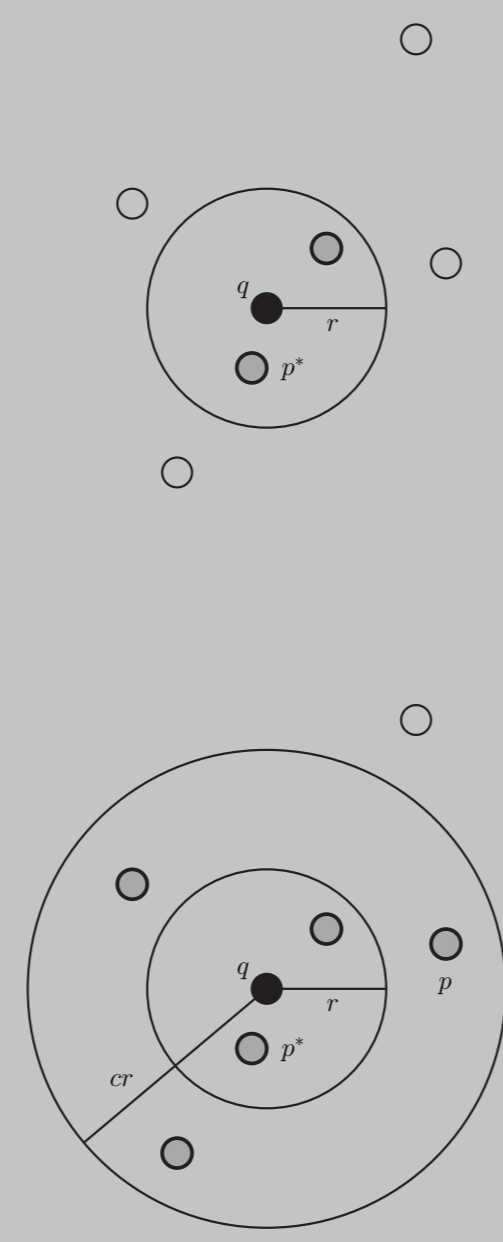


Beyond Locality-Sensitive Hashing

Alexandr Andoni (Microsoft Research), Piotr Indyk (MIT),
Huy L. Nguyen (Princeton) and Ilya Razenshteyn (MIT)

(Approximate) Near Neighbor Problem ((A)NN)

- ▶ **r-Near Neighbor Problem (NN)**
 - ▷ P — n -point subset of \mathbb{R}^d , $r > 0$
 - ▷ Given $q \in \mathbb{R}^d$ report any $p \in P$ s.t. $\|p - q\| \leq r$
 - ▷ Hard, if d is “large”
- ▶ **(c, r)-Approximate Near Neighbor Problem (ANN)**
 - ▷ In addition, we are given $c > 1$
 - ▷ Given $q \in \mathbb{R}^d$ report $p \in P$ s.t. $\|q - p\| \leq cr$, if there exists $p^* \in P$ s.t. $\|q - p^*\| \leq r$



Applications

- ▶ Pattern recognition, statistical classification, computer vision, computational geometry, databases, recommendation systems, DNA sequencing, spell checking, plagiarism detection, clustering etc
- ▶ **Approximate string matching**
 - ▷ Text T , query S : substring of T closest to S
 - ▷ Know $|S|$ in advance: can reduce to NN wrt Hamming distance
- ▶ **ML: nearest neighbor rule**
 - ▷ Use the label of the closest labeled example
- ▶ **Near-duplicate document retrieval**
 - ▷ Bag of words

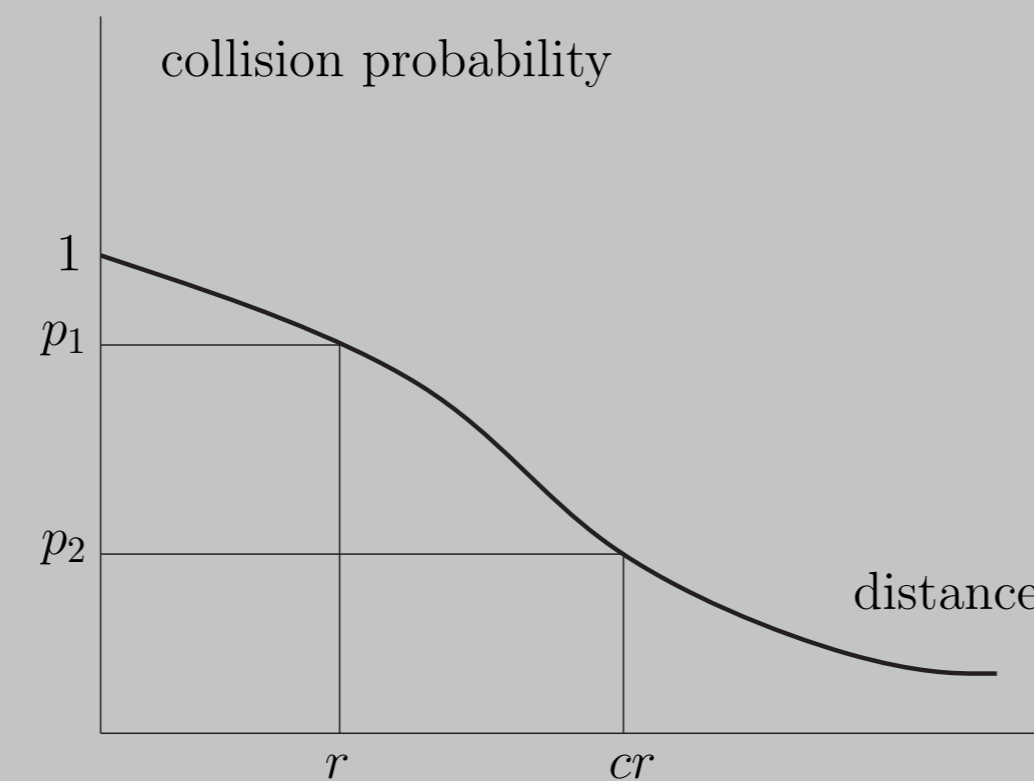
$T = \text{GAGTAACTCAATA}$
 $S = \text{AGTA}$



Document	Representation
In the beginning God created the heaven and the earth.	beginning 1
And the earth was without form, and void; and darkness was upon the face of the deep.	earth 2
And the Spirit of God moved upon the face of the waters.	God 3
And God said, Let there be light: and there was light.	

Locality-Sensitive Hashing (LSH)

- ▶ Introduced in (Indyk, Motwani 1998)
- ▶ A hash family \mathcal{H} of functions $h: \mathbb{R}^d \rightarrow \mathcal{U}$ s.t. for every $p, q \in \mathbb{R}^d$
 - ▷ if $\|p - q\| \leq r$, then $\Pr_{h \sim \mathcal{H}} [h(p) = h(q)] \geq p_1$ (“large”)
 - ▷ if $\|p - q\| \geq cr$, then $\Pr_{h \sim \mathcal{H}} [h(p) = h(q)] \leq p_2$ (“small”)
- ▶ Call such a family (r, cr, p_1, p_2) -sensitive
- ▶ Close points collide often



From LSH to ANN

- ▶ \mathcal{H} is a “reasonable” (r, cr, p_1, p_2) -sensitive family
- ▶ “Quality” of \mathcal{H}

$$\rho = \frac{\ln(1/p_1)}{\ln(1/p_2)}$$
- ▶ Can solve ANN in $O(n^{1+\rho} + d \cdot n)$ space and $O(d \cdot n^\rho)$ query time (Indyk, Motwani 1998)
- ▶ Example: $\{0, 1\}^d$, Hamming distance;
 - ▷ $\mathcal{H} = \{h_1, \dots, h_d\}$, $h_i(x) = x_i$
 - ▷ $\rho \leq 1/c$
- ▶ The best ρ for ℓ_p metric ($\|x - y\|_p = (\sum_i |x_i - y_i|^p)^{1/p}$) for $1 \leq p \leq 2$ is $\rho = 1/c^p + o(1)$ (matching lower bounds) (Indyk, Motwani 1998), (Andoni, Indyk 2006), (Motwani, Naor, Panigrahy 2007), (O’Donnell, Wu, Zhou 2011)

The main result

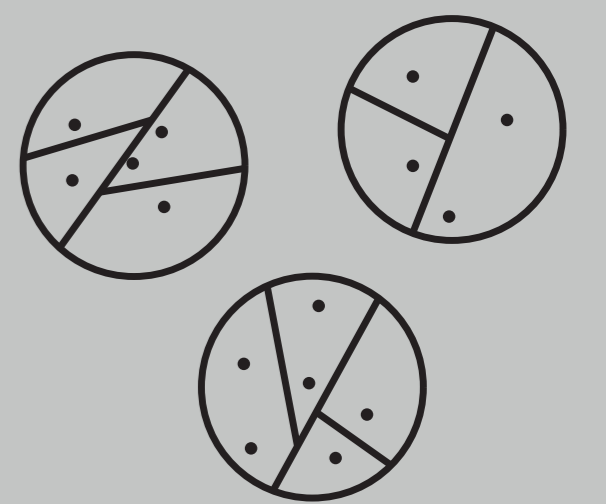
- ▶ A data structure for the (c, r) -ANN with space $O(n^{1+\rho} + d \cdot n)$ and query time $O(d \cdot n^\rho)$, where
$$\rho \leq \frac{0.73}{c^p} + O\left(\frac{1}{c^{3p/2}}\right) + o(1).$$
- ▶ Go beyond the best possible ρ for “vanilla” Locality-Sensitive Hashing
- ▶ The first improvement upon (Indyk, Motwani 1998) for ℓ_1 and (Andoni, Indyk 2006) for ℓ_2

The main approach: data-dependent hashing

- ▶ LSH is oblivious, try to be data-dependent in spirit of (Fredman, Komlós, Szemerédi 1984)
- ▶ **The definition of LSH is too strong!** Enough to satisfy the inequalities for $p \in P$ and $q \in \mathbb{R}^d$.
 - ▷ Exploit the geometry of P ?
 - ▷ **That is what we do!**
- ▶ Parallels with practice
 - ▷ PCA trees (Sproull 1991), (McNames 2001), (Verma, Kpotufe, Dasgupta 2009)
 - ▷ Spectral Hashing (Weiss, Torralba, Fergus 2008)
 - ▷ Semantic Hashing (Salakhutdinov, Hinton 2009)
 - ▷ WTA Hashing (Yagnik, Strelow, Ross, Lin 2011)

The two main steps

- ▶ Cluster P into low-diameter pieces (of diameter $O(cr)$)
- ▶ Handle each piece using a new “low-diameter” family

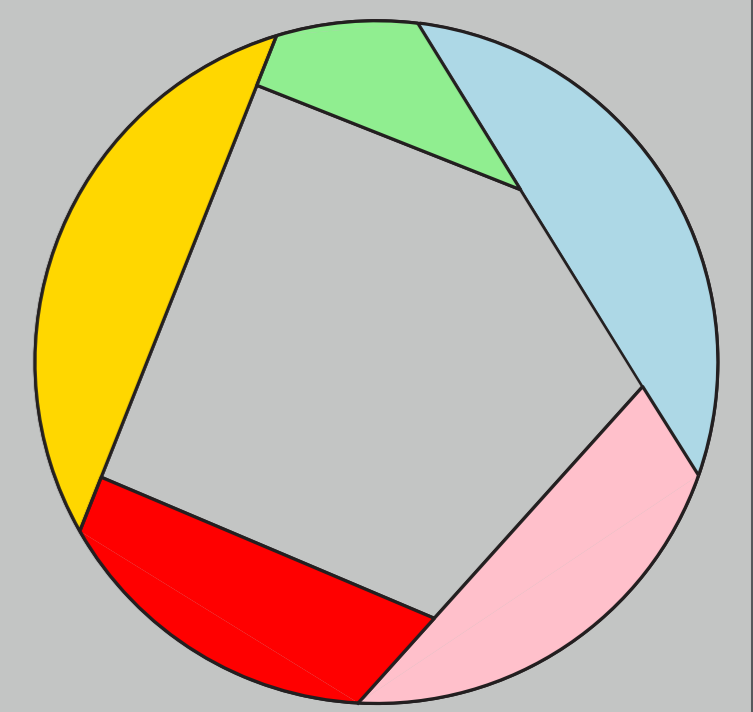


The low-diameter case

- ▶ Euclidean case: the previous record is $\rho \leq 1/c^2$
- ▶ Points and queries are on a sphere of radius $O(cr)$
- ▶ Can achieve

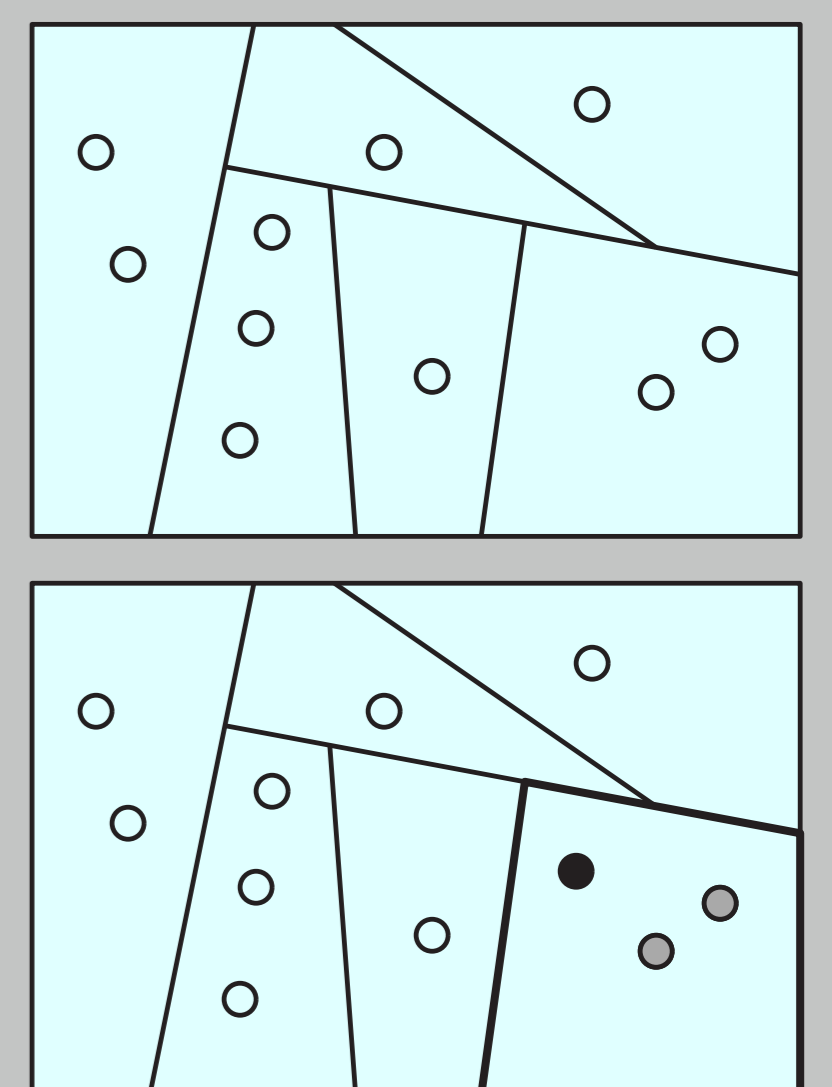
$$\rho = \frac{\ln(1/p_1)}{\ln(1/p_2)} \leq \frac{1 - \Omega(1)}{c^2}$$

using “ball carving”
(similar to (Karger, Motwani, Sudan 1998))



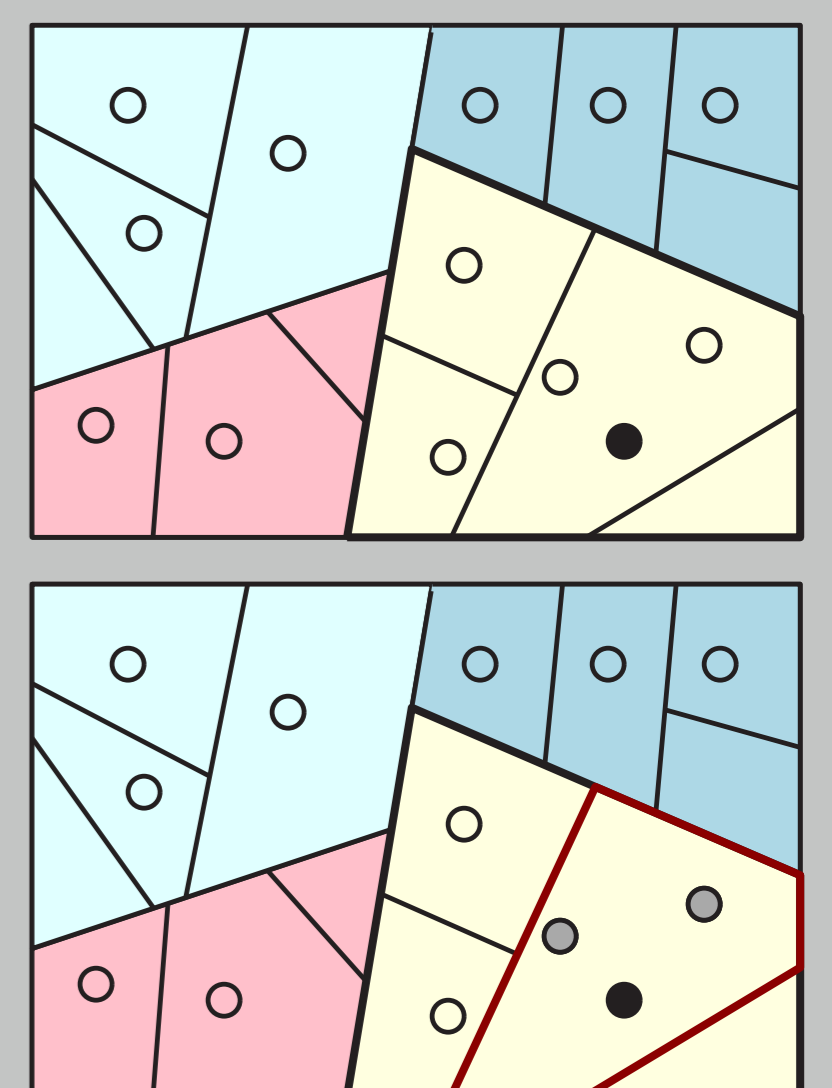
Vanilla LSH

- ▶ \mathcal{H} is a (r, cr, p_1, p_2) -sensitive family
- ▶ Hash P using a concatenation of k functions from $\mathcal{H}: x \mapsto (h_1(x), h_2(x), \dots, h_k(x))$
- ▶ Locate a query q , enumerate points from the bucket
- ▶ If choose k carefully: need n^ρ hash tables
- ▶ Overall: $n^{1+\rho}$ space, n^ρ query time



Two-level hashing

- ▶ Coarser partition (smaller k)
- ▶ Every part has a low diameter (aim at $O(cr)$)
- ▶ The “low-diameter” family to partition space finer
- ▶ “Outer” (data-independent, $\rho \leq 1/c^2$) + “inner” (data-dependent, “low-diameter” family) hash tables
- ▶ Get $\rho \leq (1 - \Omega(1))/c^2$: in inner tables better relation between p_1 and p_2 !



Conclusions

- ▶ Overcome limitations of LSH by using data-dependency
- ▶ Open problem: improve ρ
 - ▷ For a random instance from (Motwani, Naor, Panigrahy 2007) achieve $\rho \lesssim 1/(2c^2)$
 - ▷ Achieve this ρ for the general case?
 - ▷ Do even better for that random instance?
- ▶ Lower bounds?