

Code

Blame

227 lines (198 loc) · 7.82 KB

Raw

📄

📥

✎

▼

🔗

```
1  /**
2   * @brief Flips bits of single-precision floating-point number
3   *
4   * flip a float for sorting
5   * finds SIGN of fp number.
6   * if it's 1 (negative float), it flips all bits
7   * if it's 0 (positive float), it flips the sign only
8   * @param[in] f floating-point input (passed as unsigned int)
9   * @returns uint that stores the flipped version of the input
10  * @see floatUnflip
11  */
12  __device__ uint float_flip(uint f)
13  {
14      uint mask = -int(f >> 31) | 0x80000000; // -int(0): 0x00000000; -int(1): 0xffffffff
15      return f ^ mask;
16  }
17
18  /**
19   * @brief Reverses bit-flip of single-precision floating-point number
20   *
21   * flip a float back (invert FloatFlip)
22   * signed was flipped from above, so:
23   * if sign is 1 (negative), it flips the sign bit back
24   * if sign is 0 (positive), it flips all bits back
25   * @param[in] f floating-point input (passed as unsigned int)
26   * @returns uint that stores the unflipped version of the input
27   * @see floatFlip
28  */
29  __device__ uint float_unflip(uint f)
30  {
31      uint mask = ((f >> 31) - 1) | (0x80000000);
32      return f ^ mask;
33  }
34
35  template <int BitsPerPass>
36  static __global__ void histogram_kernel(const uint* in_buf, int* histogram,
37                                         const Counter* counter, const int pass, bool select_min)
38  {
39      const int num_elements = counter->len;
40      const int tid = blockIdx.x * blockDim.x + threadIdx.x;
41
42      for (int idx = tid; idx < num_elements; idx += gridDim.x * blockDim.x)
43      {
44          uint mask = float_flip(in_buf[idx]);
45          mask = mask >> (32 - BitsPerPass * (pass + 1));
46          uint bucket = mask & ((0x01 << BitsPerPass) - 1);
47          if (!select_min)
48          {
49              // bucket = ~bucket; // 不能按位取反, 正确做法是对最低8位取反
50              bucket = bucket ^ ((0x01 << BitsPerPass) - 1);
51          }
52
53          atomicAdd(&histogram[bucket], 1);
54          // atomicAggInc_cg(histogram + bucket); // 结果异常, 因为一个warp中的所有线程都是活跃的, 只是执行原子操作的bucket不同
55      }
56
57  }
58
59  template <int BitsPerPass>
60  static __global__ void scan_select_kernel(int* histogram, Counter* counter)
61  {
62      constexpr int num_buckets = 256; num_buckets < BitsPerPass();
```

```

62     constexpr int num_buckets = calc_num_buckets<bitsPerPass>(),
63     constexpr int num_warps    = num_buckets / kWarpSize;
64
65     __shared__ int sums[num_warps];
66     __shared__ int histogram_smem[num_buckets];
67     const int tid = blockIdx.x * blockDim.x + threadIdx.x;
68
69     for (int i = tid; i < num_buckets; i += blockDim.x)
70     {
71         histogram_smem[i] = histogram[i];
72         histogram[i] = 0; // 设置histogram的所有值为0, 用于初始化下一次pass
73     }
74     __syncthreads();
75
76     int lane_id = tid % kWarpSize;
77     int warp_id = threadIdx.x / kWarpSize;
78
79     int value = histogram_smem[tid];
80     // using a shfl instruction for a scan. (Hillis Steele scan)
81     for (int i = 1; i < kWarpSize; i *= 2)
82     {
83         uint mask = 0xffffffff;
84         int n = __shfl_up_sync(mask, value, i, kWarpSize);
85         if (lane_id >= i)
86         {
87             value += n;
88         }
89     }
90
91     // write the sum of the warp to smem
92     if (threadIdx.x % kWarpSize == kWarpSize - 1)
93     {
94         sums[warp_id] = value;
95     }
96     __syncthreads();
97
98     // scan sum the warp sums
99     if (warp_id == 0 && lane_id < num_warps)
100    {
101        int warp_sum = sums[lane_id];
102        uint mask = (0x01 << num_warps) - 1;
103        for (int i = 1; i < num_warps; i *= 2)
104        {
105            int n = __shfl_up_sync(mask, warp_sum, i, num_warps);
106            if (lane_id >= i)
107            {
108                warp_sum += n;
109            }
110        }
111
112        sums[lane_id] = warp_sum;
113    }
114    __syncthreads();
115
116    // perform a uniform add across warps in the block
117    // read neighbouring warp's sum and add it to threads value
118    int block_sum = 0;
119    if (warp_id > 0)
120    {
121        block_sum = sums[warp_id - 1];
122    }
123    histogram_smem[tid] = value + block_sum;
124
125    // choose bucket
126    int k = counter->k;
127    int len = counter->len;
128    __syncthreads();
129
130    for (int i = tid; i < num_buckets; i += blockDim.x)
131    {
132        int pre = (i == 0) ? 0 : histogram_smem[i - 1];
133        int cur = histogram_smem[i];
134
135        // one and only one thread will satisfy this condition, so counter is written by only one thread
136        if (pre < k && cur >= k)
137        {

```

```

138         counter->k = k - pre;        // how many values still are there to find
139         counter->len = cur - pre;    // number of values in next pass
140         counter->prev_len = len;
141         counter->bucket_bits = i;    // bucket
142     }
143 }
144
145 if (tid == 0)
146 {
147     counter->filter_cnt = 0;    // 重新设置filter_cnt的值为0
148 }
149 }
150
151 template <int BitsPerPass>
152 static __global__ void filter_kernel(const uint* in_buf, const int* in_idx,
153                                     uint* out_buf, int* out_idx,
154                                     int* out_index, Counter* counter,
155                                     const int pass, const bool select_min)
156 {
157     constexpr int num_buckets = cal_num_buckets<BitsPerPass>();
158     uint select_bucket = counter->bucket_bits;
159     uint bucket_bits = (select_min ? select_bucket : (select_bucket ^ ((0x01 << BitsPerPass) - 1)));
160
161     const int tid = blockIdx.x * blockDim.x + threadIdx.x;
162     for (int idx = tid; idx < counter->prev_len; idx += gridDim.x * blockDim.x)
163     {
164         uint mask = float_flip(in_buf[idx]);
165         mask = mask >> (32 - BitsPerPass * (pass + 1));
166         uint bits = mask & ((0x01 << BitsPerPass) - 1);
167
168         if (bits == bucket_bits)
169         {
170             // int pos = atomicAdd(&counter->filter_cnt, 1);
171             int pos = atomicAggInc_cg(&counter->filter_cnt);
172             out_buf[pos] = in_buf[idx];
173             out_idx[pos] = in_idx[idx];
174         }
175         else if ((bits < bucket_bits) == select_min)
176         {
177             // int pos = atomicAdd(&counter->output_cnt, 1);
178             int pos = atomicAggInc_cg(&counter->output_cnt);
179             out_index[pos] = in_idx[idx];
180             // 如果实际的topk的值，在数据集中有多个相等的值，（这里收集到的值肯定小于topk）
181             // 则在最后一个pass后，将out_idx的指标都附加到out_index后面，直到数据个数为topk个
182         }
183     }
184 }
185
186 static __global__ void gather_kernel(int* out_index, int* in_index,
187                                     Counter* counter, const int* num_elements,
188                                     const int topk)
189 {
190     const int output_len = MIN(*num_elements, topk);
191     const int tid = blockIdx.x * blockDim.x + threadIdx.x;
192     for (int idx = tid; idx < (output_len - counter->output_cnt); idx += gridDim.x * blockDim.x)
193     {
194         out_index[counter->output_cnt + idx] = in_index[idx];
195     }
196 }
197
198 // topk operator
199 {
200     for (int pass = 0; pass < num_passes; ++pass)
201     {
202         // 1. Calculate histogram
203         constexpr int TS_topk = 16;
204         const int GS_topk = (total_num_score_ + (BS * TS_topk) - 1) / (BS * TS_topk);
205         histogram_kernel<BitsPerPass><<<GS_topk, BS, 0, cu_stream>>>
206             (in_buf, histogram, counter, pass, select_min);
207
208         // 2. Scan the histogram (Inclusive prefix sum)
209         // 3. Choose the bucket (Find the bucket j of the histogram that the k-th value falls into)
210         scan_select_kernel<BitsPerPass><<<1, num_buckets, 0, cu_stream>>>
211             (histogram, counter);
212     }
213 }

```

```
214
215     // 4. Filtering (Input elements whose digit value <j are the top-k elements)
216     filter_kernel<BitsPerPass><<<GS_topk, BS, 0, cu_stream>>>
217         (in_buf, in_idx, out_buf, out_idx,
218         output_index, counter, pass, select_min);
219
220     in_buf  = (pass % 2 == 0) ? buf_value_1 : buf_value_2;
221     in_idx  = (pass % 2 == 0) ? buf_index_1 : buf_index_2;
222     out_buf = (pass % 2 == 0) ? buf_value_2 : buf_value_1;
223     out_idx = (pass % 2 == 0) ? buf_index_2 : buf_index_1;
224 }
225
226 gather_kernel<<<1, BS, 0, cu_stream>>>(output_index, in_idx, counter, ws1_data, total_num_out_);
227 }
```