

SODA: A Set of Fast Oblivious Algorithms in Distributed Secure Data Analytics

Xiang Li
Tsinghua University
lixiang20@mails.tsinghua.edu.cn

Yunqian Luo
Tsinghua University
luoyq19@mails.tsinghua.edu.cn

Nuozhou Sun
Tsinghua University
snz21@mails.tsinghua.edu.cn

Mingyu Gao
Tsinghua University
gaomy@tsinghua.edu.cn

1 OBLIVIOUS FILTER

Algorithm 1: Oblivious Global Filter

```

1  $\mathcal{B}_{all} \leftarrow [];$ 
2 for  $i \leftarrow 0$  to  $N_{in} - 1$  do
3    $\mathcal{B}_{all}[i] \leftarrow \text{OFilterStage1}(\mathcal{P}_{in}[i], N_{in})$ 
4 for  $i \leftarrow 0$  to  $N_{out} - 1$  do
5    $\mathcal{P}_{out}[j] \leftarrow \text{OFilterStage2}(\mathcal{B}_{all}[..][j])$ 

```

Algorithm 2: Oblivious Global Filter: Stage 1

Input: Local data partition \mathcal{P} , total number of partitions N

Output: N buckets \mathcal{B} .

```

1 foreach  $d \in \mathcal{P}$  do // randomly assign records to buckets.
2    $d.\text{bucket} \leftarrow \text{GetRandom}() \bmod N;$ 
3  $\mathcal{B} \leftarrow \text{BuildBuckets}(\mathcal{P}, N);$ 
4 return  $\mathcal{B};$ 

5 function  $\text{BuildBuckets}(\mathcal{P}, l, N):$ 
6    $\text{OSort}(\mathcal{P}, (\_.\text{bucket}, \_.\text{valid}, \_.\text{key}));$ 
7   /* Assign locations in buckets. */
8    $b \leftarrow \perp; p \leftarrow \perp;$ 
9   foreach  $d \in \mathcal{P}$  do
10     $c \leftarrow d.\text{bucket} \neq b;$  // encounter a new bucket?
11     $b \leftarrow d.\text{bucket};$ 
12     $\text{cmov}(c, p, l \times b);$  // starting location of a new bucket.
13     $d.\text{pos} \leftarrow p; p \leftarrow p + 1;$ 
14    /* Place records to assigned locations. */
15    Pad  $\mathcal{P}$  to length  $l \cdot N;$ 
16     $\text{ODistribute}(\mathcal{P}, \_.\text{pos});$  // invalid records at end of each bucket.
17     $\mathcal{B} \leftarrow \text{Chunk } \mathcal{P} \text{ into buckets of size } l;$ 
18    return  $\mathcal{B};$ 

```

Algorithm 3: Oblivious Global Filter: Stage 2

Input: Buckets \mathcal{B} collected from all nodes

Output: Partition \mathcal{P} to be processed later by this node

```

/* Note that each bucket is already sorted by  $k$  */
1  $\mathcal{D} \leftarrow \text{OMerge}(\mathcal{B}, (\neg \text{is\_valid}(\_.\text{v}), \_.\text{k}));$  // merge by  $k$ 
2  $\mathcal{P} \leftarrow$  Directly remove all dummy values at the end of  $\mathcal{D};$ 
3 return  $\mathcal{P};$ 

```

Algorithm 4: (Pseudo-)Random Communication Buckets

```

1 function  $\text{BuildBuckets}(\mathcal{P}, N):$ 
2   Input: Data partition  $\mathcal{P}$  with (pseudo-)randomly assigned
3     bucket field per record, total number of partitions  $N$ .
4   Output:  $N$  buckets  $\mathcal{B}$ .
5   /* In implementation, we fix the error rate  $Ne^{-k^2/3}$  to compute the pad
6     length */
7   Compute the padding size  $l \leftarrow \frac{|\mathcal{P}|}{N} + k\sqrt{\frac{|\mathcal{P}|}{N}};$ 
8    $\text{OSort}(\mathcal{P}, (\_.\text{bucket}, \_.\text{valid}, \_.\text{key}));$ 
9   /* Assign locations in buckets. */
10   $b \leftarrow \perp; p \leftarrow \perp;$ 
11  foreach  $d \in \mathcal{P}$  do
12     $c \leftarrow d.\text{bucket} \neq b;$  // encounter a new bucket?
13     $b \leftarrow d.\text{bucket};$ 
14     $\text{cmov}(c, p, l \times b);$  // starting location of a new bucket.
15     $d.\text{pos} \leftarrow p; p \leftarrow p + 1;$ 
16    /* Place records to assigned locations. */
17  Pad  $\mathcal{P}$  to length  $l \cdot N;$ 
18   $\text{ODistribute}(\mathcal{P}, \_.\text{pos});$  // invalid records at end of each bucket.
19   $\mathcal{B} \leftarrow \text{Chunk } \mathcal{P} \text{ into buckets of size } l;$ 
20  return  $\mathcal{B};$ 

```

2 OBLIVIOUS AGGREGATE

Algorithm 5: Oblivious Aggregate

```

/* In parallel, each node  $i$  processes its own partition. */
1  $\mathcal{B}_{all} \leftarrow []$ ;
2 for  $i \leftarrow 0$  to  $N_{in} - 1$  do
3    $\mathcal{B}_{all}[i] \leftarrow \text{OAggStage1}(\mathcal{F}, \mathcal{P}_{in}[i], i, \text{false})$ 

/* In parallel, each node  $j$  collects its partition through network and do
further processing. */
4 for  $j \leftarrow 0$  to  $N_{out} - 1$  do
5    $\mathcal{P}_{out}[j] \leftarrow \text{OAggStage2}(\mathcal{F}, \mathcal{B}_{all}[\dots][j], \text{true})$ 

```

Algorithm 6: Oblivious Aggregate: Stage 1

```

1 function OAggStage1( $\mathcal{F}, \mathcal{P}, i, N_{out}, \text{should\_assign\_loc}$ ):
   Input: Local data partition  $\mathcal{P}$ , total number of partitions  $N$ .
   Aggregate function  $\mathcal{F}$ 
   Output:  $N$  buckets  $\mathcal{B}$ .

2   OSort( $\mathcal{P}, \_key$ ) ; // sort by key.
3   if  $\text{should\_assign\_loc}$  then
4     /*  $\mathcal{P}$  has two more fields: location  $l$ , i.e., the current item index, and
       original partition id  $i$ . */
      $\mathcal{P} \leftarrow \text{AssignLoc}(\mathcal{P}, i)$ ;
5   else
6      $\mathcal{P} \leftarrow \mathcal{P}$ ;

   /* Perform aggregate and mark dummy items. */
7    $\mathcal{P} \leftarrow \text{Aggregate}(\mathcal{F}, \mathcal{P})$ ;
   /* Assign buckets. */
8    $\text{nonce} \leftarrow 0$ ;
9   foreach  $d \in \mathcal{P}$  do
10     $h \leftarrow \text{Hash}(d.key)$  ; // use real key.
11     $\text{cmov}(\neg d.valid, h, \text{Hash}(d.key \parallel \text{nonce}))$  ; // use dummy key.
12     $d.bucket \leftarrow h \bmod N$ ;
13     $\text{nonce} \leftarrow \text{nonce} + 1$ ;
14    $\mathcal{B} \leftarrow \text{BuildBuckets}(\mathcal{P}, N)$ ;
15   return  $\mathcal{B}$ ;

16 function Aggregate( $\mathcal{F}, \mathcal{D}$ ):
   Input: Aggregate function  $\mathcal{F}$ , data  $\mathcal{D}$ 
   Output: Data after aggregate  $\mathcal{D}$ 

   /* Do further aggregate */
17    $\text{acc} \leftarrow \mathcal{D}[0].v$ ;
18   for  $i \leftarrow 1$  to  $|\mathcal{D}| - 1$  do
19      $c \leftarrow \text{acc} \neq \perp \wedge \mathcal{D}[i].v \neq \perp$ ;
20      $\text{cmov}(\text{acc} = \perp, \text{acc}, \mathcal{D}[i].v)$ ;
21      $\text{cmov}(c, \text{acc}, \mathcal{F}(\text{acc}, \mathcal{D}[i].v))$ ;
22      $c \leftarrow \mathcal{D}[i-1].k \neq \mathcal{D}[i].k$ ;
23      $\text{cmov}(c, \text{acc}, \mathcal{D}[i].v)$ ;
24      $\text{cmov}(\neg c, \mathcal{D}[i-1].v, \perp)$ ;
25      $\mathcal{D}[i].v \leftarrow \text{acc}$ ;
26   return  $\mathcal{D}$ ;

```

Algorithm 7: Oblivious Aggregate: Stage 2

```

1 function OAggStage2( $\mathcal{F}, \mathcal{B}, \text{should\_remove\_dummy}$ ):
   Input: Aggregate function  $\mathcal{F}$ , buckets  $\mathcal{B}$  collected from all
   nodes
   Output: Partition  $\mathcal{P}$  to be processed later by this node

   /* Note that each bucket is already sorted by  $k$  */
2    $\mathcal{D} \leftarrow \text{OMerge}(\mathcal{B}, \_k, \text{is\_valid}(\_v))$  ; // merge by  $k$ , and invalid
   values are at the front of the group.

   /* Do further aggregate */
3    $\mathcal{D} \leftarrow \text{Aggregate}(\mathcal{F}, \mathcal{D})$ ;
4   if  $\text{should\_remove\_dummy}$  then
5     /* Remove all dummy values */
      $\mathcal{P} \leftarrow \text{OLocalFilter}(\mathcal{D})$ ;
6   else
7      $\mathcal{P} \leftarrow \mathcal{D}$ ;
8   return  $\mathcal{P}$ ;

```

3 OBLIVIOUS JOIN

Algorithm 8: Oblivious Join

```

1 Oblivious global filter  $\mathcal{P}_{A-in}$  and  $\mathcal{P}_{B-in}$ ; // shuffle.
  /* Perform oblivious aggregate, and the aggregate function is count. */
2  $\mathcal{B}_{A-agg}, \mathcal{B}_{B-agg} \leftarrow [], []$ ;
3 for  $i \leftarrow 0$  to  $N_{in} - 1$  do
4    $\mathcal{B}_{A-agg}[i], \pi_A[i] \leftarrow \text{OAggStage1}(\text{count}, \mathcal{P}_{A-in}[i], i, \text{true})$ ;
5    $\mathcal{B}_{B-agg}[i], \pi_B[i] \leftarrow \text{OAggStage1}(\text{count}, \mathcal{P}_{B-in}[i], i, \text{true})$ ;
  /* Note that  $\mathcal{P}_{A-in}[i]$  and  $\mathcal{P}_{B-in}[i]$  is sorted by  $k$  at this point. */
6  $\beta_A, \beta_B \leftarrow 0, 0$ ;
7 for  $j \leftarrow 0$  to  $N_{out} - 1$  do
8    $\mathcal{P}_{A-agg}[j] \leftarrow \text{OAggStage2}(\text{count}, \mathcal{B}_{A-agg}[\dots][j], \text{false})$ ;
9    $\mathcal{P}_{B-agg}[j] \leftarrow \text{OAggStage2}(\text{count}, \mathcal{B}_{B-agg}[\dots][j], \text{false})$ ;
  /*  $\mathcal{P}_{A-agg}$  and  $\mathcal{P}_{B-agg}$  has fields  $k, v, i, l$  */
10   $a_A, cnt_A \leftarrow \text{OJoinStage0}(\mathcal{P}_{A-agg}[j])$ ;
11   $a_B, cnt_B \leftarrow \text{OJoinStage0}(\mathcal{P}_{B-agg}[j])$ ;
  /* Need each node to broadcast their  $a_A, a_B, cnt_A$  and  $cnt_B$ . */
12   $\alpha_A, \alpha_B \leftarrow \max(\alpha_A, a_A), \max(\alpha_B, a_B)$ ;
13   $\beta_A, \beta_B \leftarrow \beta_A + cnt_A, \beta_B + cnt_B$ ;
14  $\mathcal{I} \leftarrow []$ ;
15 for  $j \leftarrow 0$  to  $N_{out} - 1$  do
16    $\text{bin\_num}, \text{bin\_size}, N_{out}', \text{acc\_prod}, \mathcal{P}_{agg}[j] \leftarrow$ 
      $\text{OJoinStage1}(\mathcal{P}_{A-agg}[j], \mathcal{P}_{B-agg}[j], \alpha_A, \alpha_B, \beta_A, \beta_B,$ 
      $N_{out}')$ ;
  /*  $\mathcal{P}_{agg}$  has fields:  $k, v, i, l, t, b$ . */
17   Append  $(\text{bin\_num}, \text{bin\_size}, \text{acc\_prod})$  to  $\mathcal{I}$ ;
18  $\mathcal{B}_{A-agg}, \mathcal{B}_{B-agg} \leftarrow [], []$ ;
19 for  $j \leftarrow 0$  to  $N_{out} - 1$  do
  /* Coordinate the last bin number and update all bin numbers */
20    $\mathcal{P}_{agg}[j], \mathcal{D}[j], \mathcal{S}[j], \mathcal{L}[j], \mathcal{S}_{rem}[j] \leftarrow$ 
      $\text{OJoinStage2}(\mathcal{P}_{agg}[j], [\alpha_A, \alpha_B], j, \mathcal{I}, N_{out}')$ ;
21 for  $j \leftarrow 0$  to  $N_{out} - 1$  do
  /* Aggregate column sum and adjust partition number by column. */
22    $\mathcal{L}_{prev} \leftarrow \perp$ ;
23   if  $j = 0$  then
24      $\mathcal{S}_g, \mathcal{S}_{rem} \leftarrow \text{AssignRemBin}(\mathcal{S}_{rem}, N_{out}')$ ; // collect  $\mathcal{S}_{rem}$ 
      from other nodes and broadcast the new one.
25    $\mathcal{S}_g, \mathcal{S}[j], \mathcal{L}[j] \leftarrow \text{OJoinStage3}(\mathcal{S}_g, \mathcal{S}[j], \mathcal{L}_{prev}, \mathcal{L}[j],$ 
      $\mathcal{S}_{rem}[j])$ ;
26    $\mathcal{L}_{prev} \leftarrow \mathcal{L}[j]$ ; // Send to node  $j + 1$ .
27   Send  $\mathcal{S}_g$  to node  $j + 1$ ;
28  $\gamma \leftarrow \sum_i \mathcal{S}_g[i].1$ ; // computed by node  $N_{out} - 1$  and broadcast.
29 for  $j \leftarrow 0$  to  $N_{out} - 1$  do
30    $\mathcal{B}_{A-agg}[j], \mathcal{B}_{B-agg}[j] \leftarrow \text{Split the buckets returned by}$ 
      $\text{OJoinStage4}(\mathcal{P}_{agg}[j], \mathcal{D}[j], \mathcal{S}[j], \mathcal{L}[j], \mathcal{S}_{rem}[j])$ ;
  /* Assign the destination partition number to items in input partitions, pad
  according to the public information, and build buckets. */
31  $\mathcal{B}_{all} \leftarrow []$ ;
32 for  $i \leftarrow 0$  to  $N_{in} - 1$  do
33    $\mathcal{B}_{all}[i] \leftarrow \text{OJoinStage5}([\mathcal{P}_{A-in}[i], \mathcal{P}_{B-in}[i]],$ 
      $[\mathcal{B}_{A-agg}[\dots][i], \mathcal{B}_{B-agg}[\dots][i]], [\beta_A, \beta_B], N_{in}, N_{out}')$ 
  /* Do join. */
34 for  $j \leftarrow 0$  to  $N_{out}' - 1$  do
35    $\mathcal{P}_{out}[j] \leftarrow \text{OJoinStage6}(\mathcal{B}_{all}[\dots][j], \gamma/N_{out}' + \alpha_A \cdot \alpha_B)$ ;

```

Algorithm 9: Oblivious Join: Stage 0

```

1 function  $\text{OJoinStage0}(\mathcal{P})$ :
  Input: Partition after aggregate  $\mathcal{P}$ 
  Output: The max count of items belonging to a key  $a$ , the
    total count of items in this partition  $cnt$ 
2    $a \leftarrow 0$ ;
3    $cnt \leftarrow 0$ ;
4   for  $i \leftarrow 0$  to  $|\mathcal{P}| - 1$  do
    /* After  $\text{OAggStage2}$ , only one value correspond to the same key is
    valid. */
5      $c \leftarrow \mathcal{P}[i].v \neq \perp$ ;
6      $\text{cmov}(c, a, \max(a, \mathcal{P}[i].v))$ ;
7      $\text{cmov}(c, cnt, cnt + \mathcal{P}[i].v)$ ;
8   return  $a, cnt$ ;

```

Algorithm 10: Oblivious Join: Stage 1

```

1 function OJoinStage1( $\mathcal{P}_A, \mathcal{P}_B, \alpha_A, \alpha_B, \beta_A, \beta_B, N_{out}$ ):
   Input: Partitions of two tables  $\mathcal{P}_A$  and  $\mathcal{P}_B$ , statistical
         information  $\alpha_A, \alpha_B, \beta_A, \beta_B$ , expected # of nodes after
         join  $N_{out}$ 
   Output: # of nodes after join  $N_{out}'$ , union partition  $\mathcal{P}$ 
2   Add  $t$  field denoting table id to all items in  $\mathcal{P}_A$  and  $\mathcal{P}_B$ ;
   // distinguish the items with the same key in two tables.
3   Alias  $v$  field as  $v.0/v.1$  for A/B, and pad  $v.1/v.0$  respectively;
   /* Items in table B are behind those in table A with the same keys. */
4    $\mathcal{P} \leftarrow \text{OMerge}([\mathcal{P}_A, \mathcal{P}_B], (\_k, \_t, \text{is\_valid}(\_v)))$ ; //  $\mathcal{P}_A$  and
    $\mathcal{P}_B$  are already sorted by  $k$ .
5   if  $(\alpha_A + \alpha_B) / (\beta_A + \beta_B) > 1/N_{out}$  then
6      $N_{out}' \leftarrow \lceil (\beta_A + \beta_B) / (\alpha_A + \alpha_B) \rceil$ ;
7   else
8      $N_{out}' \leftarrow N_{out}$ ;
   /* Invalid items which do not have counterparts to join. */
9   for  $i \leftarrow 0$  to  $|\mathcal{P}| - 1$  do
10     $c \leftarrow \mathcal{P}[i-1].k \neq \mathcal{P}[i].k$ ; // for out-of-bound  $i$  it is true.
11     $\text{cmov}(c \wedge \text{acc}.0 * \text{acc}.1 > 0, \mathcal{P}[i-1].v, \text{acc})$ ;
12     $\text{cmov}(c, \text{acc}, (\perp, \perp))$ ;
13     $\text{acc} \leftarrow \text{acc} + \mathcal{P}[i].v$ ; //  $\perp$  can be seen as 0.
14     $\mathcal{P}[i].v \leftarrow (\perp, \perp)$ ;
15   $\text{cmov}(\text{acc}.0 * \text{acc}.1 > 0, \mathcal{P}[|\mathcal{P}| - 1].v, \text{acc})$ ;
   /* First-level assignment. */
16   $\text{cap} \leftarrow \alpha_A + \alpha_B$ ;
17   $\text{bin\_num} \leftarrow 0$ ;
18   $\text{bin\_size} \leftarrow 0$ ;
19   $\text{acc\_prod} \leftarrow 0$ ;
20  for  $i \leftarrow |\mathcal{P}| - 1$  to 0 do
21     $c \leftarrow \mathcal{P}[i].k \neq \mathcal{P}[i+1].k$ ; // for out-of-bound  $i$  it is true.
22     $\text{cmov}(\neg c, \mathcal{P}[i].v, \mathcal{P}[i+1].v)$ ; // fill the vacant.
   /* Avoid repeated count. */
23     $c_0 \leftarrow c \wedge \mathcal{P}[i].v \neq (\perp, \perp)$ ;
24     $\text{bin\_num}_{\text{new}}, \text{bin\_size} \leftarrow \text{NextFit}(\text{cap}, c_0,$ 
       $\mathcal{P}[i].v.0 + \mathcal{P}[i].v.1, \text{bin\_num}, \text{bin\_size})$ ;
25     $c_1 \leftarrow \text{bin\_num}_{\text{new}} > \text{bin\_num}$ ; // a new bin is created
26     $\mathcal{P}[i].b, \text{bin\_num} \leftarrow \text{bin\_num}_{\text{new}}, \text{bin\_num}_{\text{new}}$ ;
   /* Compute accumulated product for normal items. */
27     $\mathcal{P}[i].a \leftarrow 0$ ; // add new field  $a$ , equal to 0 by default.
28     $\text{cmov}(c_0 \wedge c_1, \mathcal{P}[i+1].a, \text{acc\_prod})$ ;
29     $\text{cmov}(c_0 \wedge c_1, \text{acc\_prod}, \mathcal{P}[i].v.0 \cdot \mathcal{P}[i].v.1)$ ;
30     $\text{cmov}(c_0 \wedge \neg c_1, \text{acc\_prod}, \text{acc\_prod} + \mathcal{P}[i].v.0 \cdot \mathcal{P}[i].v.1)$ ;
31  return  $\text{bin\_num}, \text{bin\_size}, N_{out}', \text{acc\_prod}, \mathcal{P}$ ;
32 function NextFit( $\text{cap}, c, w, \text{bin\_num}, \text{bin\_size}$ ,
    $n = \text{bin\_num} + 1$ ):
   Input: Bin capacity  $\text{cap}$ , condition when bin state changes  $c$ ,
         item weight  $w$ , bin info  $\text{bin\_num}$  and  $\text{bin\_size}$ , and set
         bin number  $n$ , equal to  $\text{bin\_num} + 1$  by default
   Output: changed  $\text{bin\_num}$  and  $\text{bin\_size}$ 
   /* Try to place in the old bin, the bin capacity is  $\text{cap}$ . */
33   $\text{tmp} \leftarrow \text{bin\_size} + w$ ; //  $\perp$  can be viewed as 0 here.
34   $\text{cmov}(c \wedge \text{tmp} \leq \text{cap}, \text{bin\_size}, \text{tmp})$ ;
   /* Open a new bin if space limited. */
35   $\text{cmov}(c \wedge \text{tmp} > \text{cap}, \text{bin\_num}, n)$ ;
36   $\text{cmov}(c \wedge \text{tmp} > \text{cap}, \text{bin\_size}, w)$ ;
37  return  $\text{bin\_num}, \text{bin\_size}$ ;

```

Algorithm 11: Oblivious Join: Stage 2

```

1 function OJoinStage2( $\mathcal{P}, \alpha, j, I, N_{out}$ ):
   Input:  $\mathcal{P}, \alpha$ , and  $N_{out}$ , same meaning as in OJoinStage1,
         current agg partition id  $j$ , last bin information of all
         partitions  $I$ 
   Output: Partition after modification  $\mathcal{P}$ , intermediate result for
         sum  $\mathcal{D}$ , product sum of columns  $\mathcal{S}$ , location of last
         bin  $l$ , last row of bins  $\mathcal{S}_{\text{rem}}$ 
2    $\text{cap} \leftarrow \alpha[0] + \alpha[1]$ ;
3    $\mathcal{P} \leftarrow \text{CoordBinNum}(\mathcal{P}, \text{cap}, j, I)$ ;
   /* Compute accumulated product for last bins. */
4    $j_s, j_{sl}, \text{acc\_prod} \leftarrow 0, 0, 0$ ;
5   for  $i \leftarrow 0$  to  $|I| - 1$  do
6      $\text{cmov}(I[i].0 \neq I[i-1].0, j_s, i)$ ; // false for out-of-bound  $i$ .
7      $\text{cmov}(j_s = i \wedge j_s \leq j, \text{acc\_prod}, I[i].2)$ ;
8      $\text{cmov}(j_s \neq i \wedge j_s \leq j, \text{acc\_prod}, \text{acc\_prod} + I[i].2)$ ;
9      $\text{cmov}(j_s \leq j, j_{sl}, j_s)$ ;
10   $\text{cmov}(j_{sl} \neq j, \text{acc\_prod}, -\text{acc\_prod})$ ;
11   $\mathcal{P}[0].a \leftarrow \text{acc\_prod}$ ; // items with non-positive number does not
   involve 2nd assignment.
   /* Items with the same  $b$  share the same  $a$ . */
12  for  $i \leftarrow 1$  to  $|\mathcal{P}| - 1$  do
13     $\text{cmov}(\mathcal{P}[i].b = \mathcal{P}[i-1].b, \mathcal{P}[i].a, \mathcal{P}[i-1].a)$ ;
   /* Begin second level assignment to balance the results after join. */
14   $\text{idx} \leftarrow 0$ ; // mark and count the representatives of bins.
15   $l \leftarrow -1$ ; // whether the last bin involves second assignment.
16  for  $i \leftarrow 0$  to  $|\mathcal{P}| - 1$  do
17     $c \leftarrow \mathcal{P}[i].b \neq \mathcal{P}[i-1].b \wedge \mathcal{P}[i].a > 0$ ;
18     $\mathcal{D}[i] \leftarrow (i, \perp, 0)$ ; // true for out-of-bound  $i$ .
19     $\text{cmov}(c, \mathcal{D}[i], (i, \text{idx}, \mathcal{P}[i].a))$ ;
20     $\text{cmov}(c, \text{idx}, \text{idx} + 1)$ ;
   /* Summarize the accumulated product by column. */
21   $\mathcal{S} \leftarrow [(0, 0), (1, 0), (2, 0), (3, 0), \dots, (N_{out} - 1, 0)]$ ;
22   $\text{cur} \leftarrow 0$ ;
23   $\text{tmp} \leftarrow [(0, \perp, \perp), (1, \perp, \perp), \dots, (N_{out} - 1, \perp, \perp)]$ ;
24  for  $i \leftarrow 0$  to  $\lceil |\mathcal{D}| / N_{out} \rceil - 1$  do
25     $\text{chunk} \leftarrow \mathcal{D}[i \times N_{out} : (i+1) \times N_{out}]$ ; // reference
26    foreach  $d \in \text{chunk}$  do
27       $\text{cmov}(d.1 \neq \perp, d.1, \text{cur mod } N_{out})$ ;
28       $\text{cmov}(d.1 \neq \perp, \text{cur}, \text{cur} + 1)$ ;
29    OSort( $\text{chunk}, \_1$ );
30    OPlace( $\text{chunk}, \_1$ );
31    for  $j \leftarrow 0$  to  $N_{out} - 1$  do
32       $\text{cmov}(\text{chunk}[j].1 \neq \perp \wedge \text{tmp}[j].1 = \perp, \text{tmp}[j].1,$ 
         $\text{chunk}[j].2)$ ;
33    OSort( $\mathcal{S}, \_1$ ); // assume it is in descending order.
34    OSort( $\text{tmp}, \_1$ ); // assume it is in ascending order, sort dummy
        to the end.
35     $c \leftarrow \text{cur} \geq N_{out}$ ;
36    for  $j \leftarrow 0$  to  $N_{out} - 1$  do
37       $\text{cmov}(c, \mathcal{S}[j].1, \mathcal{S}[j].1 + \text{tmp}[j].1)$ ;
38       $\text{cmov}(c, \text{tmp}[j].2, \mathcal{S}[j].0)$ ;
39    OSort( $\text{tmp}, \_0$ ); // assume it is in ascending order.
40     $\text{cur} \leftarrow \text{cur mod } N_{out}$ ;
41    for  $j \leftarrow 0$  to  $N_{out} - 1$  do
42       $\text{cmov}(c \wedge j < \text{cur}, \text{tmp}[j].1, \text{chunk}[j].2)$ ;
43       $\text{cmov}(c \wedge j \geq \text{cur}, \text{tmp}[j].1, \perp)$ ;
44       $\text{cmov}(c, \text{chunk}[j].2, \text{tmp}[j].2)$ ;
45       $\text{cmov}(\neg c, \text{chunk}[j].2, \perp)$ ;
46       $\text{cmov}(c \wedge \mathcal{P}[0].a > 0 \wedge l = -1, l, \text{chunk}[0].2)$ ;
        // record the original column num for the last bin.
47   $\text{cmov}(\mathcal{P}[0].a > 0 \wedge l = -1, l, -2)$ ; // the last bin is in  $\mathcal{S}_{\text{rem}}$ .
48   $\mathcal{S}_{\text{rem}} \leftarrow \text{tmp}[0 : N_{out}]$ ; // strip field 2.
49  return  $\mathcal{P}, \mathcal{D}, \mathcal{S}, l, \mathcal{S}_{\text{rem}}$ ;

```

Algorithm 12: Assign Remainder Bins

```
1 function AssignRemBin( $S_{rem}, N_{out}$ ):  
   Input: Remaining bin info  $S_{rem}$ , number of nodes after join  
            $N_{out}$   
   Output: Initial column sum info  $S$ , remaining bin info after  
           assignment  $S_{rem}$   
2    $S \leftarrow [(0, 0), (1, 0), \dots, (N_{out} - 1, 0)]$ ;  
3    $S_{rem} \leftarrow S_{rem}.zipWithIndex.flatten$ ; // add field  $\_2$ , the  
           partition id.  
4    $OSort(S_{rem}, \_1)$ ; // in ascending order, sorted  $\perp$  to the end.  
5   for  $i \leftarrow 0$  to  $|S_{rem}|$  do  
6      $k \leftarrow i \bmod N_{out}$ ;  
7      $cmov(S_{rem}[i].1 \neq \perp, S[k].1, S[k].1 + S_{rem}[i].1)$ ;  
8      $cmov(S_{rem}[i].1 \neq \perp, S_{rem}[i].1, k)$ ; // reinterpret the  $\_1$   
           field to record the final assignment result.  
9    $OSort(S_{rem}, (\_2, \_0))$ ;  
10   $S_{rem} \leftarrow \text{Split } S_{rem} \text{ by different } \_2 \text{ and strip field } \_2$ ;  
11  return  $S, S_{rem}$ ;
```

Algorithm 13: Oblivious Join: Stage 3

```
1 function OJoinStage3( $S_r, S_o, l_r, l_o, S_{rem}$ ):  
   Input: Received and owned column sum info  $S_r$  and  $S_o$ , last  
           bin location  $l_r$  and  $l_o$ , remaining bin info  $S_{rem}$   
   Output: To-be-sent and owned column sum info  $S_r$  and  $S_o$ ,  
           updated bin location  $l_o$   
2    $N_{out} \leftarrow |S_o|$ ;  
3    $OSort(S_r, \_1)$ ; // descending.  
4    $OSort(S_o, \_1)$ ; // ascending.  
5    $S_r.1 \leftarrow S_o.1 + S_r.1$ ; // vectorized add.  
6    $S_o.1 \leftarrow S_r.0$ ; // record the reassignment result.  
7    $OSort(S_o, \_0)$ ; // restore the original order.  
   /* Adjust partition number for the last bin. */  
8    $l \leftarrow l_o$ ;  
9   foreach  $s \in S_o$  do  
10     $c \leftarrow l \geq 0 \wedge l = s.0$ ;  
11     $cmov(c, l_o, s.1)$ ; // if the last bin is in regular bins.  
12   $cmov(l = -2, l_o, S_{rem}[0].1)$ ; // if the last bin is in remaining bins.  
13   $cmov(l = -1, l_o, l_r)$ ;  
14  return  $S_r, S_o, l_o$ ;
```

Algorithm 14: Coordinate Bin Number

```
1 function CoordBinNum( $\mathcal{P}, cap, j, I$ ):  
   Input: Partition after aggregate  $\mathcal{P}$ , bin capacity  $cap$ , partition  
           id  $j$ , last bin information of all partitions  $I$   
   Output: Partition after adjusting bin number  $\mathcal{P}$   
   /* Decide the first and the last bin number of this partition. */  
2    $bin\_num \leftarrow I[0].0$ ;  
3    $bin\_size \leftarrow I[0].1$ ;  
4    $n \leftarrow bin\_num + 1$ ;  
5    $n_{last} \leftarrow 0$ ;  
6   for  $i \leftarrow 1$  to  $|I| - 1$  do  
   /* Try to place in the old bin, the bin capacity is  $cap$ . */  
7      $cur\_bin\_num \leftarrow I[i].0$ ;  
8      $num, size \leftarrow \text{NextFit}(cap, true, I[i].1, bin\_num,$   
            $bin\_size, n + cur\_bin\_num)$ ;  
9     if  $i = j$  then  
10       $n_{last} \leftarrow n$ ;  
11      $cmov(bin\_size + I[i].1 > cap, n, n + 1)$ ; // as a new bin is  
           created.  
12      $n \leftarrow n + cur\_bin\_num$ ;  
13      $bin\_num, bin\_size \leftarrow num, size$ ;  
14      $I[i].0 \leftarrow bin\_num$ ;  
   /* Assign new bin number according to old bin number. */  
15    $tmp\_num \leftarrow \mathcal{P}[0].b$ ;  
16   for  $i \leftarrow 0$  to  $|\mathcal{P}| - 1$  do  
17      $c \leftarrow \mathcal{P}[i].b \neq tmp\_num$ ;  
18      $cmov(c, \mathcal{P}[i].b, (n_{last} + \mathcal{P}[i].b))$ ;  
19      $cmov(\neg c, \mathcal{P}[i].b, I[j].0)$ ;  
20  return  $\mathcal{P}$ ;
```

Algorithm 16: Oblivious Join: Stage 5

```
1 function OJoinStage5( $\mathcal{P}, \mathcal{B}, \beta, N_{in}, N_{out}$ ):
   Input: Partition  $\mathcal{P}$ , buckets with aggregate information  $\mathcal{B}$ ,
         total number of items  $\beta$ , # of nodes before group by
          $N_{in}$ , # of nodes after group by  $N_{out}$ 
   Output: Buckets with original data  $\mathcal{B}'$ 

2   for  $t \leftarrow 0$  to 1 do
3      $\mathcal{P}[t] \leftarrow \text{PatchPartNum}(\mathcal{P}[t], \mathcal{B}[t], N_{in})$ ;
4   compute the padding length  $l$ ;
5    $\mathcal{B}' \leftarrow \text{BuildBuckets}(\mathcal{P}.\text{flatten}(), l, N_{out})$ 
6   return  $\mathcal{B}'$ ;

7 function PatchPartNum( $\mathcal{P}, \mathcal{B}, N_{in}$ ):
8    $\mathcal{D} \leftarrow \text{OMerge}(\mathcal{B}, \_l)$ ;
9   for  $i \leftarrow 0$  to  $|\mathcal{D}| - 1$  do
10     $\mathcal{P}[i].b \leftarrow \mathcal{D}[i].b$ ;
    /*  $\mathcal{P}$  has three fields:  $k, v, b$ . The last one in the items with the same key
       always has the accurate new partition number. Use it to fix others
       with the same key. */
11   for  $i \leftarrow |\mathcal{P}| - 2$  to 0 do
12      $c \leftarrow \mathcal{P}[i].k = \mathcal{P}[i+1].k$ ;
13      $\text{cmov}(c, \mathcal{P}[i].b, \mathcal{P}[i+1].b)$ ;
14   return  $\mathcal{P}$ ;
```

Algorithm 17: Oblivious Join: Stage 6

```
1 function OJoinStage6( $\mathcal{B}, len$ ):
   Input: Buckets collected from all nodes  $\mathcal{B}$ , pad length for join
         result  $len$ 
   Output: Partition after join  $\mathcal{P}$ 

   /* Note that each bucket is already sorted by  $k$  */
2    $\mathcal{P} \leftarrow \text{OMerge}(\mathcal{B}, (\neg \text{is\_valid}(v), \_k))$ ; // merge by  $k$ 
3    $\mathcal{P} \leftarrow \text{Apply single-node oblivious join in [1] on } \mathcal{P}$ . During the
         oblivious expansion, pad both to  $len$ ;
4   return  $\mathcal{P}$ ;
```

Algorithm 15: Oblivious Join: Stage 4

```
1 function OJoinStage4( $\mathcal{P}, \mathcal{D}, \mathcal{S}, l, S_{rem}$ ):
   Input: Partition after aggregate  $\mathcal{P}$ , intermediate result for sum
          $\mathcal{D}$ , product sum of columns  $\mathcal{S}$ , last bin location  $l$ , last
         row of bins  $S_{rem}$ 
   Output: Buckets  $\mathcal{B}$ 

   /* Prepare the patch of new partition numbers. */
2    $N_{out} \leftarrow |\mathcal{S}|$ ;
3    $cur \leftarrow 0$ ;
4   for  $i \leftarrow 0$  to  $|\mathcal{D}| - 1$  do
5      $\text{cmov}(\mathcal{D}[i].1 \neq \perp, cur, (cur + 1) \bmod N_{out})$ ;
6    $tmp \leftarrow S_{rem}$ ; //  $S_{rem}$  is already sorted by the index (mod  $N_{out}$ )
7   for  $i \leftarrow \lceil |\mathcal{D}|/N_{out} \rceil - 1$  to 0 do
8      $chunk \leftarrow \mathcal{D}[i \times N_{out} : (i+1) \times N_{out}]$ ; // reference
9      $buf \leftarrow chunk$ ; // a copy.
10     $should\_update \leftarrow \text{whether } \_1 \neq \perp \text{ is for all in chunk}$ ;
11     $last \leftarrow cur$ ;
12    for  $j \leftarrow N_{out} - 1$  to 0 do
13       $c \leftarrow chunk[j].1 \neq \perp \wedge chunk[j].1 < last$ ;
14       $\text{cmov}(c, chunk[j].2, tmp[j].1)$ ; // assign partition
        number to current chunk
15       $\text{cmov}(c, cur, cur - 1)$ ;
16       $buf[j].1 \leftarrow j$ ; // record the original order.
17     $\text{OSort}(buf, \_2)$ ; // By old column number. If it contains dummy,
        it will not be used. So no OPlace is needed.
18     $buf[\_2].2 \leftarrow \mathcal{S}[\_2].1$ ; // vectorized assignment.
19     $\text{OSort}(buf, \_1)$ ; // restore the previous order.
20     $\text{cmov}(cur = 0 \wedge should\_update, cur, N_{out})$ ;
21    for  $j \leftarrow N_{out} - 1$  to 0 do
22       $\text{cmov}(should\_update, tmp[j].1, buf[j].2)$ ; // update
         $tmp$  for patching next chunk.
23       $c \leftarrow chunk[j].1 \neq \perp \wedge chunk[j].1 \geq last$ ;
24       $\text{cmov}(c, chunk[j].2, tmp[j].1)$ ; // assign partition
        number to next chunk
25       $\text{cmov}(c, cur, cur - 1)$ ;
26     $\text{OSort}(chunk, \_0)$ ; // the original order.

   /* Assign new partition number. */
27    $\mathcal{P}[0].b \leftarrow l$ ;
28   for  $i \leftarrow 1$  to  $|\mathcal{P}| - 1$  do
29      $\text{cmov}(\mathcal{D}[i].1 \neq \perp, \mathcal{P}[i].b, \mathcal{D}[i].2)$ ;
30      $\text{cmov}(\mathcal{D}[i].1 = \perp, \mathcal{P}[i].b, \mathcal{P}[i-1].b)$ ;

   /* Prepare to patch new partition number back to original items. */
31    $\text{OSort}(\mathcal{P}, (\_2, \_1, \_i, \_l))$ ; // sort by  $(t, i, l)$ 
32    $\mathcal{B} \leftarrow \text{Split } \mathcal{P}$  by different  $t, i$  and strip fields  $t, i$ ;
33   return  $\mathcal{B}$ ;
```

REFERENCES

- [1] Simeon Krastnikov, Florian Kerschbaum, and Douglas Stebila. 2020. Efficient Oblivious Database Joins. *Proc. VLDB Endow.* 13, 12 (jul 2020), 2132–2145. <https://doi.org/10.14778/3407790.3407814>