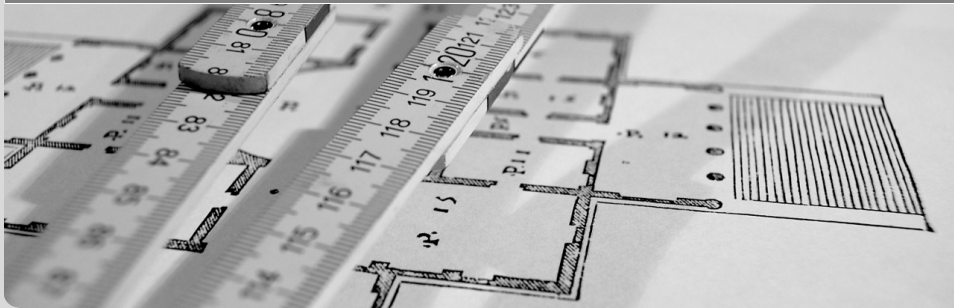


# Parallel Algorithm for Closest Pair Problem

Ge Wu | July 25, 2013

INSTITUTE FOR THEORETICAL INFORMATICS



## Closest Pair Problem

- Given  $n$  **different unordered** points  $P = \{p_1, p_2, \dots, p_n\}$  in **unit square**

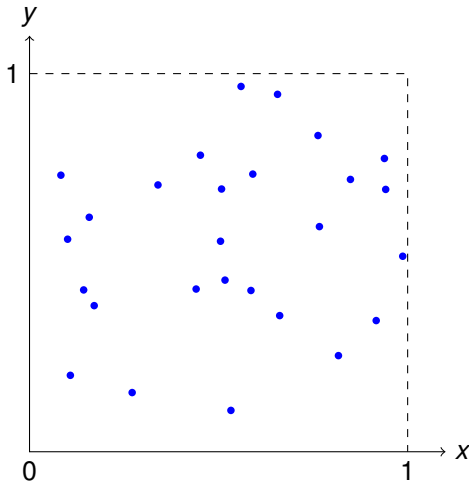
$$p_i = (x_i, y_i) \in (0, 1) \times (0, 1) \subset \mathbb{R}^2$$

- Find the shortest euclidean distance between any two points

- $O(n \log n)$  lower bound in comparison tree model
- **Bentley and Shamos 1976**  
 $O(n \log n)$  algorithm using divide and conquer
- **Rabin 1976**  
 $O(n)$  randomized algorithm with  $O(1)$  floor function
- **Fortune and Hopcroft 1978**  
 $O(n \log \log n)$  deterministic algorithm with  $O(1)$  floor function
- **Khuller and Matias 1995**  
Another  $O(n)$  randomized algorithm

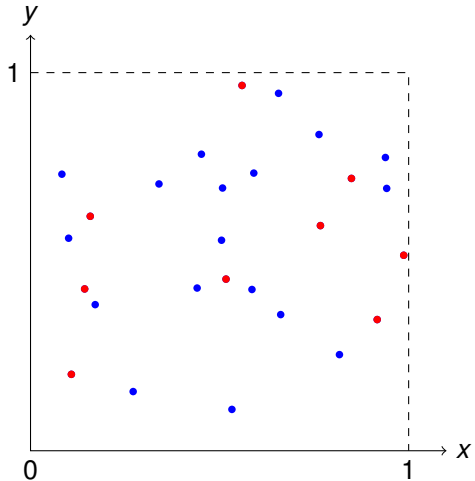
# Algorithm

- Sample ( $n^{\frac{2}{3}}$  Points from  $P$ )
- Partition
- Compute

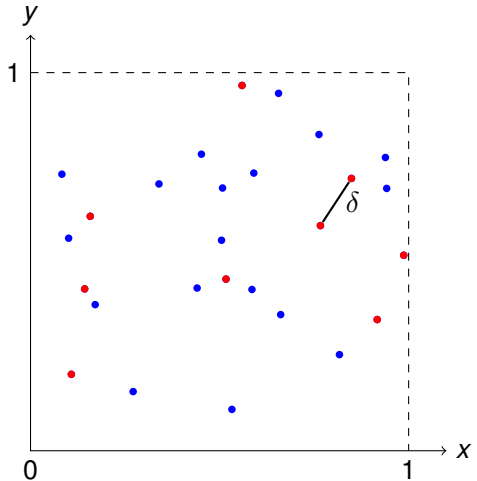


# Algorithm

- Sample ( $n^{\frac{2}{3}}$  Points from  $P$ )
- Partition
- Compute

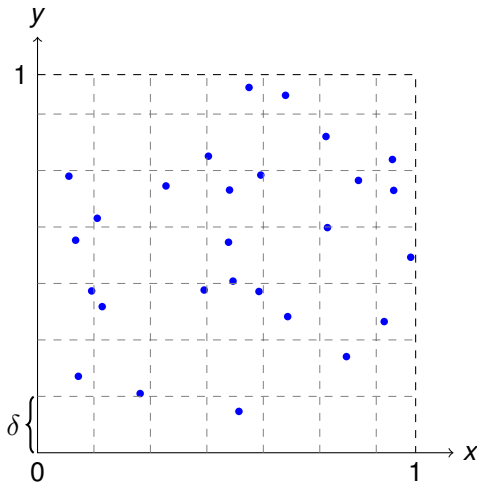


- Sample ( $n^{\frac{2}{3}}$  Points from  $P$ )
- Partition
- Compute



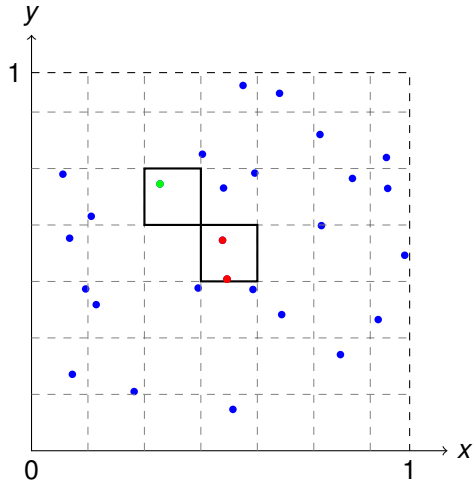
# Algorithm

- Sample ( $n^{\frac{2}{3}}$  Points from  $P$ )
- Partition
- Compute



# Algorithm

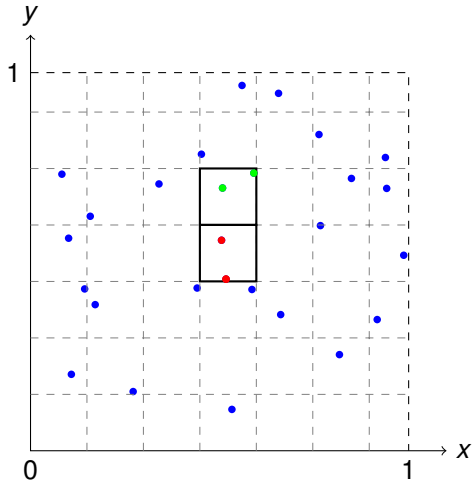
- Sample ( $n^{\frac{2}{3}}$  Points from  $P$ )
- Partition
- Compute



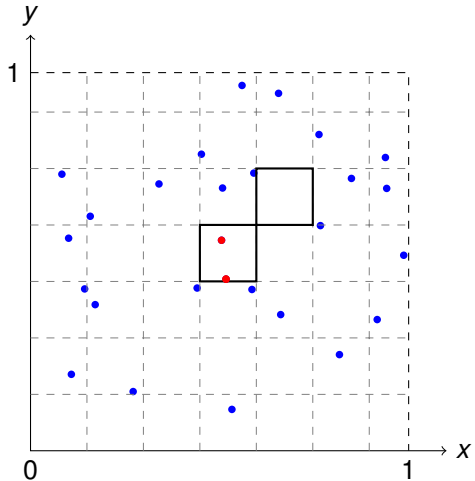


# Algorithm

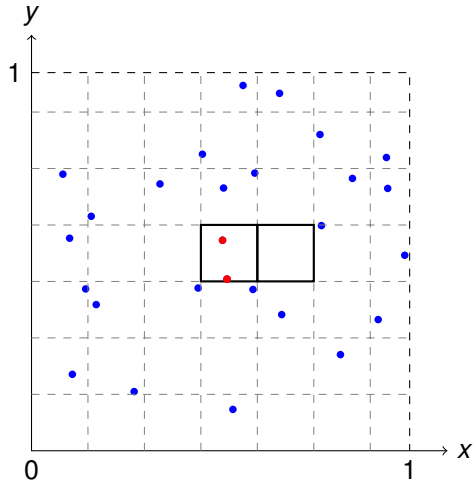
- Sample ( $n^{\frac{2}{3}}$  Points from  $P$ )
- Partition
- Compute



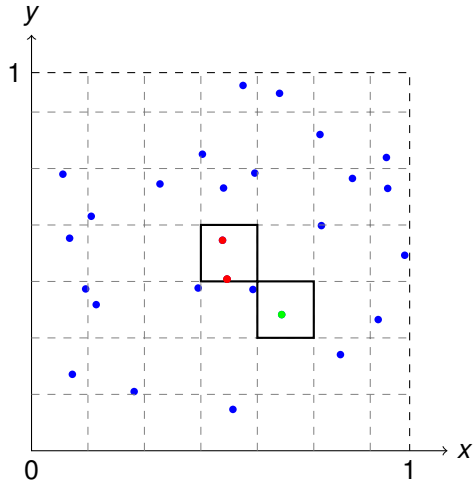
- Sample ( $n^{\frac{2}{3}}$  Points from  $P$ )
- Partition
- Compute



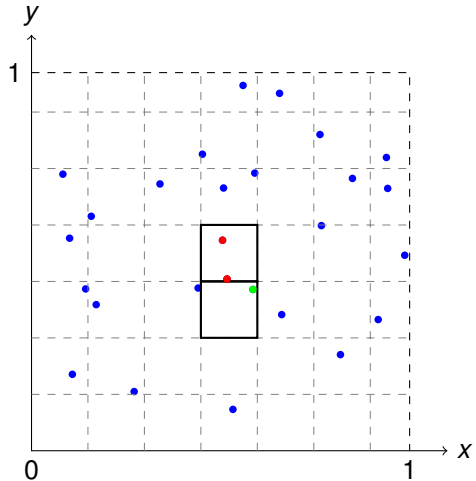
- Sample ( $n^{\frac{2}{3}}$  Points from  $P$ )
- Partition
- Compute



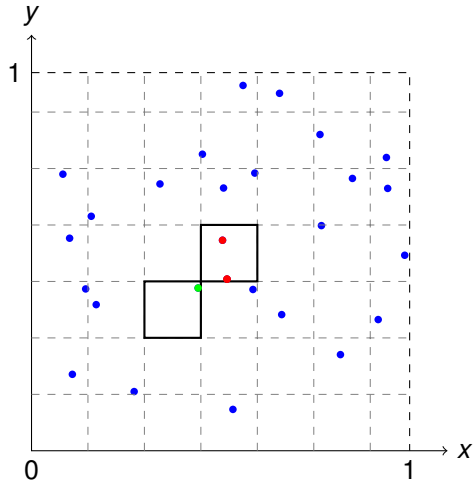
- Sample ( $n^{\frac{2}{3}}$  Points from  $P$ )
- Partition
- Compute



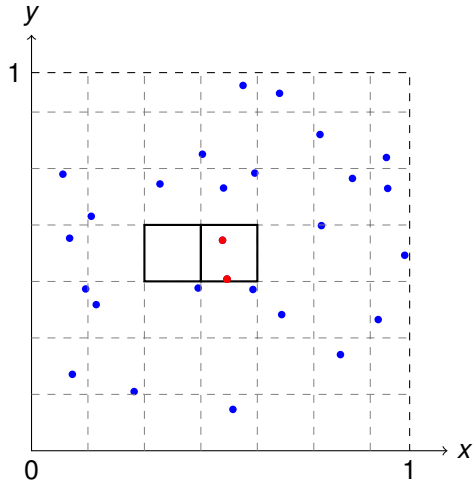
- Sample ( $n^{\frac{2}{3}}$  Points from  $P$ )
- Partition
- Compute



- Sample ( $n^{\frac{2}{3}}$  Points from  $P$ )
- Partition
- Compute

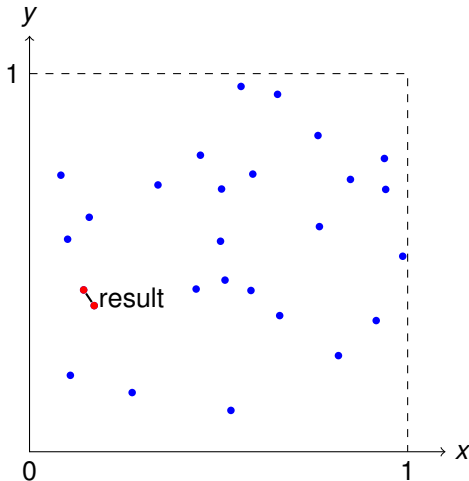


- Sample ( $n^{\frac{2}{3}}$  Points from  $P$ )
- Partition
- Compute



# Algorithm

- Sample ( $n^{\frac{2}{3}}$  Points from  $P$ )
- Partition
- Compute





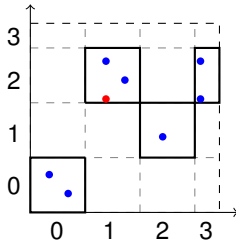
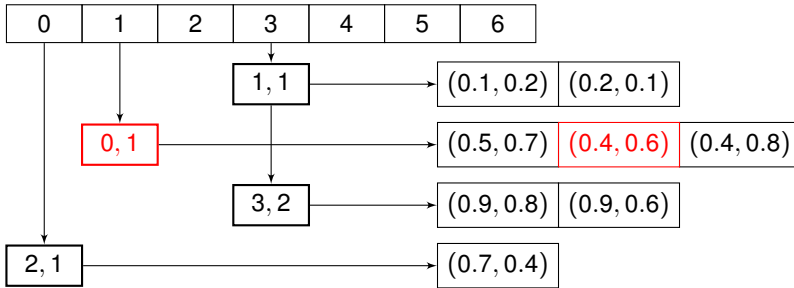
- Compute all pair of distances:  $O((n^{\frac{2}{3}})^2) = O(n^{\frac{4}{3}})$
- Divide & Conquer:  $O(n^{\frac{2}{3}} \log n^{\frac{2}{3}}) = O(n^{\frac{2}{3}} \log n) \subset O(n)$   
More samples possible:  $O(\frac{n}{\log n} \log \frac{n}{\log n}) \subset O(n)$
- Other approach?

- Map points to the cells
- Divide the coordinates by cell length and **truncate** them to integer
- Partition  $G = \{G_1, G_2, \dots, G_k\}$  with:

$$\bigcup_{i=1}^k G_i = P, G_i \cap G_j = \emptyset, k \leq n$$

- $O(n)$  with hashing

# Partition



$$H(cell_{x\_idx}, cell_{y\_idx}) = (cell_{x\_idx} + 2 \cdot cell_{y\_idx}) \bmod 7$$

## Theorem

Let  $G_i$  with  $1 \leq i \leq k$  be the ***i-th*** non-empty cell and  $c, d$  some positive constants, then

$$\text{Prob} \left( \sum_{i=1}^k |G_i|^2 \leq c \cdot n \right) \geq 1 - \frac{1}{2^{n^d}}$$

Runtime for computing all pairs distances

$$O\left(\sum_{i=1}^k |G_i| \cdot |G_i \cup \text{neighbor}(G_i)|\right) = O\left(\sum_{i=1}^k |G_i|^2\right) = O(n)$$

Recursively calculate the closest distance on samples  $S = \{s_1, \dots, s_{n^{\frac{2}{3}}}\}$

- Sample another  $n^{\frac{4}{9}}$  points  $S'$  from  $S$
- Get closest distance on  $S'$  by calculating all pair distances:  $O(n^{\frac{8}{9}})$
- Partition and compute the closest distance on  $S$

```
0 Cell set  $H = \emptyset$ ,  $result = \infty$ 
1 Choose set  $S$  of  $n^{\frac{2}{3}}$  samples  $O(n^{\frac{2}{3}})$ 
   $\delta = \min_{x,y \in S, x \neq y} (dist(x, y))$   $O(n^{\frac{8}{9}})$ 
2 for  $i = 1$  to  $n$  do  $O(n)$ 
   $c = cell(p_i)$ 
  if not  $H.findCell(c)$  then
     $H.addCell(\{c, idx = \{c.x, c.y\}\})$ 
     $H.findCell(c).addPoint(p_i)$ 
3 foreach  $c_1$  in  $H$  do  $O(n)$ 
  foreach  $c_2$  in  $\{c_1\} \cup neighbor(c_1)$  do
    foreach  $(p_i, p_j)$  in  $c_1 \times c_2$ 
       $result = \min(result, dist(p_i, p_j))$ 
```

Runtime:  $O(n)$  with high probability, worst case:  $O(n^2)$

## ■ PRAM, CREW model

```
1: procedure SAMPLE(in  $P$ ,  $id$ , out  $S$ )
2:    $total \leftarrow 0$ 
3:   for  $i \leftarrow \lceil \frac{n}{\#PE} \rceil \cdot id$  to  $\min(\lceil \frac{n}{\#PE} \rceil \cdot (id + 1) - 1, n - 1)$  do     $\triangleright O(\frac{n}{p})$ 
4:      $chosen[i] \leftarrow (\text{Random}([0..1]) < n^{\frac{2}{3}}/n)$ 
5:      $total \leftarrow total + chosen[i]$ 
6:   end for
7:    $pos@id \leftarrow \text{PrefixSum}(total@id)$                                  $\triangleright O(\log p)$ 
8:   for  $i \leftarrow L \cdot id$  to  $\min(L \cdot (id + 1) - 1, n - 1)$  do           $\triangleright O(\frac{n}{p})$ 
9:     if  $chosen[i]$  then
10:       $S[--pos] \leftarrow P[i]$ 
11:    end if
12:  end for
13: end procedure
```

- Concurrent hash table (possibly with variable size)
- Concurrent read, exclusive write (block the chain)
- At most  $n$  add operation, constant blocking time for each chain, if elements evenly distributed

```
1: procedure PARTITION(in  $P, \delta$ , out  $H$ )
2:    $H \leftarrow \emptyset$ 
3:   for  $i \leftarrow \lceil \frac{n}{\#PE} \rceil \cdot id$  to  $\min(\lceil \frac{n}{\#PE} \rceil \cdot (id + 1) - 1, n - 1)$  do  $\triangleright O(\frac{n}{p})$ 
4:      $cell \leftarrow (\lfloor \frac{p_i.x}{\delta} \rfloor, \lfloor \frac{p_i.y}{\delta} \rfloor)$ 
5:     if  $H.findCell(cell)$  then
6:        $H.addCell(cell)$   $\triangleright$  constant times for each chain
7:     end if
8:      $H.findCell(cell).addPoint(p_i)$ 
9:   end for
10: end procedure
```



```

1: procedure COMPUTE(in  $H, k, id$ , out  $result$ )
2:    $pairs \leftarrow 0, pcnt \leftarrow 0$ 
3:   for  $i \leftarrow \lceil \frac{k}{\#PE} \rceil \cdot id$  to  $\min(\lceil \frac{k}{\#PE} \rceil \cdot (id + 1) - 1, k - 1)$  do  $\triangleright O(\frac{k}{p})$ 
4:      $pairs \leftarrow pairs + |H[i]| \cdot |\mathbf{Neighbor}(H[i])|$ 
5:   end for
6:    $total \leftarrow \mathbf{ReduceSum}(pairs@id)$   $\triangleright O(\log p)$ 
7:    $pcnt \leftarrow pairs > 0 ? 0 : \min(1, \lfloor \frac{pairs}{total} \rfloor) \cdot \#PE$ 
8:    $pre\_pcnt@id \leftarrow \mathbf{PrefixSum}(pcnt@id)$   $\triangleright O(\log p)$ 
9:    $id' \leftarrow \mathbf{BinarySearch}(id, pre\_pcnt@[0..\#PE - 1])$   $\triangleright O(\log p)$ 
10:   $P \leftarrow$  the  $(pre\_pcnt[id'] - id + 1)$ -th portion of pairs from
     $H[\lceil \frac{k}{\#PE} \rceil \cdot id']$  to  $H[\min(\lceil \frac{k}{\#PE} \rceil \cdot (id + 1) - 1, k - 1)]$   $\triangleright O(\frac{n}{p})$ 
11:   $r \leftarrow \mathbf{ShortestDistance}(P)$   $\triangleright O(\frac{n}{p})$ 
12:   $result \leftarrow \mathbf{ReduceMin}(r@id)$   $\triangleright O(\log p)$ 
13: end procedure

```

0 Cell set  $H = \emptyset$ ,  $result = \infty$  **##  $H$  concurrent Hashmap**

1 Choose set  $S$  of  $n^{\frac{2}{3}}$  samples **## allocate  $p$  PEs**

$\delta = \min_{x,y \in S, x \neq y} (dist(x, y))$  **## recurse once**

2 **for**  $i = 1$  **to**  $n$  **do** **## allocate  $p$  PEs**

$c = cell(p_i)$

**if not**  $H.findCell(c)$  **then**

$H.addCell(\{c, idx = \{c.x, c.y\}\})$  **## maximal  $n$  times**

$H.findCell(c).addPoint(p_i)$

3 **foreach**  $c_1$  **in**  $H$  **do**

**foreach**  $c_2$  **in**  $\{c_1\} \cup neighbor(c_1)$  **do**

**foreach**  $(p_i, p_j)$  **in**  $c_1 \times c_2$

$result = \min(result, dist(p_i, p_j))$

$$O\left(\frac{n}{p} + \log p\right)$$

$$O\left(\frac{n^{\frac{8}{9}}}{p} + \log p\right)$$

$$O\left(\frac{n}{p}\right)$$

$$O\left(\frac{n}{p} + \log p\right)$$

Total runtime:  $O\left(\frac{n}{p} + \log p\right)$

- [BS76] Jon Louis Bentley and Michael Ian Shamos.  
“Divide-and-conquer in multidimensional space”. In:  
*Proceedings of the eighth annual ACM symposium on Theory of computing*. STOC '76. Hershey, Pennsylvania, USA: ACM, 1976, pp. 220–230. DOI: 10.1145/800113.803652. URL: <http://doi.acm.org/10.1145/800113.803652>.
- [CG89] Benny Chor and Oded Goldreich. “On the power of two-point based sampling”. In: *Journal of Complexity* 5.1 (1989), pp. 96–106.
- [Die+97] Martin Dietzfelbinger et al. “A reliable randomized algorithm for the closest-pair problem”. In: *Journal of Algorithms* 25.1 (1997), pp. 19–51.

## Reference II

- [FH78] Steven Fortune and John E Hopcroft. *A note on Rabin's nearest-neighbor algorithm*. Tech. rep. Cornell University, 1978.
- [Goe03] Brian Goetz. *Java theory and practice: Building a better HashMap*. 2003. URL: <http://www.ibm.com/developerworks/java/library/j-jtp08223/index.html/>.
- [KM95] Samir Khuller and Yossi Matias. “A simple randomized sieve algorithm for the closest-pair problem”. In: *Information and Computation* 118.1 (1995), pp. 34–37.
- [KT06] J. Kleinberg and E. Tardos. “Algorithm Design”. In: Pearson Education, 2006. Chap. 13 Randomized Algorithms.
- [Lip09] Richard J. Lipton. *Rabin Flips a Coin*. 2009. URL: <http://rjlipton.wordpress.com/2009/03/01/rabin-flips-a-coin/> (visited on 07/15/2013).

- [Rab76] Michael Oser Rabin. “Probabilistic algorithms”. In: *Algorithms and Complexity: New Directions and Recent Results*. Ed. by Joseph Frederick Traub. Academic Press, 1976, pp. 21–39.
- [Weba] URL: [http://www.cse.ust.hk/tcsc/comp670r/Class\\_1\\_Notes.ppt/](http://www.cse.ust.hk/tcsc/comp670r/Class_1_Notes.ppt/).
- [Webb] *Closest Pair of Points Problem*. URL: [http://en.wikipedia.org/wiki/Closest\\_pair\\_of\\_points\\_problem/](http://en.wikipedia.org/wiki/Closest_pair_of_points_problem/).

## Theorem

Let  $G_i$  with  $1 \leq i \leq k$  be the ***i-th*** non-empty cell and  $c, d$  some positive constants, then

$$\text{Prob} \left( \sum_{i=1}^k |G_i|^2 \leq c \cdot n \right) \geq 1 - \frac{1}{2^{n^d}}$$

## Sampling Lemma

Let  $G = \{G_1, G_2, \dots, G_k\}$  be a partition of set  $P$ , for which  $N(D) \geq n$ , where

$$N(G) = \sum_{i=1}^k \frac{|G_i| \cdot (|G_i| - 1)}{2}$$

If  $n^{\frac{2}{3}}$  pairwise different elements are drawn at random from  $P$  then the probability, that two elements will be chosen from the same  $G_i$ , is at least  $1 - 2^{-n^c}$  for some positive constant  $c$

Proof see

- **Rabin 1976**
- **Dietzfelbinger et al. 1997**

Estimate the probability using Chebyshev's inequality

## Lemma 2

Let  $G_\delta$  be a grid with gap  $\delta$  and  $G_{2\delta}$  another grid with gap  $2\delta$ , which is obtained by ignoring every second line of  $G_\delta$ , then

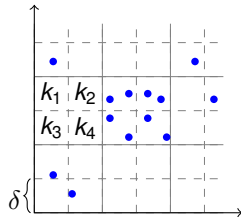
$$N(G_{2\delta}) \leq 4N(G_\delta) + \frac{3}{2}n$$

Let  $k = \sum_{i=1}^4 k_i$

$f(x) = x(x-1)$  convex  $\Rightarrow f(\frac{1}{4}k) \leq \frac{1}{4} \sum_{i=1}^4 f(k_i)$

$$\begin{aligned} \frac{1}{2}k(k-1) &= 8 \cdot \frac{1}{4}k(\frac{1}{4}k-1) + \frac{3}{2}k \\ &\leq 8 \cdot \frac{1}{4} \sum_{i=1}^4 k_i * (k_i-1) + \frac{3}{2}k \end{aligned}$$

$$\Rightarrow N(G_{2\delta}) \leq 8 \cdot \frac{1}{2}N(G_\delta) + \frac{3}{2}n$$



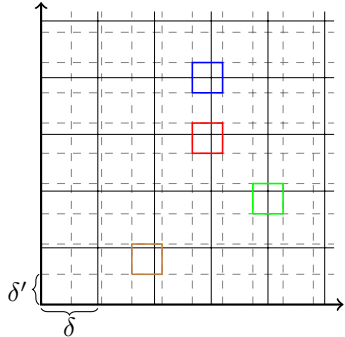


## Lemma 3

For any grid  $G_\delta$ ,  $G_{\delta'}$  with  $\delta' \leq \delta$  the following applies

$$N(G_{\delta'}) \leq 16N(G_\delta) + 6n$$

$$\begin{aligned} N(G_{\delta'}) &\leq \sum_{i=1}^4 N(G_{2\delta}^i) \\ &\leq \sum_{i=1}^4 (4N(G_\delta) + \frac{3}{2}n) \\ &= 16N(G_\delta) + 6n \end{aligned}$$

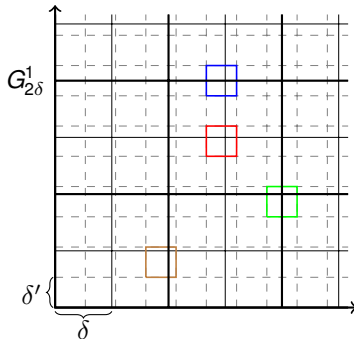


## Lemma 3

For any grid  $G_\delta$ ,  $G_{\delta'}$  with  $\delta' \leq \delta$  the following applies

$$N(G_{\delta'}) \leq 16N(G_\delta) + 6n$$

$$\begin{aligned}
 N(G_{\delta'}) &\leq \sum_{i=1}^4 N(G_{2\delta}^i) \\
 &\leq \sum_{i=1}^4 (4N(G_\delta) + \frac{3}{2}n) \\
 &= 16N(G_\delta) + 6n
 \end{aligned}$$

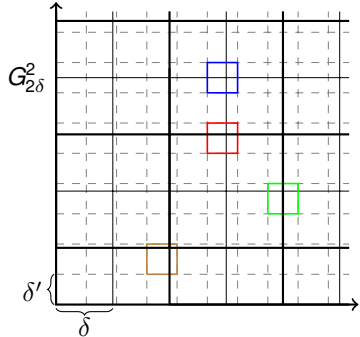


## Lemma 3

For any grid  $G_\delta$ ,  $G_{\delta'}$  with  $\delta' \leq \delta$  the following applies

$$N(G_{\delta'}) \leq 16N(G_\delta) + 6n$$

$$\begin{aligned}
 N(G_{\delta'}) &\leq \sum_{i=1}^4 N(G_{2\delta}^i) \\
 &\leq \sum_{i=1}^4 (4N(G_\delta) + \frac{3}{2}n) \\
 &= 16N(G_\delta) + 6n
 \end{aligned}$$

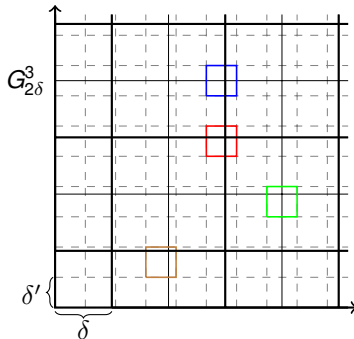


## Lemma 3

For any grid  $G_\delta$ ,  $G_{\delta'}$  with  $\delta' \leq \delta$  the following applies

$$N(G_{\delta'}) \leq 16N(G_\delta) + 6n$$

$$\begin{aligned}
 N(G_{\delta'}) &\leq \sum_{i=1}^4 N(G_{2\delta}^i) \\
 &\leq \sum_{i=1}^4 (4N(G_\delta) + \frac{3}{2}n) \\
 &= 16N(G_\delta) + 6n
 \end{aligned}$$

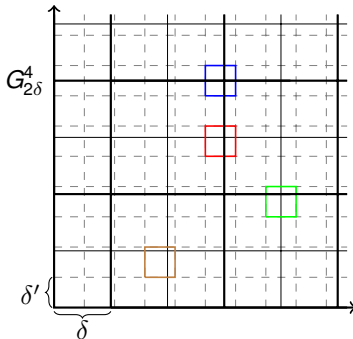


## Lemma 3

For any grid  $G_\delta$ ,  $G_{\delta'}$  with  $\delta' \leq \delta$  the following applies

$$N(G_{\delta'}) \leq 16N(G_\delta) + 6n$$

$$\begin{aligned}
 N(G_{\delta'}) &\leq \sum_{i=1}^4 N(G_{2\delta}^i) \\
 &\leq \sum_{i=1}^4 (4N(G_\delta) + \frac{3}{2}n) \\
 &= 16N(G_\delta) + 6n
 \end{aligned}$$



- Two of  $n^{\frac{2}{3}}$  samples will be very likely in the same cell, if  $N(G) \geq n$
- $N(G_{2\delta}) \leq 4N(G_\delta) + \frac{3}{2}n$
- $N(G_{\delta'}) \leq 16N(G_\delta) + 6n$  if  $\delta' \leq \delta$

## Proof.

Let  $\delta^*$  be the grid gap, which makes  $n \leq N(G_{\delta^*}) < 5.5n$ . (It exists!)  
If  $n^{\frac{2}{3}}$  samples are randomly chosen, two of them will be in one cell with high probability. Let  $\delta$  be the distance between them with  $\delta < 2\delta^*$ , then

$$\begin{aligned} N(G_\delta) &\leq 16N(G_{2\delta^*}) + 6n \\ &\leq 16(4N(G_{\delta^*}) + \frac{3}{2}n) + 6n \in O(n) \end{aligned}$$

