

# 15-826: Multimedia (Databases) and Data Mining

Lecture#5: Multi-key and
Spatial Access Methods – II – z-ordering

C. Faloutsos



#### Must-read material

- MM-Textbook, Chapter 5.1
- Ramakrinshan+Gehrke, Chapter 28.4
- J. Orenstein, <u>Spatial Query Processing in an</u> <u>Object-Oriented Database System</u>, Proc. ACM SIGMOD, May, 1986, pp. 326-336, Washington D.C.



#### **Outline**

Goal: 'Find similar / interesting things'

Intro to DB



- Indexing similarity search
- Data Mining



## Indexing - Detailed outline

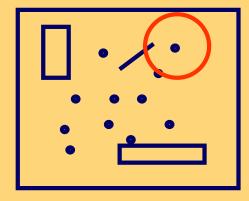
- primary key indexing
- secondary key / multi-key indexing



- spatial access methods
  - problem dfn
  - z-ordering
  - R-trees
  - **—** ...
- text
- •



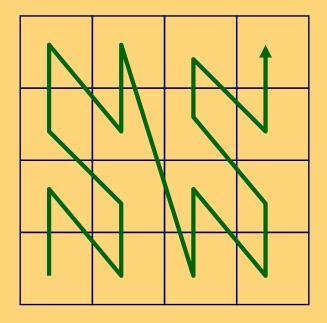
- Given a collection of geometric objects (points, lines, polygons, ...)
- Find cities within 100mi from Pittsburgh



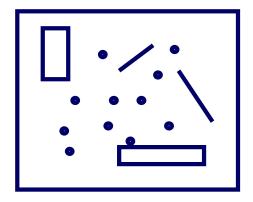


## Solution#1: z-ordering

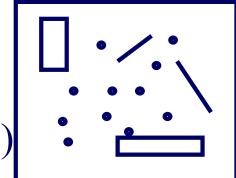
A: z-ordering/bit-shuffling/linear-quadtrees



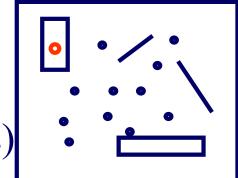
- Given a collection of geometric objects (points, lines, polygons, ...)
- organize them on disk, to answer spatial queries (like??)



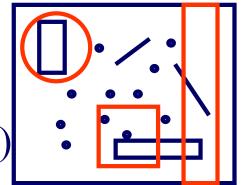
- Given a collection of geometric objects (points, lines, polygons, ...)
- organize them on disk, to answer
  - point queries
  - range queries
  - k-nn queries
  - spatial joins ('all pairs' queries)



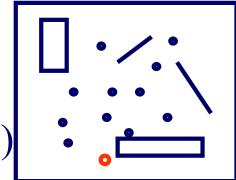
- Given a collection of geometric objects (points, lines, polygons, ...)
- organize them on disk, to answer
  - point queries
  - range queries
  - k-nn queries
  - spatial joins ('all pairs' queries)



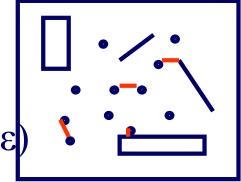
- Given a collection of geometric objects (points, lines, polygons, ...)
- organize them on disk, to answer
  - point queries
  - range queries
  - k-nn queries
  - spatial joins ('all pairs' queries)



- Given a collection of geometric objects (points, lines, polygons, ...)
- organize them on disk, to answer
  - point queries
  - range queries
  - k-nn queries
  - spatial joins ('all pairs' queries)



- Given a collection of geometric objects (points, lines, polygons, ...)
- organize them on disk, to answer
  - point queries
  - range queries
  - k-nn queries
  - spatial joins ('all pairs' within ε





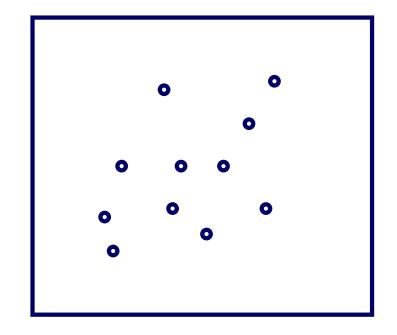
• Q: applications?

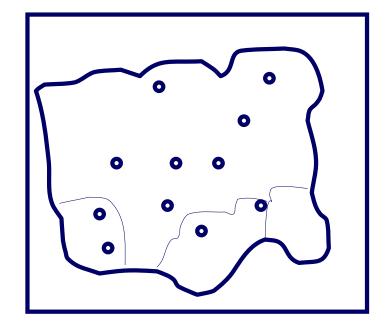


traditional DB

GIS

age





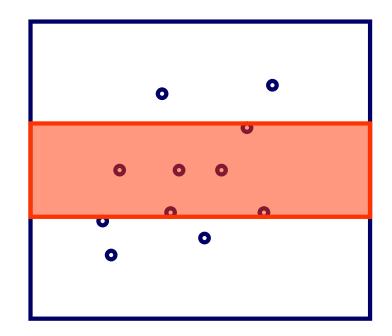
salary

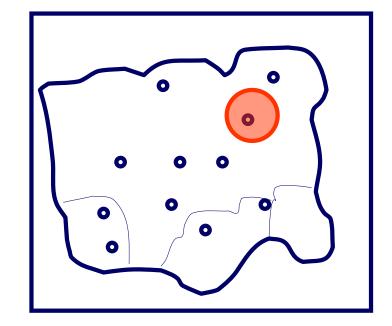


traditional DB

**GIS** 

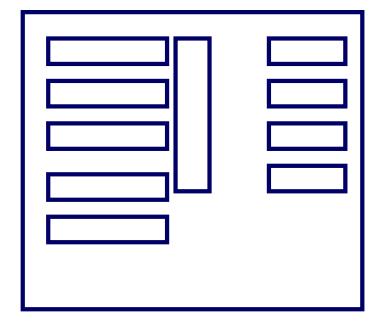
age





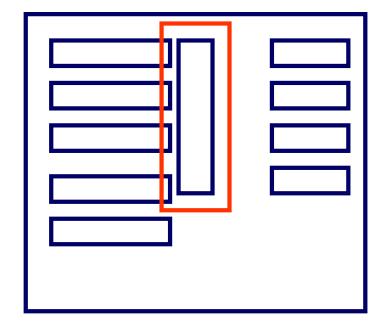
salary

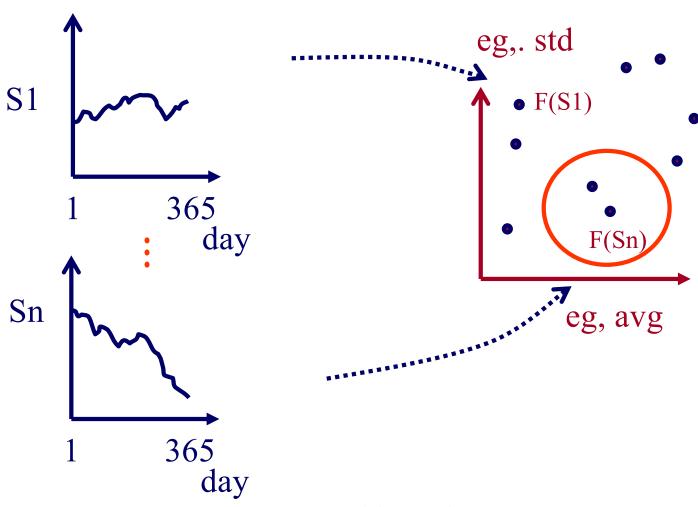
#### CAD/CAM



find elements too close to each other

#### CAD/CAM





15-826

Copyright: C. Faloutsos (2025)



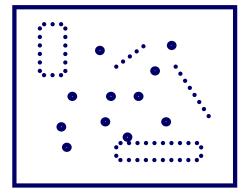
## Indexing - Detailed outline

- primary key indexing
- secondary key / multi-key indexing
- spatial access methods
  - problem dfn
- z-ordering
  - R-trees
  - **—** ...
  - text
  - ...

#### **SAMs: solutions**

- z-ordering
- R-trees
- (grid files)

Q: how would you organize, e.g., *n*-dim points, on disk? (*C* points per disk page)



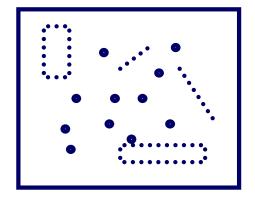
Q: how would you organize, e.g., *n*-dim points, on disk? (*C* points per disk page)

Hint: reduce the problem to 1-d points (!!)

Q1: why?

**A**:

Q2: how?



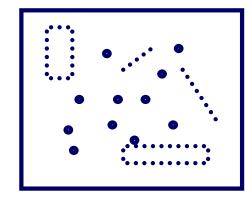
Q: how would you organize, e.g., *n*-dim points, on disk? (*C* points per disk page)

Hint: reduce the problem to 1-d points (!!)

Q1: why?

A: B-trees!

Q2: how?



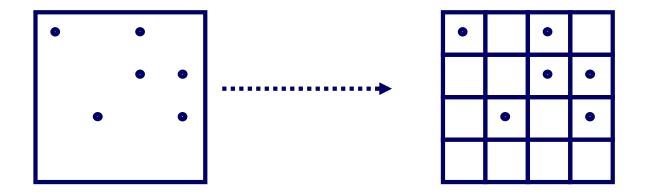
Q2: how?

A: assume finite granularity; z-ordering = bitshuffling = N-trees = Morton keys = geocoding = ...

Q2: how?

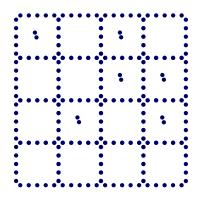
A: assume finite granularity (e.g.,  $2^{32}x2^{32}$ ; 4x4 here)

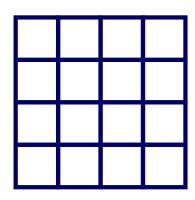
Q2.1: how to map n-d cells to 1-d cells?





Q2.1: how to map *n*-d cells to 1-d cells?

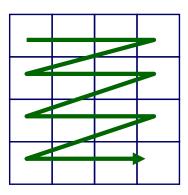




Q2.1: how to map *n*-d cells to 1-d cells?

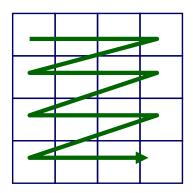
A: row-wise

Q: is it good?



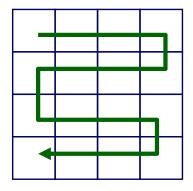
Q: is it good?

A: great for 'x' axis; bad for 'y' axis



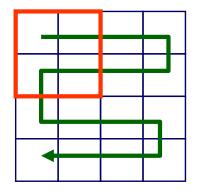


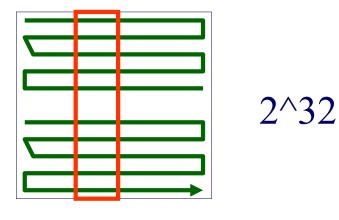
Q: How about the 'snake' curve?



Q: How about the 'snake' curve?

A: still problems:



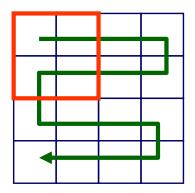


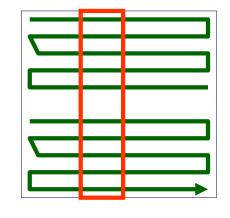
2^32

Q: Why are those curves 'bad'?

A: no distance preservation (~ clustering)

Q: solution?



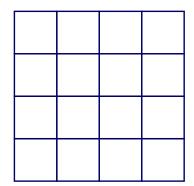


2^32

2^32

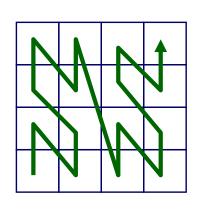


Q: solution? (w/ good clustering, and easy to compute, for 2-d and *n*-d?)



Q: solution? (w/ good clustering, and easy to compute, for 2-d and *n*-d?)

A: z-ordering/bit-shuffling/linear-quadtrees



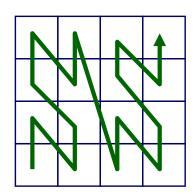
'looks' better:

- few long jumps;
- scoops out the whole quadrant before leaving it
- a.k.a. space filling curves

z-ordering/bit-shuffling/linear-quadtrees

Q: How to generate this curve (z = f(x,y))?

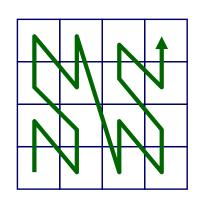
A: 3 (equivalent) answers!



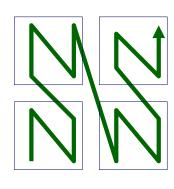
z-ordering/bit-shuffling/linear-quadtrees

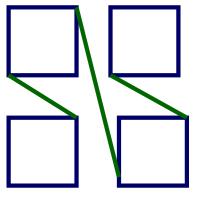
Q: How to generate this curve (z = f(x,y))?

A1: 'z' (or 'N') shapes, RECURSIVELY









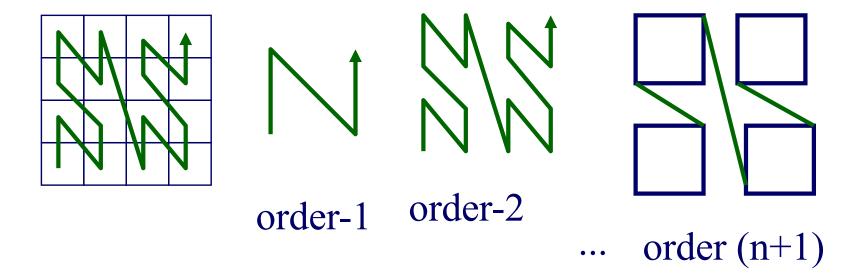
order-1

order-2

order (n+1)

#### Notice:

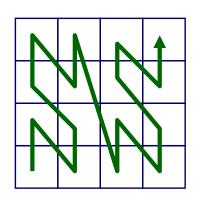
- self similar (we'll see about fractals, soon)
- method is hard to use: z = ? f(x,y)



z-ordering/bit-shuffling/linear-quadtrees

Q: How to generate this curve (z = f(x,y))?

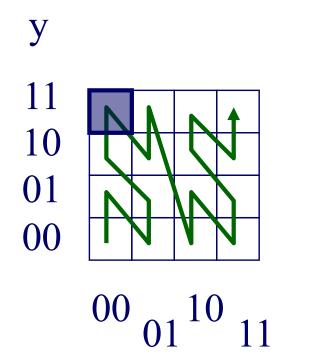
A: 3 (equivalent) answers!



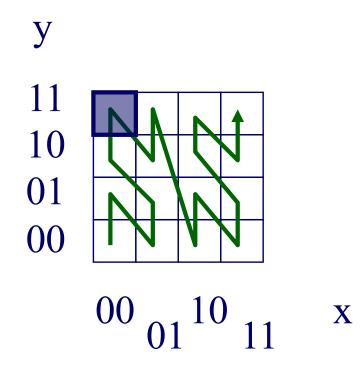
Method #2?

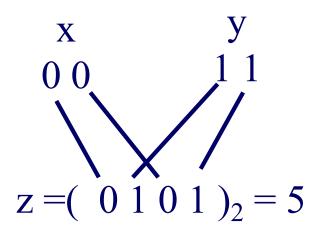
#### bit-shuffling



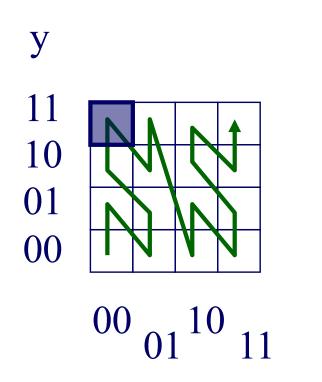


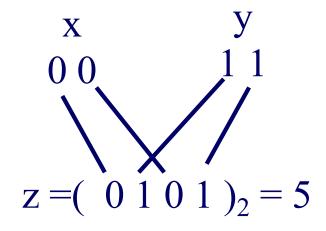
#### bit-shuffling





#### bit-shuffling



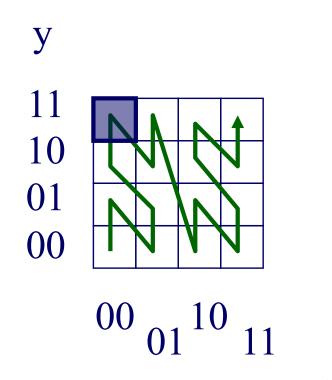


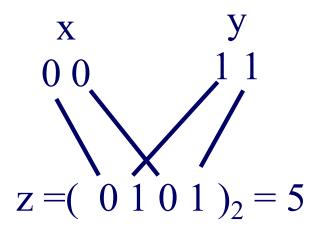
How about the reverse:

$$(x,y)=g(z)?$$

X

#### bit-shuffling



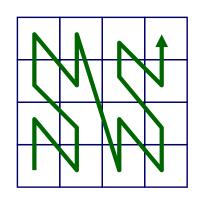


How about *n*-d spaces?

z-ordering/bit-shuffling/linear-quadtrees

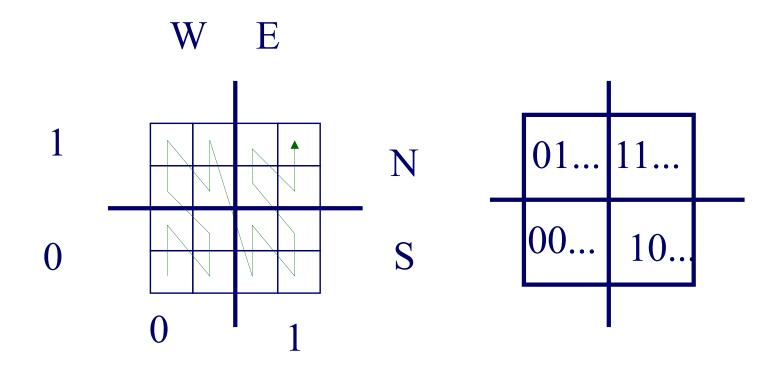
Q: How to generate this curve (z = f(x,y))?

A: 3 (equivalent) answers!



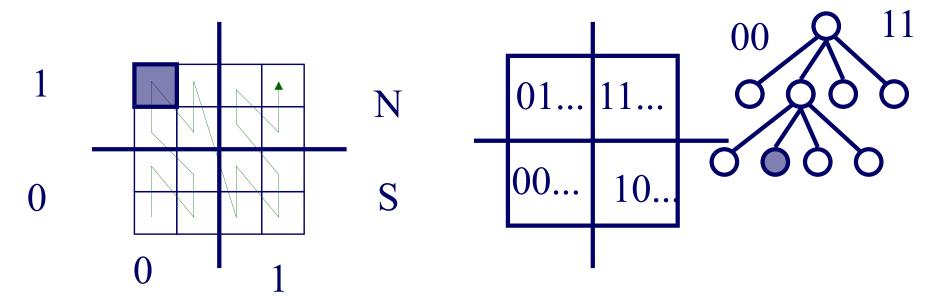
Method #3?

linear-quadtrees : assign N->1, S->0 e.t.c.

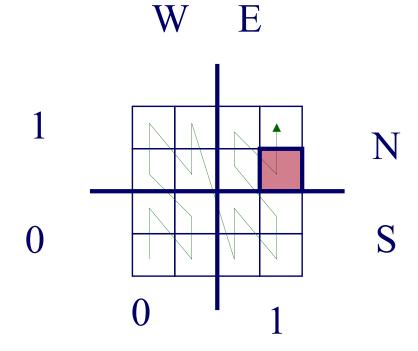


... and repeat recursively. Eg.: z<sub>blue-cell</sub>=

$$WN;WN = (0101)_2 = 5$$

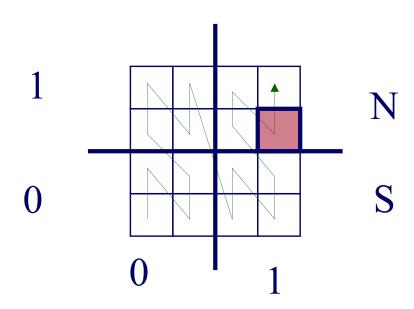


Drill: z-value of magenta cell, with the three methods?



Drill: z-value of magenta cell, with the three methods?

W E



method#1: 14

method#2: shuffle(11;10)=

$$(1110)_2 = 14$$

Drill: z-value of magenta cell, with the three methods?

W E

N
S

method#1: 14

method#2: shuffle(11;10)=

 $(1110)_2 = 14$ 

method#3: EN;ES = ... = 14

#### z-ordering - Detailed outline

- spatial access methods
  - z-ordering
    - main idea 3 methods



- non-point (eg., region) data
- analysis; variations
- R-trees

<del>-</del> ...

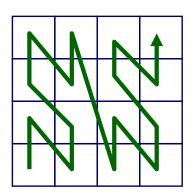




Q1: How to store on disk?

A:

Q2: How to answer range queries etc



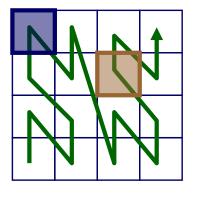


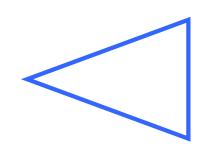
Q1: How to store on disk?

A: treat z-value as primary key; feed to B-tree

#### **PGH**

SF





Z	cname	etc
5	SF	
12	PGH	

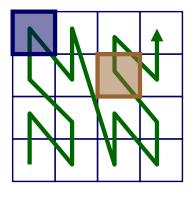


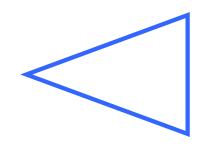
#### MAJOR ADVANTAGES w/ B-tree:

- already inside commercial systems (no coding/debugging!)
- concurrency & recovery is ready

#### **PGH**

SF





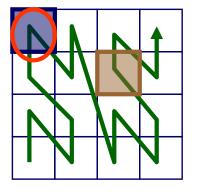
<u>Z</u>	Chame	elc
5	SF	
12	PGH	

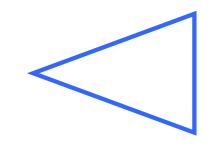


Q2: queries? (eg.: *find city at (0,3) )?* 



SF





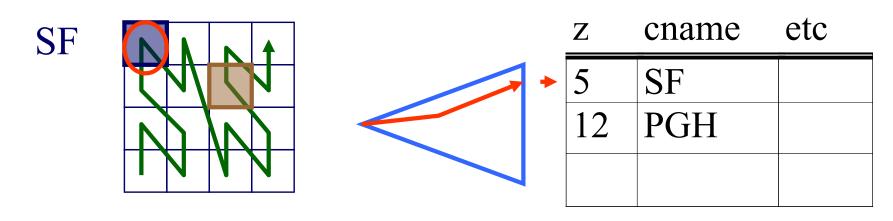
<u>Z</u>	chame	eic
5	SF	
12	PGH	



Q2: queries? (eg.: *find city at (0,3) )?* 

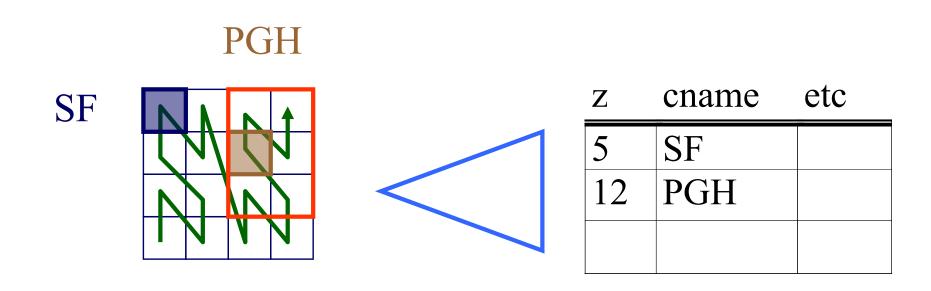
A: find z-value; search B-tree

#### **PGH**





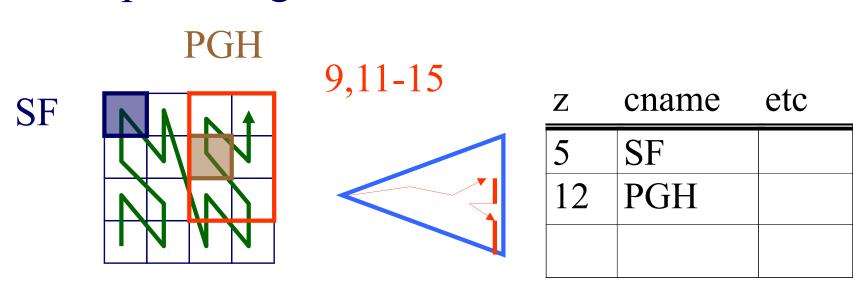
Q2: range queries?





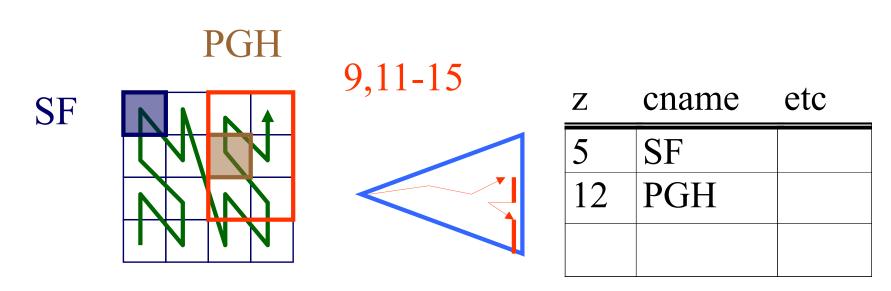
Q2: range queries?

A: compute ranges of z-values; use B-tree





Q2': range queries - how to reduce # of qualifying of ranges?



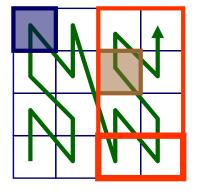


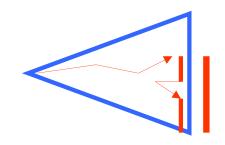
Q2': range queries - how to reduce # of qualifying of ranges?

A: Augment the query!

**PGH** 

SF





5	SF	
12	PGH	

cname

etc



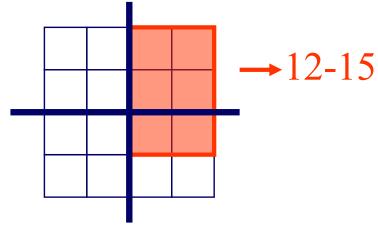
Q2'': range queries - how to break a query into ranges?





Q2'': range queries - how to break a query into ranges?

A: recursively, quadtree-style; decompose only non-full quadrants

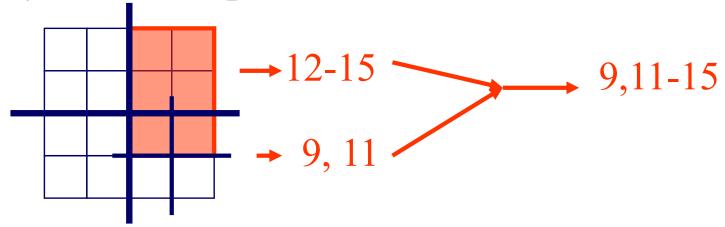


9,11-15



Q2'': range queries - how to break a query into ranges?

A: recursively, quadtree-style; decompose only non-full quadrants



#### z-ordering - Detailed outline

- spatial access methods
  - z-ordering
    - main idea 3 methods

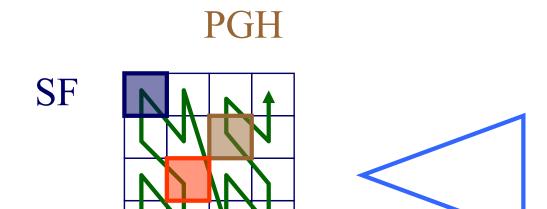


- use w/ B-trees; algorithms (range, knn queries ...)
- non-point (eg., region) data
- analysis; variations
- R-trees

<del>-</del> ...



Q3: k-nn queries? (say, 1-nn)?



Z	Chame	elc
5	SF	
12	PGH	

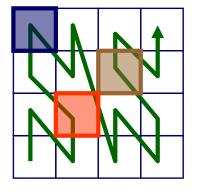


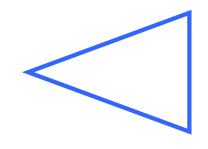
Q3: k-nn queries? (say, 1-nn)?

A: traverse B-tree; find nn wrt z-values and ...

#### **PGH**



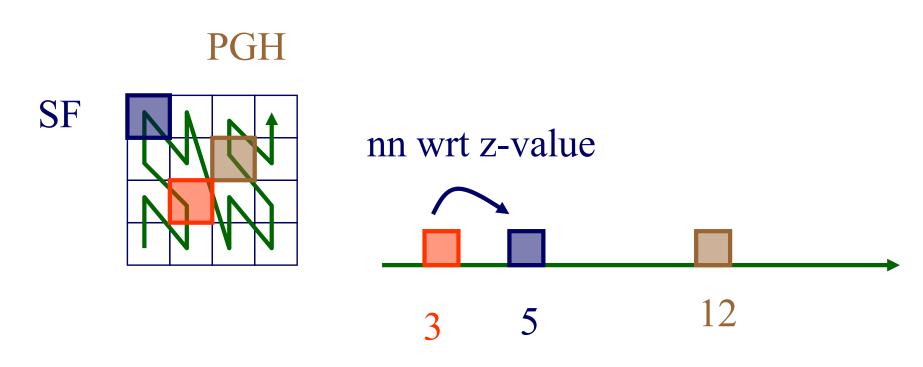




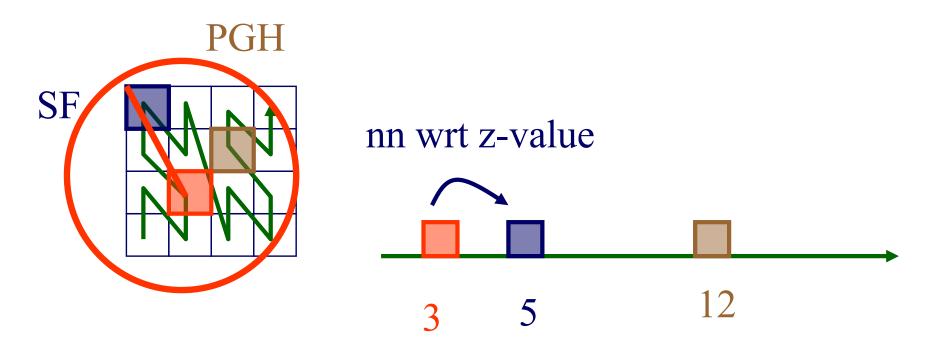
<u>Z</u>	cname	eic
5	SF	
12	PGH	



... ask a range query.



... ask a range query.

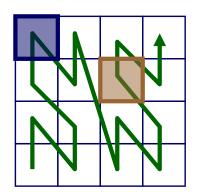




Q4: all-pairs queries? (all pairs of cities within 10 miles from each other?)

#### **PGH**

SF



(we'll see 'spatial joins' later: find all PA counties that intersect a lake)



#### z-ordering - Detailed outline

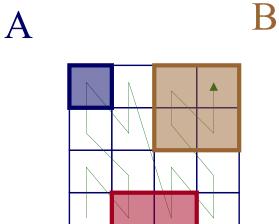
- spatial access methods
  - z-ordering
    - main idea 3 methods
    - use w/B-trees; algorithms (range, knn queries ...)
    - non-point (eg., region) data
    - analysis; variations
  - R-trees

**—** ...



Q: z-value for a region?

$$z_B = ??$$
 $z_C = ??$ 

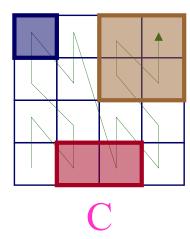


Q: z-value for a region?

A: 1 or more z-values; by quadtree decomposition

A



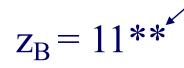


$$z_{\rm R} = ??$$

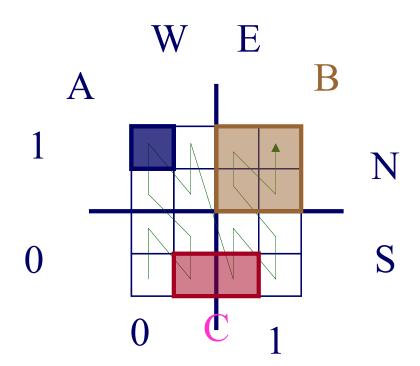
$$z_{B} = ??$$
 $z_{C} = ??$ 

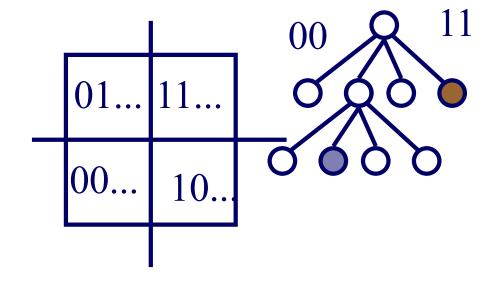
"don't care"

Q: z-value for a region?



$$z_{\rm C} = ??$$



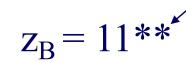


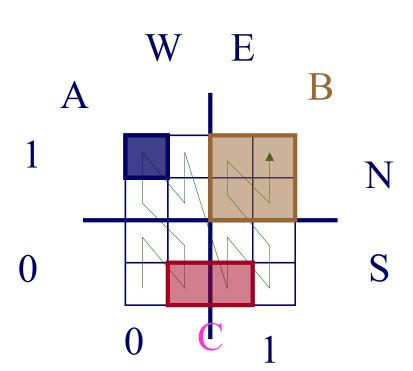
#### Carnegie Mellon

#### z-ordering - regions

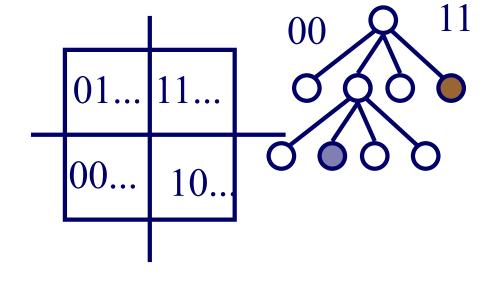
"don't care"

Q: z-value for a region?





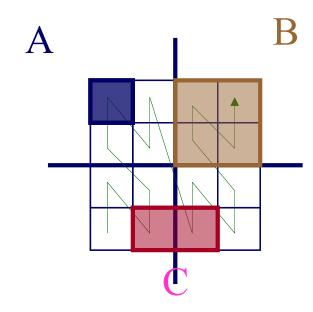
$$z_C = \{0010; 1000\}$$





Q: How to store in B-tree?

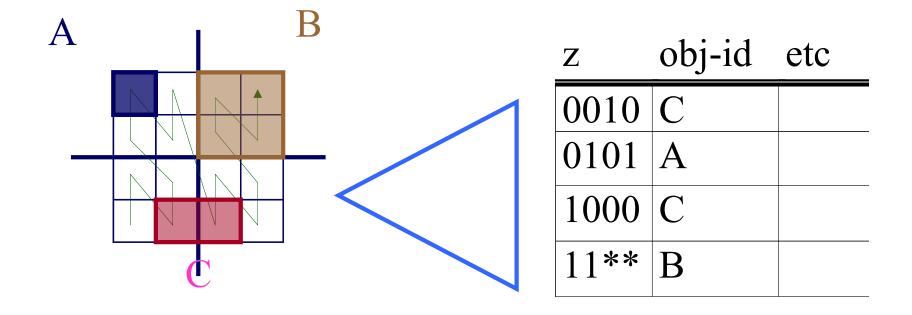
Q: How to search (range etc queries)



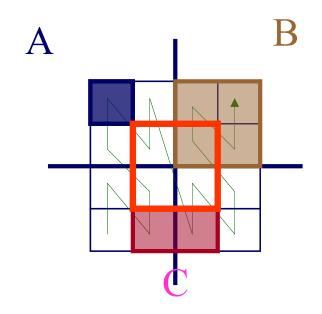


Q: How to store in B-tree? A: sort (\*<0<1)

Q: How to search (range etc queries)



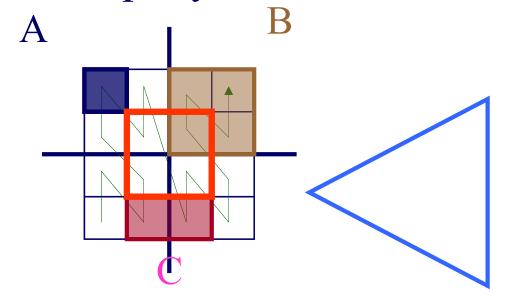
Q: How to search (range etc queries) - eg 'red' range query



Z	obj-id	etc
0010	C	
0101	A	
1000	С	
11**	В	

Q: How to search (range etc queries) - eg 'red' range query

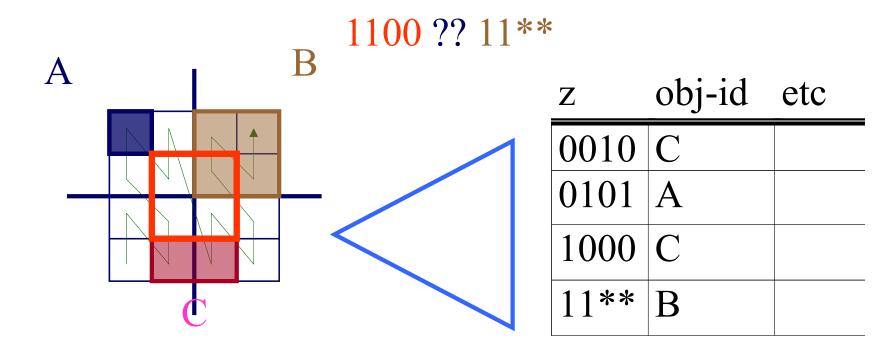
A: break query in z-values; check B-tree



Z	obj-id	etc
0010	C	
0101	A	
1000	С	
11**	В	



Almost identical to range queries for point data, except for the "don't cares" - i.e.,



Almost identical to range queries for point data, except for the "don't cares" - i.e.,

$$z1 = 1100 ?? 11** = z2$$

Specifically: does z1 contain/avoid/intersect z2?

Q: what is the criterion to decide?

z1 = 1100 ?? 11\*\* = z2

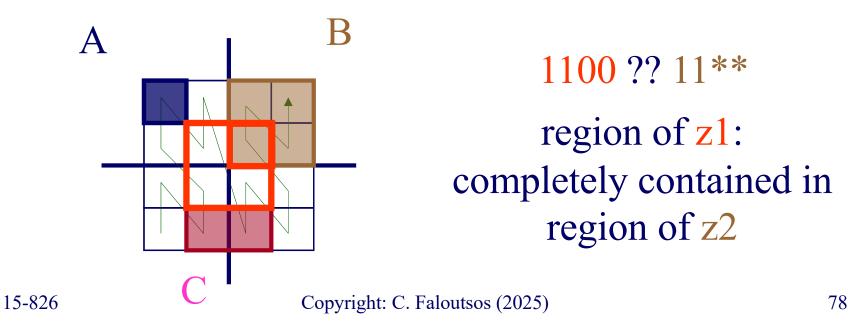
Specifically: does z1 contain/avoid/intersect z2?

Q: what is the criterion to decide?

A: **Prefix property**: let r1, r2 be the corresponding regions, and let r1 be the smallest (=> z1 has fewest '\*'s). Then:



- r2 will either contain completely, or avoid completely r1.
- it will contain r1, if z2 is the prefix of z1



Drill (True/False). Given:

- z1 = 011001\*\*
- z2=01\*\*\*\*\*
- z3 = 0100\*\*\*\*

T/F r2 contains r1

T/F r3 contains r1

T/F r3 contains r2

Drill (True/False). Given:

- z1 = 011001\*\*
- z<sub>2</sub>= 01\*\*\*\*\*
- z3 = 0100\*\*\*\*

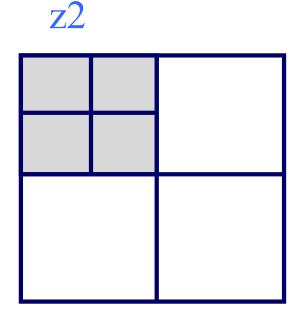
T/F r2 contains r1 - TRUE (prefix property)

T/F r3 contains r1 - FALSE (disjoint)

T/F r3 contains r2 - FALSE (r2 contains r3)

#### Drill (True/False). Given:

- z1 = 011001\*\*
- $z^2 = 01******$
- z3=0100\*\*\*\*



#### Drill (True/False). Given:

• 
$$z1 = 011001**$$

• 
$$z^2 = 01******$$

• 
$$z3 = 0100****$$

z3

 $\mathbf{z}^2$ 

T/F r2 contains r1 - TRUE (prefix property)

T/F r3 contains r1 - FALSE (disjoint)

T/F r3 contains r2 - FALSE (r2 contains r3)

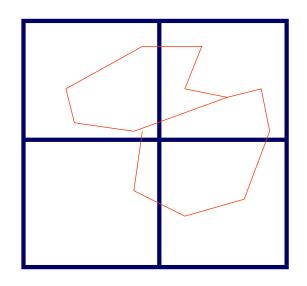


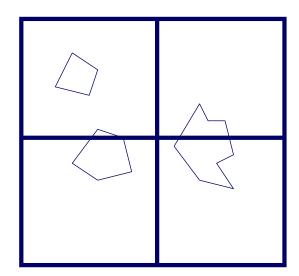
Spatial joins: find (quickly) all

counties

intersecting

lakes





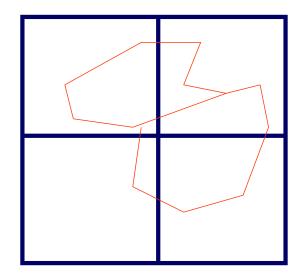


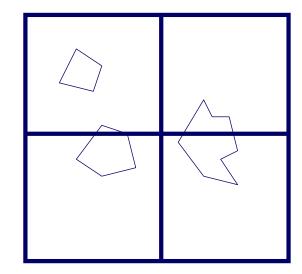
Spatial joins: find (quickly) all

counties

intersecting

lakes





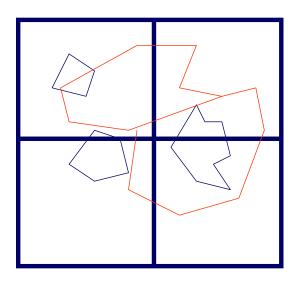


Spatial joins: find (quickly) all

counties

intersecting

lakes





Spatial joins: find (quickly) all

counties intersecting lakes

Naive algorithm: O( N \* M)

Something faster?

Spatial joins: find (quickly) all counties intersecting lakes

z obj-id etc

0010 ALG

...

1000 WAS

11\*\* ALG

Z	obj-id	etc
0011	Erie	
0101	Erie	
• • •		
10**	Ont.	



Spatial joins: find (quickly) all

counties intersecting lakes

Solution: merge the lists of (sorted) z-values, looking for the prefix property

footnote#1: '\*' needs careful treatment

footnote#2: need dup. elimination



#### z-ordering - Detailed outline

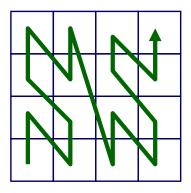
- spatial access methods
  - z-ordering
    - main idea 3 methods
    - use w/ B-trees; algorithms (range, knn queries ...)
    - non-point (eg., region) data
    - analysis; variations
  - R-trees

\_\_\_\_





Q: is z-ordering the best we can do?

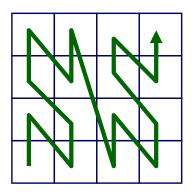




Q: is z-ordering the best we can do?

A: probably not - occasional long 'jumps'

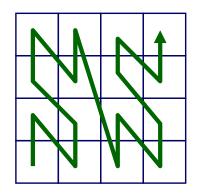
Q: then?

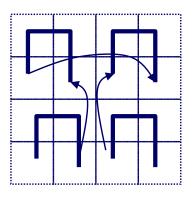


Q: is z-ordering the best we can do?

A: probably not - occasional long 'jumps'

Q: then? A1: Gray codes







• Ingenious way to spot flickering LED – binary:

Officially.	000	0
	001	1
2 511	010	2
F. Gray. <i>Pulse code communication</i> , March 17, 1953	011	3
	100	4
	101	5
U.S. Patent 2,632,058	110	6
	111	7

Ingenious way to spot flickering LED

0

1

Ingenious way to spot flickering LED

0

 $\cdot 0$ 

1

. 1

• •

• •

Ingenious way to spot flickering LED

0 .0 1 .1 ...

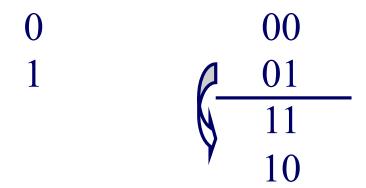
Ingenious way to spot flickering LED

 $\begin{array}{c}
0 \\
1
\end{array}$   $\begin{array}{c}
.1 \\
\hline
0
\end{array}$ 

#### Carnegie Mellon

# (Gray codes)

Ingenious way to spot flickering LED

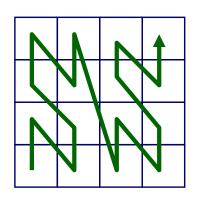


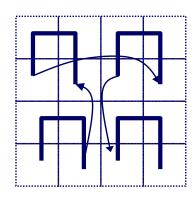
Ingenious way to spot flickering LED

Q: is z-ordering the best we can do?

A: probably not - occasional long 'jumps'

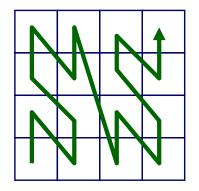
Q: then? A1: Gray codes – CAN WE DO BETTER?

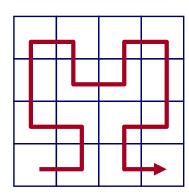






A2: Hilbert curve! (a.k.a. Hilbert-Peano curve)





#### Carnegie Mellon

#### (break)



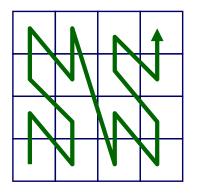
David Hilbert (1862-1943)

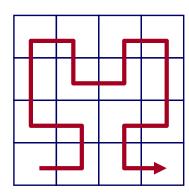


Giuseppe Peano (1858-1932)

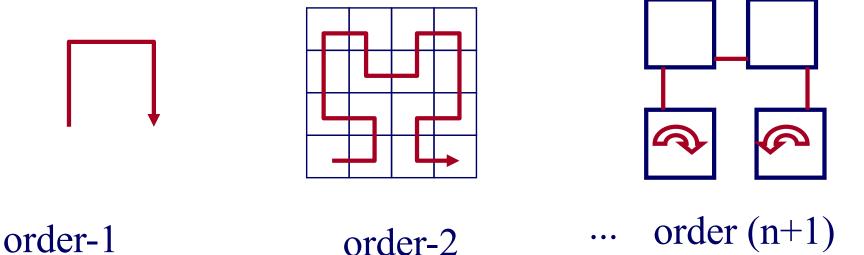


'Looks' better (never long jumps). How to derive it?





'Looks' better (never long jumps). How to derive it?



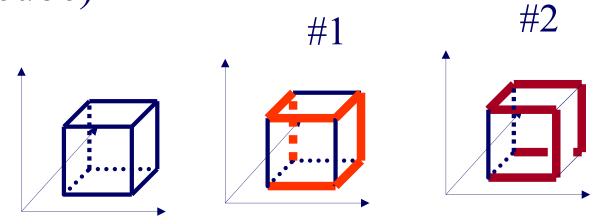
Q: function for the Hilbert curve (h = f(x,y))?

A: bit-shuffling, followed by post-processing, to account for rotations. Linear on # bits.

See textbook, for pointers to code/algorithms (eg., [Jagadish, 90])

Q: how about Hilbert curve in 3-d? n-d?

A: Exists (and is not unique!). Eg., 3-d, order-1 Hilbert curves (Hamiltonian paths on cube)





### z-ordering - Detailed outline

- spatial access methods
  - z-ordering
    - main idea 3 methods
    - use w/ B-trees; algorithms (range, knn queries ...)
    - non-point (eg., region) data
    - analysis; variations
  - R-trees

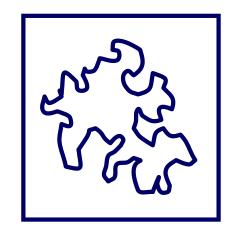


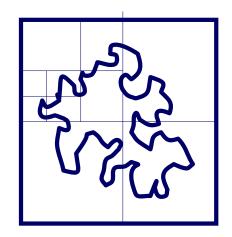


### z-ordering - analysis

Q: How many pieces ('quad-tree blocks') per region?

A: proportional to perimeter (surface etc)

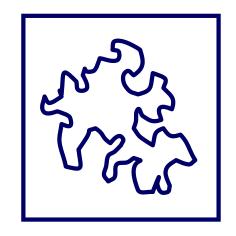


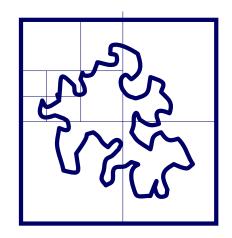




(How long is the coastline, say, of England?

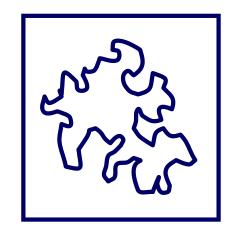
Paradox: The answer changes with the yardstick -> fractals ...)

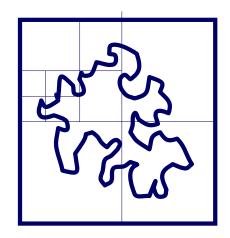






Q: Should we decompose a region to full detail (and store in B-tree)?

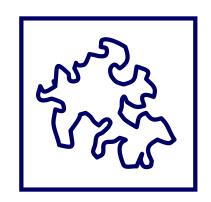


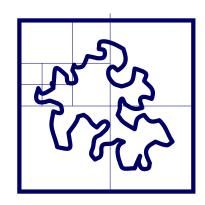




Q: Should we decompose a region to full detail (and store in B-tree)?

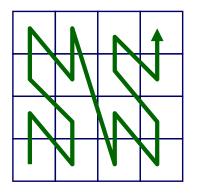
A: NO! approximation with 1-3 pieces/z-values is best [Orenstein90]

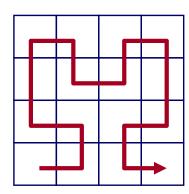






Q: how to measure the 'goodness' of a curve?

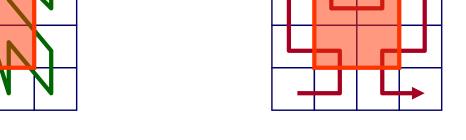






Q: how to measure the 'goodness' of a curve?

A: e.g., avg. # of runs, for range queries



4 runs 3 runs (#runs ~ #disk accesses on B-tree)

Q: So, is Hilbert really better?

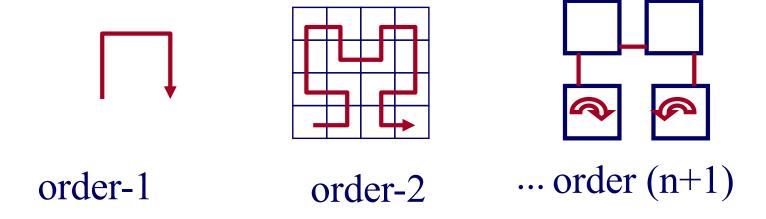
A: 27% fewer runs, for 2-d (similar for 3-d)

Q: are there formulas for #runs, #of quadtree blocks etc?

A: Yes ([Jagadish; Moon+ etc] see textbook)

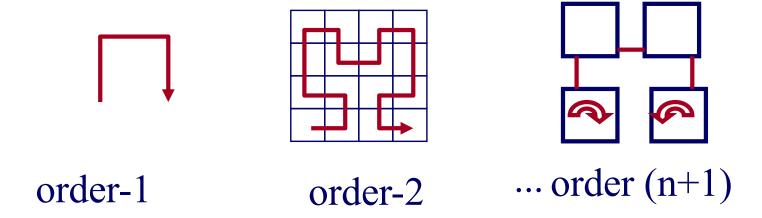


Hilbert and z-ordering curves: "space filling curves": eventually, they visit every point in n-d space - therefore:





... they show that the plane has as many points as a line (-> headaches for 1900's mathematics/topology). (fractals, again!)

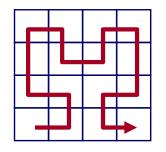


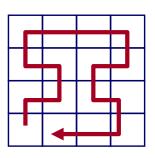


Observation #2: Hilbert (like) curve for video encoding [Y. Matias+, CRYPTO '87]:

Given a frame, visit its pixels in randomized hilbert order; compress; and transmit







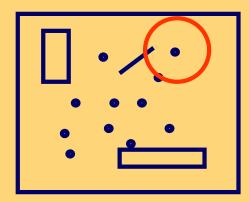


In general, Hilbert curve is great for preserving distances, clustering, vector quantization etc



# Spatial Access Methods - problem

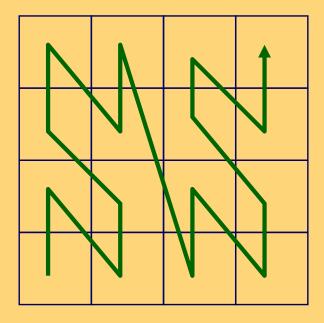
- Given a collection of geometric objects (points, lines, polygons, ...)
- Find cities within 100mi from Pittsburgh





# Solution#1: z-ordering

A: z-ordering/bit-shuffling/linear-quadtrees



#### **Conclusions**

- z-ordering is a great idea (n-d points -> 1-d points; feed to B-trees)
- used
- ... by <u>TIGER</u> of US Bureau of Census
- ... by <u>AWS Aurora</u>
- ... and (probably) by other GIS products
- works great with low-dim points