	A	C
8	B	A
1	B	A
4	C	A
6	D	A
9	R	B
2	R	B





Search  
Interval:  
(0,11)

Search  
String:  
ABRA

11	\$	A	1 <sup>st</sup> R
10	A	R	
7	A	D	
0	A	\$	2 <sup>nd</sup> R
3	A	R	
5	A	C	
8	B	A	
1	B	A	
4	C	A	
6	D	A	
9	R	B	
2	R	B	



Search  
Interval:  
(10,11)

Search  
String:  
ABRA

11	\$	A	1 <sup>st</sup> R
10	A	R	
7	A	D	
0	A	\$	2 <sup>nd</sup> R
3	A	R	
5	A	C	
8	B	A	
1	B	A	
4	C	A	
6	D	A	
9	R	B	
2	R	B	



Search  
Interval:  
(10,11)

Search  
String:  
ABRA

11	\$	A	
10	A	R	
7	A	D	
0	A	\$	
3	A	R	
5	A	C	
8	B	A	
1	B	A	
4	C	A	
6	D	A	
9	R	B	1 <sup>st</sup> B
2	R	B	2 <sup>nd</sup> B



Search  
Interval:  
(6,7)

Search  
String:  
ABRA

11	\$	A
10	A	R
7	A	D
0	A	\$
3	A	R
5	A	C
8	B	A
1	B	A
4	C	A
6	D	A
9	R	B
2	R	B



Search  
Interval:  
(6,7)

Search  
String:  
ABRA

11	\$	A
10	A	R
7	A	D
0	A	\$
3	A	R
5	A	C
8	B	A
1	B	A
4	C	A
6	D	A
9	R	B
2	R	B

2<sup>nd</sup> A

3<sup>rd</sup> A



Search  
Interval:  
(2,3)

Search  
String:  
ABRA

Two occurrences found  
At indexes 0 and 7

0 1 2 3 4 5 6 7  
A B R A C A D A B R A

11	\$	A
10	A	R
7	A	D
0	A	\$
3	A	R
5	A	C
8	B	A
1	B	A
4	C	A
6	D	A
9	R	B
2	R	B



# Improvements

- This algorithm can be improved:
  - We can make lookup tables for the “rank” of each symbol in BWT
  - We can make lookup tables for the location of each symbol in the middle array
  - Lookup is  $O(n)$  (with a very small constant)
  - Memory footprint is just 3 arrays for each character in the original string!



# Improvements

- This algorithm can be expanded
  - Allow for mismatches





# HAVE A GOOD BREAK!

- GO HAVE FUN!!!!!!

