

# Cheaper by the Dozen: Batched Algorithms\*

*Ben Gum*<sup>†</sup> and *Richard Lipton*<sup>‡</sup>

## 1 Introduction

While computing power and memory size have been steadily increasing as predicted by Moore's Law, they are still dwarfed by the size of massive data sets resultant from a number of applications. Many problems arising from astrophysics, computational biology, telecommunications, and the Internet often have an amount of accompanying data in the terabyte range. The analysis of this data by classical algorithms is often prohibitively expensive. Thus new ideas are necessary to create algorithms to deal with these massive data sets.

In this paper we develop the idea of *batching*, processing several queries at a time, for more efficient algorithms for several query problems. The advantages of our algorithms, over the classical approach of putting the massive dataset into a data structure, are threefold: improved asymptotic performance, significantly smaller data structures, and a number of I/O's which is linear in the size of the massive dataset. We use two techniques, query data structures and sampling, in the design of our batched algorithms. In addition, we believe that batched algorithms have many practical implications. Consider a webpage that answers queries on a large data set. Instead of answering these queries one at a time, which can result in a substantial bottleneck, we wait for several queries to accumulate, and then apply a batched algorithm that can answer them significantly faster.

To illustrate the idea of batched algorithms, we consider the dictionary problem. Suppose we begin with  $n$  unsorted items. If we have only one query, it does not make sense to place the  $n$  items in a data structure; the best we can do is the brute force method of comparing the query with all  $n$  items. Now suppose we have  $b$  queries. If  $b$  is large enough and we have enough space, it makes sense to build a data structure such as a binary tree or perfect hash table.

However, if  $1 < b \ll n$ , we can do better. We simply sort the list of the  $b$

---

\*Ben Gum's work was supported in part by an NSF Graduate Fellowship

<sup>†</sup>Princeton University, Department of Computer Science, gum@cs.princeton.edu

<sup>‡</sup>Georgia Institute of Technology, College of Computing, rjl@cc.gatech.edu

queries and for each dictionary item, perform binary search to see if the dictionary item matches a query. Thus we can check the  $b$  queries in  $O(n \log b + b \log b)$  time and using only  $O(b)$  storage space. Compare this to  $O(bn)$  time for the case of no preprocessing, and  $O(n \log n + b \log n)$  time and  $O(n)$  space, for the case of a binary search tree. For this problem, and a number of others we will show, batched algorithms 'fill in the blank' between the brute force approach for one query, and the classical data structure approach for  $\Omega(n)$  queries.

In the sections that follow, we show how batched algorithms can be applied to other problems. In section 2, we compare the idea of batching to other ideas in algorithms, including external memory algorithms and stream computing. We give an  $O(nbd^{0.3})$  time algorithm for the problem of  $d$ -dimensional Hamming nearest neighbor in section 3. In section 4, we show how query data structures can give us batched algorithms which answer  $b$  queries in  $O(n \log b)$  time for 1D Nearest Neighbor, String Matching, 2D Range Searching, 2D Nearest Neighbor, and 2D In hull. In section 5, by adding the idea of sampling to the query data structure, we give  $O(n \log b)$  time batched algorithms for Selection and Halfplane Intersection. Without the sampling phase, the algorithms would be I/O intensive and require the maintenance of a size- $O(n)$  data structure. Finally, in section 6, we discuss some directions for future work.

## 2 Related Work

We get the name *batching* from cryptography. There are several works which apply the idea of batching to problems in cryptography, including signature verification [2], exponentiation [10], and RSA [4]. In these works, batched algorithms are proposed to solve  $b$  related problems with less than  $b$  times the calculation necessary for solving one. A similar idea, called *Mass Production*, is mentioned in [12] which gives a more efficient algorithm for evaluating circuits.

Our batched algorithms on classical problems bear similarity to those of Karp, Motwani and Raghavan in [8]. In the cases of in hull and halfplane intersection, we derived our algorithms from their "Deferred Data Structure" algorithms. However, while their objective is that queries can appear online, we allow ourselves to see and manipulate the set of queries offline, but constrain ourselves to  $o(n)$  extra storage (where  $n$  is the size of the massive data set) and at most a constant number of passes over the massive data set.

Algorithms on large datasets are discussed in detail in the literature on External Memory Algorithms [11]. However most of the External Memory Algorithms are based on either previously sorted data, or on creating a data structure for the external data. We go in a different direction in this paper. We consider the case when placing our massive dataset in a data structure is either impossible due to its sheer size, or uneconomical because our query set is significantly smaller.

Stream computing [6] is another algorithmic idea related to our work. Stream computing algorithms try to compute functions of large datasets with just one pass through the data. In our work we also try to minimize the number of I/O's, but we allow ourselves an additional *sampling pass* in some cases. In addition, we focus on

algorithms for answering queries,  $b$  at a time, on the massive datasets.

It is interesting to note that batched algorithms can be thought of as the opposite of parallel algorithms. A batched algorithm answers several queries on one machine, while a parallel algorithm solves one problem on several machines. A batched algorithm for  $b$  queries will take at most  $b$  times as much computation as an algorithm for a single query, whereas a parallel algorithm run on  $m$  machines can achieve at most a factor of  $m$  savings. In fact, in the batched algorithms that we describe in future section, we will be able to answer  $b$  queries in only  $\log b$  times the computation required for a single query.

### 3 Hamming Nearest Neighbor

We consider the problem of  $d$ -dimensional nearest neighbor on an unorganized data set of size  $n$ . In many cases the best known approach for answering a query is simply the brute force method, of comparing it with every data item, which takes  $O(nd)$  time.

Notice however that the Hamming distance between two vectors  $\vec{a}$  and  $\vec{b}$  is  $\vec{a} \cdot \vec{a} + \vec{b} \cdot \vec{b} - 2\vec{a} \cdot \vec{b}$ . Using this fact and fast matrix multiplication we can create a batched algorithm, which multiplies two  $d \times d$  matrices in  $O(d^{2.3})$  time, to answer  $b$   $d$ -dimensional Hamming nearest neighbor queries in  $O(nbd^{0.3})$ , where  $b \geq d$ . To see this, partition the queries and the data items into sets of size  $d$ , and then perform  $(n/d)(b/d) = nb/d^2$   $d \times d$  matrix multiplications.

This relatively simple algorithm illustrates the power of batched algorithms for difficult problems on massive data sets. We believe similar ideas can be used to create algorithms for different nearest neighbor problems. To further explore the power of batching, in the next sections we give batched algorithms for a number of classical problems.

## 4 Query Data Structures

In this section we present batched algorithms for several problems which rely on the technique of creating a query data structure with the  $b$  query items and using it to process the  $n$  data items. These batched algorithms typically take  $\text{poly}(b)$  time to build the query data structure,  $O(n \log b)$  time to process the data, and  $\text{poly}(b)$  additional time to process information gathered from the massive data set. We chose  $b$  small enough so that the entire runtime is  $O(n \log b)$ . Additional advantages are that these algorithms require just  $\text{poly}(b)$  space, usually  $O(b)$  or  $O(b^2)$ , and just make one pass through the massive data set.

### 4.1 1D Nearest Neighbor

For this problem the query data structure is a small modification of the sorted list of queries. First we sort the queries to get the list  $q_1 \leq \dots \leq q_b$ . We also maintain the lists of variables  $l_1 \leq \dots \leq l_b$ , which are initially set to  $\infty$ , and  $r_1 \leq \dots \leq r_b$ , which are initially set to  $-\infty$ .

We then process the  $n$  unsorted data items as follows. For each data item  $d$

1. Do a binary search on the query set to find  $i$  such that  $q_i \leq d \leq q_{i+1}$ .
2. If  $d < r_i$  then set  $r_i = d$
3. If  $d > l_{i+1}$  then set  $l_{i+1} = d$

After all the  $n$  data items have been processed, we can examine the final  $l$ 's and  $r$ 's, in  $O(b)$  time, and find the nearest neighbor of each query item. This algorithm takes  $O(n \log b)$  time, since  $b \ll n$ . Also just  $O(b)$  space is required and each of the  $n$  unsorted items only needs to be examined once.

## 4.2 String Matching

We consider the standard string matching problem. We are given a  $0,1$  string of length  $n$ , and  $b$   $0,1$  query strings, each of length  $m$ , where  $mb \ll n$ . We use a modification of the Rabin-Karp fingerprinting algorithm [7]. Let our text be  $X = x_1x_2\dots x_n$  where  $x_1, x_2, \dots, x_n$  are bits. Let  $Y_1, Y_2, \dots, Y_b$  be the  $b$  query strings, each of length  $m$ . Let each length- $m$  substring of  $X$  be referred to as  $X(j) = x_{j+1}x_{j+2}\dots x_{j+m}$  for all  $0 \leq j \leq n - m$ .

The *fingerprinting* consists of looking at the residues of length- $m$  bit strings mod a random prime. Although fingerprinting introduces a small probability of a false match, it greatly reduces the time necessary to compare to check for matches. From now on we think of length- $m$  bit strings as  $m$ -bit integers. We randomly select a prime  $p$  from all the primes less than  $T$  (we will choose  $T$  later to make the probability of a mismatch small). Let the fingerprint of  $Z$ ,  $F_p(Z)$ , be  $Z \pmod{p}$ . Notice that  $F_p(X(j+1)) = [2F_p(X(j)) - 2^m x_{j+1} + x_{j+m+1}] \pmod{p}$ . So once  $F_p(X(j))$  is calculated,  $F_p(X(j+1))$  can be found with  $O(1)$  operations.

Suppose we have two  $m$ -bit integers  $Z_1 \neq Z_2$ . The probability that  $F_p(Z_1) = F_p(Z_2)$  is the same as the probability that  $p$  is a factor of  $Z_1 - Z_2$  which is at most  $m/\pi(T)$ , where  $\pi(T)$  is the number of primes less than  $T$ . (The  $m$  in the numerator is because the number of distinct prime factors of any positive integer  $k$  is at most  $\log_2 k$ ). Thus the probability that  $F_p(X(j)) = F_p(Y_i)$ , given  $X(j) \neq Y_i$ , for at least one  $0 \leq j \leq n - m$ ,  $1 \leq i \leq b$  is at most  $nbm/\pi(T) = O(1/nb)$  if we set  $T = n^2 b^2 m \log(n^2 bm)$ .

In this batched algorithm, our query data structure is simply a sorted list of the  $b$  fingerprints of the query strings. Thus for each  $X(j)$ , we can check in  $O(\log b)$  whether the fingerprint of  $X(j)$  is equal to the fingerprint of a query string. In total this takes a preprocessing time of  $O(mb + b \log b)$  to calculate the fingerprints of the query strings and sort those fingerprints. Since  $bm \ll n$ , the total runtime of the algorithm is  $O(n \log b)$  and the total storage used is  $O(b)$ . (Note that this is a Monte Carlo algorithm; it can be converted to a Las Vegas algorithm by verifying each match, and if there is a mismatch, then the rest is checked by the  $O(nbm)$  brute force algorithm).

### 4.3 2D Range Searching

Consider the standard 2D range query problem. We have an unorganized set of  $n$  points,  $S$ , in  $R^2$  and an axis parallel rectangle, defined by  $x_1 < x_2$  and  $y_1 < y_2$ . We would like to find out which points in  $S$  lie in the rectangle defined by  $(x_1, x_2, y_1, y_2)$ . For one query rectangle we need  $O(n)$  comparisons. We will give a batched algorithm for  $b$  query rectangles which runs in  $O(n \log b + k)$ , where  $k$  is the number of points reported, takes  $O(b^2 + k)$  space, and examines each of the  $n$  data points only once. For convenience, we assume that none of the data points lie on a line which defines the query rectangles. (If this is not the case, we can adapt our algorithm, but it is slightly cumbersome).

We place all the  $x_1$ 's and  $x_2$ 's in a  $O(\log b)$ -depth binary search tree. At each leaf of this tree we place a  $O(\log b)$ -depth binary search tree made of the  $y_1$ 's and  $y_2$ 's. This gives us a size  $O(b^2)$  binary tree of depth  $O(\log b)$  such that each leaf corresponds to a rectangle (possibly unbounded on up to 2 sides) in  $R^2$ . Each of the query rectangles can be partitioned into rectangles represented by the leaf nodes of this tree. For each query rectangle, we maintain a list (initially empty) of all the data points contained in it. All this takes  $O(b^2)$  time.

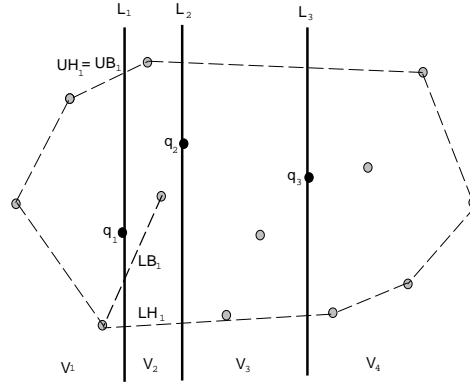
We process each data point in by finding the leaf rectangle that it lies in and adding it to the list of each query rectangle that the leaf rectangle lies in. It takes  $O(n \log b)$  time for all of the  $n$  data points to find the leaf rectangles that they lie in and  $k$  time to update the query rectangle lists. Thus for small enough  $b$  we see that the total runtime is  $O(n \log b + k)$ .

### 4.4 2D Nearest Neighbor

For each of  $b$  query points  $q_1, \dots, q_b$ , we would like to find the closest of  $n$  unorganized data points to it. The classical algorithm, builds a Voronoi diagram with the  $n$  points in  $O(n \log n)$  and then finds which region each of the  $b$  query points lies in. We would like to find a batched algorithm to answer  $b$  queries in  $O(n \log b)$ , without the cost and space requirements of creating a Voronoi Diagram with the  $n$  data points. However, it is not clear what a good query data structure would be.

A reasonable starting point is to put the  $b$  query points in a Voronoi Diagram. For each query point  $q_i$ , we find  $c_i$ , the closest of the of the  $n$  data points to  $q_i$  that lies in  $q_i$ 's Voronoi Region. It takes  $O(n \log b)$  to compute the  $c_i$ 's. Unfortunately, for any  $q_i$  there is no guarantee that  $c_i$  is  $q_i$ 's nearest neighbor, or that  $c_i$  even exists. However, for each  $q_i$  some  $c_j$  is a 3-approximate nearest neighbor. To see this let  $nn_i$  be the nearest neighbor of  $q_i$ . If  $nn_i$  is equal to  $c_i$  or some  $c_j$ , such that  $i \neq j$ , then we are done. Otherwise let  $nn_i$  lie in the Voronoi region of  $q_j$ . Since  $nn_i$  lies in the Voronoi region of  $q_j$ ,  $dist(q_i, nn_i) \geq dist(q_j, nn_i)$ . Also since  $c_j \neq nn_i$ ,  $dist(q_j, nn_i) \geq dist(q_j, c_j)$ . So  $dist(q_i, c_j) \leq dist(q_i, nn_i) + dist(nn_i, q_j) + dist(q_j, c_j) \leq 3 * dist(q_i, nn_i)$ .

Thus in  $O(n \log b)$  we can compute the 3-approximate 2D-nearest neighbor of  $b$  query points. It is easy to find an example that shows this algorithm is no better than 3-approximate. We can improve this to 2-approximate if we partition each Voronoi region into six  $60^\circ$  wedges around the  $q_i$  in that region and then save 6



**Figure 1.** *Batched Inhull Algorithm*

points (the closest point to  $q_i$  in each of the six  $60^\circ$  wedges in its Voronoi region). There is a simple example that shows this improved algorithm is no better than 2-approximate. However, we believe these approximate nearest neighbor algorithms are promising starting points for exact nearest neighbor algorithms.

## 4.5 2D Inhull

Using ideas similar to [8] we give a batched algorithm for deciding whether or not each of  $b$  query points lie inside the convex hull of a set  $P$  of  $n$  data points.

First we sort the  $b$  query points. Let  $q_i$  be the query point with the  $i$ th smallest  $x$  coordinate. Let  $L_i$  be the vertical line passing through  $q_i$ . The  $L_i$ 's partition the plane into  $b + 1$  vertical strips. Let  $V_i$  be the  $i$ th from the left.

For each  $L_i$  we would like to find the edges of the convex hull of  $P$  (if any) that intersect  $L_i$ . Let  $UH_i$  be the edge of the upper hull which intersects  $L_i$ , and let  $LH_i$  be the edge of the lower hull which intersects  $L_i$ . Clearly  $q_i$  lies in the convex hull of  $P$  if and only if  $LH_i$  and  $UH_i$  exist and  $q_i$  lies between them.

We can find all the  $UH_i$ 's and  $LH_i$ 's in  $O(n \log b)$  time as follows. We do a binary search for each data point  $p$  to find out which vertical strip it lies in (though that is not exactly what we are interested in), for each  $L_i$  that we compare  $p$  to, we also check to see if  $p$  lies in either  $UB_i$  or  $LB_i$ , the upper and lower edges, which intersect  $L_i$ , of the convex hull of the points of  $P$  which are compared to  $L_i$ . In the figure above,  $L_1$  is compared to all of the data points to the left of  $L_2$ .  $LB_1$  is the lower edge, which intersects  $L_1$ , of the convex hull of the points compared with  $L_1$ . We can find all of the  $LB$ 's and  $UB$ 's in  $O(n \log b)$  time. Since the  $LH$ 's are a subset of the  $LB$ 's (and the  $UH$ 's are a subset of the  $UB$ 's) we can then compute the  $UH_i$ 's and  $LH_i$ 's in  $O(b^2)$  additional time (in fact  $O(b \log b)$  time is sufficient). Once we have computed the  $UH_i$ 's and  $LH_i$ 's, we can then check if each of the query points lies in the hull on  $O(b)$  total time. Thus we have answered the  $b$  queries in  $O(n \log b)$  time.

## 5 Query Data Structures Plus Sampling

In this section we present batched algorithms for Selection and Halfplane Intersection. For these problems it is possible to use straightforward query data structures to answer  $b$  queries in  $O(n \log b)$  time. However these algorithms have two major disadvantages in that they use  $O(n)$  storage and require that each data item be examined  $\log b$  times. To remedy these drawbacks, we give sampling-based algorithms which, with high probability, solve the same problem in  $O(n \log b)$  time, but only require  $o(n)$  space and only look at each of the  $n$  unsorted items at most a constant number of times.

### 5.1 Batched Randomized Selection

In this section we give two algorithms for batched selection. The first algorithm uses a query data structure in a straightforward manner. The second algorithm combines sampling with a similar query data structure and answers the queries much more efficiently.

Recall the classical selection problem: given an unsorted dataset  $S$ , of size  $n$ , and an integer  $k$ , such that  $1 \leq k \leq n$ , return the  $k$ th smallest element of  $S$ . Blum et al [3] gave an  $O(n)$  deterministic algorithm for selection in the early 1970's and a number of people since then have given algorithms with a better coefficient. We define the batched version of selection as follows. Let  $S_{(i)}$  be the  $i$ th smallest elements in a dataset  $S$ . We are given queries:  $k_1, \dots, k_b$  and a dataset  $S$  and asked to find  $S_{(k_1)}, \dots, S_{(k_b)}$ . For notational convenience, assume  $k_1 < \dots < k_b$ . In fact, a special case of this problem, finding  $b$  quantiles of an unsorted dataset [1], has many applications.

We can use the algorithm by Blum et al in a divide and conquer fashion to find an  $O(n \log b)$  algorithm for batched selection. We select  $k_{\lceil b/2 \rceil}$  in  $O(n)$  time, use it to partition  $S$  in two parts, then recurse on each half. This algorithm corresponds to creating a partial binary search tree of  $S$ , where the  $b$  interior nodes are  $S_{(k_1)}, \dots, S_{(k_b)}$ . This takes time  $T(n, b) \leq O(n) + T(k_{\lceil b/2 \rceil}, b/2) + T(n - k_{\lceil b/2 \rceil}, b/2)$ , where  $T(n, 1) = O(n)$ . This recurrence is satisfied by  $T(n, b) = O(n \log b)$ .

However this algorithm has some serious drawbacks. It requires storage space for the size  $n$  data structure, which may be impossible, or prohibitively expensive. Even if there is sufficient space in external memory, each element of  $S$  must be moved around as many as  $O(\log b)$  times, which may be very expensive. Finally even if  $S$  fits in main memory, the algorithm either has a large leading coefficient of its runtime, or is very complicated to implement (depending on whether an old deterministic or a new new deterministic select algorithm is used).

To remedy these problems, we give an improved algorithm which uses sampling to create a much smaller, faster data structure. Our sampling-based algorithm, which uses ideas from the Floyd, Rivest randomized algorithm for selection, first presented in [5] and simplified in [9], can be used to select  $b$  items using  $n \log_2 \lceil 2b + 1 \rceil + o(n)$  comparisons, when  $b = O(n^{1/6})$ .

**Algorithm:**

1. Sample  $n^{5/6}$  elements from  $S$  with replacement, call the sample  $R$ . Sort  $R$ .

(Sampling with replacement is easier to do and does not hurt our analysis).

2. For each  $k_i$ , let  $y_i = \frac{k_i}{n}n^{5/6} = k_in^{-1/6}$  ( $y_i$  is the expected number of items in  $R$  which are less than  $S_{(k_i)}$ ).
3. Let  $x_i = R_{(y_i - \sqrt{n})}$ , if  $y_i > \sqrt{n}$  otherwise let  $x_i = S_{(1)}$ . Let  $z_i = R_{(y_i + \sqrt{n})}$ , if  $y_i < n - \sqrt{n}$  otherwise let  $z_i = S_{(n)}$ . (We will show that with high probability, for all  $i$ ,  $x_i \leq S_{(k_i)} \leq z_i$ ).
4. Sort the set containing all the  $x_i$ 's and  $z_i$ 's. The  $x_i$ 's and  $z_i$ 's partition the rest of the set  $S$  into at most  $2b + 1$  parts (according to which two an item of  $S$  lies between). Using binary search, find the place of each of the  $n$  elements of  $S$  among the  $x_i$ 's and  $z_i$ 's. (This binary search part is the most costly phase of the algorithm).
5. Let  $I_i = \{s \in S \mid x_i \leq s \leq z_i\}$ . Sort the elements in each interval  $I_i$ . (We will show that with high probability each interval  $I_i$  contains at most  $4n^{2/3}$  elements).
6. Now to find each  $S_{(k_i)}$ , if  $S_{(k_i)}$  lies in  $I_i$  then from step 4 we know how many elements of  $S$  that are less than  $S_{(k_i)}$  and from step 5 we have sorted  $I_i$  so we can quickly find  $S_{(k_i)}$ .

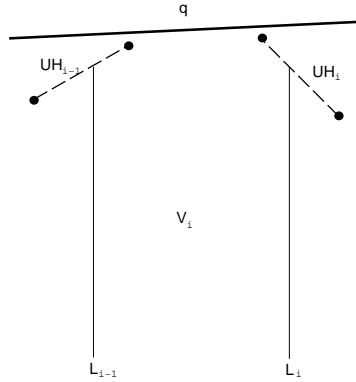
Through a straightforward use of Chernoff-Hoeffding bounds, we can show that with high probability the following conditions hold:

1. for all  $i$ ,  $S_{(k_i)}$  lies in  $I_i$
2. for all  $i$ ,  $I_i$  contains at most  $4n^{2/3}$  elements

We need the first condition so we can perform step 6 of our algorithm and find all the  $S_{(k_i)}$ 's. The second condition implies that step 5 of our algorithm takes  $o(n)$  time. The binary search phase of step 4 requires  $n \log_2 [2b + 1]$  comparisons. Since the other steps take  $o(n)$  comparisons, the total number of comparisons used by the algorithm is  $n \log_2 [2b + 1] + o(n)$ . Thus with high probability our algorithm correctly answers  $b$  select queries on an unsorted dataset of size  $n$  in  $O(n \log b)$  time, using  $o(n)$  extra storage, and examining each data item at most twice.

## 5.2 Halfplane Intersection

Finally we consider the question of whether a query point lies in the intersection of  $n$  unsorted halfplanes. As in our batched in hull algorithm, we use ideas from [8]. If the intersection is empty, the problem is trivial. Otherwise we can transpose the halfplanes in  $O(n)$  time so that they contain the origin. Then the problem dualizes to the problem of whether a query line does not intersect a convex hull. Let  $P$  be the set of points which are the duals of the  $n$  halfplanes and let  $q$  be a line which is the dual of one of the original query points. Observe that for any vertical line  $L$ , if we know the segments of the convex hull that intersect  $L$  then we can say one of four things about a query line  $q$ :



**Figure 2.** *Batched Half Plane Intersection Algorithm: If  $q$  intersects the convex hull, it does so in  $V_i$ .*

1.  $q$  does not intersect the convex hull
2.  $q$  does intersect the convex hull
3. if  $q$  intersects the convex hull, it does so to the left of  $L$
4. if  $q$  intersects the convex hull, it does so to the right of  $L$

Using this observation we can come up with an  $O(n \log b)$  batched algorithm for answering  $b$  queries. We first need to divide the plane up into  $b$  vertical strips each containing  $O(n/b)$  data points. This can be done deterministically using recursive applications of the linear-time median find algorithm. However, this would use a data structure of size  $O(n)$  and require each query to be examined  $\log b$  times.

We can instead do this using a sampling phase which requires only  $O(b^3)$  storage (which is advantageous provided  $b = o(n^{1/3})$ ). We, randomly with replacement, take a sample  $S$ , of  $b^3$  points, from  $P$ . Then we find the vertical lines  $L_1, \dots, L_{b-1}$  which partition the plane into  $b$  vertical strips, each containing  $b^2$  points of  $S$ . Using Chernoff-Hoeffding bounds we can show that with high probability, at most  $3n/b$  points of  $P$  lie in each vertical strip.

Now we find the edges of the convex hull of  $P$  which intersect the  $L_i$ 's, as in the batched in hull algorithm above, this step takes  $O(n \log b)$  time. Once we have these edges, for each query line  $q$ , by comparing  $q$  to the hull edges we have found we can say one of the following:

1.  $q$  intersects the convex hull of  $P$
2.  $q$  does not intersect the convex hull of  $P$
3. If  $q$  intersects the convex hull of  $P$  then it does so in the strip  $V_i$

In the last case, since each strip contains  $O(n/b)$  data points, we can check this for all such query lines in  $O(n \log b)$  time. Thus we have an  $O(n \log b)$ -time batched

algorithm for answering whether  $b$  query points lie in the intersection of  $n$  unsorted halfplanes.

## 6 Future Work

In the previous sections, we have explored the idea of batching and the power of its applications to algorithms for query problems. We gave a batched algorithm for answering  $b$  queries of  $d$ -dimensional Hamming nearest neighbor on a corpus of size  $n$  in  $O(nbd^{0.3})$  time, for  $b \geq d$ . In addition, using two techniques, query data structures and sampling, we have given batched algorithms for the classical problems of: string matching, range searching, selection, (1D and approximate 2D) nearest neighbor, in hull, and halfplane intersection. These batched algorithms answer  $b$  queries on a dataset of size  $n \gg b$  in  $O(n \log b)$  time using  $o(n)$  storage, and at most  $2n$  I/O's. In future work we hope to show the practicality of these algorithms through experiments.

There are numerous problems which would benefit from efficient batched algorithms. High dimensional nearest neighbor and computational biology are two very promising directions for work in batched algorithms. However we believe that nearly all query problems have the potential for efficient batched algorithms. We hope this paper will inspire future work in this area.

# Bibliography

- [1] Alsabti, K., Ranka, S., Singh, V. *A One-Pass Algorithm for Accurately Estimating Quantiles for Disk-Resident Data*, 23rd VLDB Conference. (1997).
- [2] Bellare, M., Garay, J., Rabin, T. *Fast Batch Verification for Modular Exponentiation and Digital Signatures*, Eurocrypt. (1998).
- [3] Blum, M., Floyd, R., Pratt, V., Rivest, R., Tarjan, R. *Time Bounds for Selection*, Journal of Computer and System Sciences, 7 (1973), 448–461.
- [4] Fiat, A. *Batch RSA* Journal of Cryptography, Vol 10, No 2 (1997), 75–88.
- [5] Floyd, R., Rivest, R. *Expected Time Bounds for Selection*, Communications of the ACM, 18 (1975), 165–177.
- [6] Henzinger, M., Raghavan, P., Rajagopalan, S. *Computing on Data Streams* SRC Technical Note, 1998-011. (1998).
- [7] Karp, R., Rabin, M. *Efficient randomized pattern-matching algorithms*, IBM Journal of Research and Development, 31 (1987), 249–260.
- [8] Richard M. Karp, Rajeev Motwani, and Prabhakar Raghavan. *Deferred data structuring*. SIAM Journal on Computing, 17(5):883-902, October 1988.
- [9] Motwani, R., Raghavan, P. *Randomized Algorithms*, Cambridge University Press. (1995).
- [10] M'Raihi, D., Naccache, D. *Batch exponentiation - A fast DLP based signature generation strategy* 3rd ACM Conference on Computer and Communications
- [11] Vitter, J. *External Memory Algorithms and Data Structures: Dealing with MASSIVE DATA* ACM Computing Surveys. (2000).
- [12] Wegener, I. *The Complexity of Boolean Functions* John Wiley and Sons. (1987).