

Algorithm 1 Creates $H \in \mathbb{R}^{(3m) \times (3n)}$

```
1: procedure HSTA_KERNEL
2:    $j := \text{blockIdx}.x$ ;  $i := \text{blockIdx}.y$ ;  $x := \text{threadIdx}.x$ ;  $y := \text{threadIdx}.y$ 
3:    $l := \text{laneId}$ ;  $w := \text{warpId}$ ;  $\text{numWarps} = \lceil g^2/32 \rceil$ 
4:   Initialize  $Hbuffer$  in shared memory with 0
5:    $Hlocal \leftarrow \text{GenerateMatrixHsta}(i, j, x, y)$  ▷ Return  $3 \times 3$  matrix
6:    $\text{shfl\_down}(Hlocal)$  ▷ Reduces matrices in the same warp
7:   if  $l = 0$  then
8:      $Hbuffer[\text{warpId}] = Hlocal$ 
9:   if  $x < 3$  and  $y < 3$  then
10:     $Helem[x][y] \leftarrow \text{thrust}::\text{reduce}(Hbuffer[x][y][0], Hbuffer[x][y][\text{numWarps}])$ 
11:     $H[i][j] \leftarrow Helem$  ▷  $H[i][j]$  is a  $3 \times 3$  matrix
12: procedure RIGID_KERNEL( $H, Hdiag$ )
13:    $t := \text{blockDim}.x \times \text{blockIdx}.x + \text{threadIdx}.x$ 
14:   if  $t < m$  then
15:     for  $k := 1, n$  do
16:        $Hdiag[t] \leftarrow Hdiag[t] - H[t][k]$ 
17: procedure HSTA_ASSEMBLY
18:   Move data to GPU memory
19:   Allocate  $H \in \mathbb{R}^{(3m) \times (3n)}$  in GPU memory
20:   Allocate  $Hdiag \in \mathbb{R}^{m \times 3 \times 3}$  in GPU memory and initialize with 0
21:   Allocate  $Hbuffer$  in GPU shared memory, buffer of matrices  $3 \times 3$  of size  $\lceil g^2/32 \rceil$ 
22:   Run Hsta_kernel with  $m \times n$  blocks and  $g \times g$  threads. Await for return
23:   Run Rigid_kernel with 128 threads and  $\lceil m/128 \rceil$  blocks. Await for return
24:   Keep  $Hdiag$  in GPU memory and free  $H$  from GPU memory.
```
