

[Technical Report] Benchmarking Sketch-Based Data Stream Algorithms

[Experiments and Analyses]

Chenxingyu Zhao¹, Tong Yang¹, Yucheng Huang¹, Yinda Zhang¹
Qianhui Liu¹, Bin Cui¹, Tilman Wolf²

¹Department of Computer Science, Peking University

²Department of Electrical and Computer Engineering, University of Massachusetts Amherst

{dkzcx, lqianhui, bin.cui}@pku.edu.cn, wolf@umass.edu
{yangtongemail, huangyc96, hgdkgszyd}@gmail.com

1. PROGRAMMING MODEL

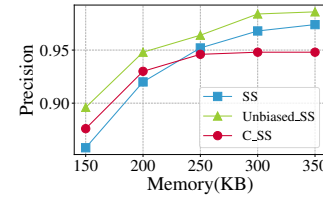
We describe the programming model using c++ style *class* shown as the code block *Class Sketch*. The class *Sketch()* has a member *Sketches* which is an probabilistic data structure (PDS) used by the algorithm. And the class has one *Update* function to update *Sketches* when a new item arrives, and a set of task-related *Query* functions to answer queries over data streams.

```
class Sketch()
{
    //probabilistic data structures
    PDS sketches;
    //customized parameters
    parameter phi, epsilon, ...;
    //other customized members
    hash_function Hash();
    hash_table HT;
    auxiliary_data_structure heap;
    .....
public:
    //To update sketches
    void Update(string Key, int V);
    //To report estimated frequency
    int frequencyQuery(string Key);
    //To report topk items
    vector< KV pair > topkQuery(int k);
    ...
    //other query functions
}
```

DSAB Programming Model for Users

2. COLDFILTER-SPACESAVING

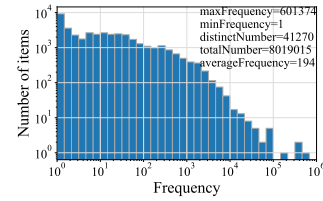
For the comparison between M-P algorithms and A-F algorithms, we compare SS and C-SS.



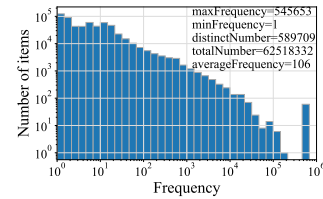
(a) IPtrace

Figure 11: SS augmented by Cold-Filter.

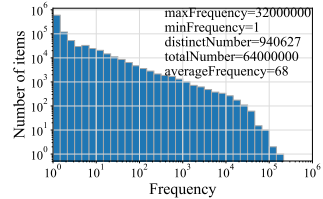
3. DATASET DISTRIBUTION



(a) kosarak



(b) IPtrace



(c) webdocs

Figure 11: Distribution of Realworld Datasets

4. PSUEDOCODE

Algorithm 1: CM

```

1 array of counters  $C[d][w]$ ; //initially 0
2  $d$  : the number of hash functions;
3  $w$  : the number of per-row counters;
4 Function Update( $Key, V$ ):
5   for  $i \leftarrow 1$  to  $d$  do
6      $C[i][h_i(Key)] \leftarrow C[i][h_i(Key)] + V$ ;
7   if have heap to maintain top-k then
8      $\hookrightarrow$  update heap

```

Algorithm 2: CM-CU

```

1 array of counters  $C[d][w]$ ; //initially 0
2  $d$  : the number of hash functions;
3  $w$  : the number of per-row counters;
4 Function Update( $Key, V$ ):
5    $minFreq \leftarrow \min_{1 \leq j \leq d} \{C[j][h_j(Key)]\}$ ;
6   for  $i \leftarrow 1$  to  $d$  do
7     if  $C[i][h_i(Key)] == minFreq$  then
8        $C[i][h_i(Key)] \leftarrow C[i][h_i(Key)] + V$ ;
9   if have heap to maintain top-k then
10     $\hookrightarrow$  update heap

```

Algorithm 3: COUNT

```

1 array of counters  $C[d][w]$ ; //initially 0
2  $w$  : the number of per-row counters;
3  $d$  hash functions :  $h_1, h_2, \dots, h_d$ ; //  $h_i : Key \rightarrow [1 : w]$ 
4  $d$  hash functions :  $g_1, g_2, \dots, g_d$ ; //  $g_i : Key \rightarrow +1, -1$ 
5 Function Update( $Key, V$ ):
6   for  $i \leftarrow 1$  to  $d$  do
7      $C[i][h_i(Key)] \leftarrow C[i][h_i(Key)] + V \times g_i(Key)$ ;
8   if have heap to maintain top-k then
9      $\hookrightarrow$  update heap

```

Algorithm 4: SPACE SAVING AND VARIANTS

```

1 Stream-Summary  $SS$ ; // list of ( $item, counter$ ) pairs
2  $c$  : the capacity of  $SS$ ;
3  $p$  : the probability of replacing minimal item;
4 Function Update( $Key, V$ ):
5   if  $item\ Key\ in\ SS$  then
6      $SS[Key].counter \leftarrow SS[Key].counter + V$ ;
7   else if  $SS.size < c$  then
8      $SS.add((Key, V))$ ;
9   else
10     $(Key_{min}, U_{min}) \leftarrow$ 
11    the pair with the minimal counter;
12     $Increment\ U_{min} \leftarrow U_{min} + V$ ;
13     $Replace\ key_{min} \leftarrow Key$  with the probability  $p$  ;

```

Algorithm 5: LOSSY-COUNTING

```

1 set of entries  $E$ ; // entry form ( $item, counter, \Delta$ )
2  $w$  : width of window;
3  $n$  : number of items (not distinct); //initially 0
4 Function Update( $Key, V$ ):
5    $n \leftarrow n + 1$ ;
6   if  $item\ Key\ in\ E$  then
7      $E[Key].counter \leftarrow E[Key].counter + V$ ;
8   else
9      $E \leftarrow E + (Key, 1, \lfloor \frac{n}{w} \rfloor)$ ;
10  if  $n \equiv 0 \bmod w$  then
11    for each entry  $e$  in  $E$  do
12      if  $e.counter \leq \frac{n}{w}$  then
13         $E \leftarrow E - e$ ;

```

Algorithm 6: HEAVYGUARDIAN

```

1 hash table  $HT$ ;
2 hash function  $h(\cdot), g(\cdot)$ ;
3  $HT[h(\cdot)].heavypart$  : array of KV pair ( $item, counter$ );
4  $HT[h(\cdot)].lightpart$  : array of counters;
5 Function Update( $Key, V$ ):
6   if  $item\ Key\ in\ HT[h(Key)].heavypart$  then
7      $\hookrightarrow$  update the KV pair of  $Key$ 
8   else if  $HT[h(Key)].heavypart$  has empty cell then
9      $HT[h(Key)] \leftarrow HT[h(Key)] + (Key, V)$ ;
10  else
11    Exponential Decay;
12    if  $HT[h(Key)].heavypart$  has zero-counter cell then
13       $HT[h(Key)] \leftarrow HT[h(Key)] + (Key, V)$ ;
14    else
15       $HT[h(Key)][g(Key)]_l \leftarrow$ 
16       $HT[h(Key)][g(Key)]_l + V$ ;

```

Algorithm 7: ASKETCH

```

1 filter  $F$ ; //array of ( $item, old\_count, new\_count$ )
2 classical sketch  $S_c$ ;
3 Function Update( $Key, V$ ):
4   if  $item\ Key\ in\ F$  then
5      $F[Key].new\_count \leftarrow F[Key].new\_count + V$ ;
6   else if  $F$  not full then
7      $F \leftarrow F + (Key, 0, V)$ ;
8   else
9      $S_c.insert(Key, V)$ ;
10     $V_{S_c} \leftarrow S_c.FrequencyQuery(Key)$  ;
11    if  $V_{S_c} > minF.new\_count$  then
12       $K_{min} \leftarrow$  minimum new_count item in  $F$ ;
13      if  $F[K_{min}].new\_count - F[K_{min}].old\_count > 0$ 
14      then
15         $S_c.insert(K_{min}, F[K_{min}].new\_count -$ 
16         $F[K_{min}].old\_count)$ ;
17         $F \leftarrow F - (K_{min}, old\_count, new\_count)$ ;
18         $F \leftarrow F + (Key, V_{S_c}, V_{S_c})$ ;

```

Algorithm 8: UNIVMON

```
1 d hash functions :  $h_1(\cdot), h_2(\cdot), \dots, h_d(\cdot)$ ;  
  //  $h_i(\cdot) : \text{item} \rightarrow \{0, 1\}$   
2 classical sketches  $S_c[d]$ ;  
3 Function Update( $Key, V$ ):  
4    $S_c[i].\text{insert}(Key, V)$ ;  
5   for  $i \leftarrow 2$  to  $d$  do  
6     if  $h_i(Key) == 1$  then  
7        $S_c[i].\text{insert}(Key, V)$ ;  
8     else  
9       break;
```

Algorithm 9: COLD-FILTER

```
1 classical sketch  $S_c$ ;  
2  $d_1$  hash functions :  $h_i(\cdot)$ ;  $d_2$  hash functions :  $g_j(\cdot)$ ;  
3 two layers of the filter:  $L_1[\omega_1], L_2[\omega_2]$ ;  
4  $\omega_1, \omega_2$  : number of counters;  
5 Function Update( $Key, V$ ):  
6   for  $i \leftarrow 1$  to  $d_1$  do  
7      $\text{update the corresponding counters } L_1[h_i(Key)]$   
8   if  $L_1[h_i(Key)]$  concurrently overflow then  
9     for  $j \leftarrow 1$  to  $d_2$  do  
10       $\text{update the corresponding counters } L_2[g_j(Key)]$   
11      if  $L_2[g_j(Key)]$  concurrently overflow then  
12         $S_c.\text{Update}(Key, V)$ ;
```

Algorithm 10: PYRAMID

```
1 classical sketch  $S_c$ ;  
2 pyramid layers of counters  $L[n_L]$ ;  
3  $n_L$  : number of layers in pyramid ;  
4 Function Update( $Key, V$ ):  
5    $S_c.\text{insert}(Key, V)$ ;  
6   Boolean carryin  $\leftarrow \text{False}$ ;  
7   if counter overflow happens during  $S_c.\text{Update}(Key, V)$   
8     then  
9        $\text{carryin} \leftarrow \text{True}$ ;  
10   $i \leftarrow 1$ ;  
11  while carryin do  
12     $L[i].\text{increment}$ ;  
13    if counter overflow happens during  $L[i].\text{increment}$   
14      then  
15         $\text{carryin} \leftarrow \text{True}$ ;  
16    else  
17       $\text{carryin} \leftarrow \text{False}$ ;
```
