

ca6f96b62a ▾

...

tensorflow / tensorflow / lite / kernels / internal / utils / sparsity_format_converter.cc



rino20 Fix dependency from tf/lite/kernels to tf/lite/tools ... ✖

History

1 contributor

392 lines (348 sloc) | 14 KB

...

```

1  /* Copyright 2020 The TensorFlow Authors. All Rights Reserved.
2
3  Licensed under the Apache License, Version 2.0 (the "License");
4  you may not use this file except in compliance with the License.
5  You may obtain a copy of the License at
6
7      http://www.apache.org/licenses/LICENSE-2.0
8
9  Unless required by applicable law or agreed to in writing, software
10 distributed under the License is distributed on an "AS IS" BASIS,
11 WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
12 See the License for the specific language governing permissions and
13 limitations under the License.
14 =====*/
15 #include "tensorflow/lite/kernels/internal/utils/sparsity_format_converter.h"
16
17 #include <stdint>
18 #include <utility>
19 #include <vector>
20
21 namespace tf lite {
22 namespace internal {
23 namespace sparsity {
24
25 namespace {
26 uint64_t GetFlattenedIndex(const std::vector<int>& indices,
27                           const std::vector<int>& shape) {
28     uint64_t index = 0;
29     int sub_elements = 1;

```

```

30     for (int i = shape.size() - 1; i >= 0; i--) {
31         index += indices[i] * sub_elements;
32         sub_elements *= shape[i];
33     }
34     return index;
35 }
36
37 std::vector<int> TfLiteIntArrayToVector(const TfLiteIntArray* int_array) {
38     std::vector<int> values;
39     if (!int_array) {
40         return values;
41     }
42
43     values.resize(int_array->size);
44     for (size_t i = 0; i < int_array->size; i++) {
45         values[i] = int_array->data[i];
46     }
47
48     return values;
49 }
50
51 } // namespace
52
53 template <typename T>
54 FormatConverter<T>::FormatConverter(
55     const std::vector<int>& shape, const std::vector<int>& traversal_order,
56     const std::vector<TfLiteDimensionType>& format,
57     const std::vector<int>& block_size, const std::vector<int>& block_map)
58     : dense_shape_(shape),
59       traversal_order_(traversal_order),
60       block_size_(block_size),
61       block_map_(block_map) {
62     dense_size_ = 1;
63     int block_dim = 0;
64     blocked_shape_.resize(shape.size());
65     format_.resize(shape.size() + block_map.size());
66     for (int i = 0; i < shape.size(); i++) {
67         format_[i] = format[traversal_order[i]];
68         dense_size_ *= shape[i];
69         if (block_dim < block_map.size() && block_map[block_dim] == i) {
70             blocked_shape_[i] = shape[i] / block_size[block_dim];
71             block_dim++;
72         } else {
73             blocked_shape_[i] = shape[i];
74         }
75     }
76
77     // Only dense blocks are supported.
78     for (int i = 0; i < block_map.size(); i++) {

```

```

79     format_[i + shape.size()] = kTfLiteDimDense;
80 }
81 }
82
83 template <typename T>
84 TfLiteStatus FormatConverter<T>::DenseToSparse(const T* src_data) {
85     int num_original_dims = dense_shape.size();
86     int num_block_dims = block_map.size();
87     int num_expanded_dims = num_original_dims + num_block_dims;
88     std::vector<int> expanded_shape(num_expanded_dims);
89     for (int i = 0; i < num_expanded_dims; i++) {
90         if (i < num_original_dims) {
91             expanded_shape[i] = blocked_shape[i];
92         } else {
93             expanded_shape[i] = block_size[i - num_original_dims];
94         }
95     }
96
97     std::vector<int> shape_offset(num_original_dims);
98     shape_offset[shape_offset.size() - 1] = 1;
99     for (int i = num_original_dims - 1; i > 0; --i) {
100         shape_offset[i - 1] = shape_offset[i] * dense_shape[i];
101     }
102
103     std::vector<int> expanded_shape_offset(num_expanded_dims);
104     for (int i = 0; i < num_original_dims; ++i) {
105         expanded_shape_offset[i] = shape_offset[i];
106     }
107     for (int i = 0; i < num_block_dims; ++i) {
108         int mapped_dim = block_map[i];
109         expanded_shape_offset[num_original_dims + i] = shape_offset[mapped_dim];
110         expanded_shape_offset[mapped_dim] *= block_size[i];
111     }
112
113     std::vector<int> dst_ordered_offset(num_expanded_dims);
114     for (int i = 0; i < num_expanded_dims; ++i) {
115         dst_ordered_offset[i] = expanded_shape_offset[traversal_order_[i]];
116     }
117
118     std::vector<bool> dst_dim_has_nonzeroes(num_expanded_dims);
119     std::fill(dst_dim_has_nonzeroes.begin(), dst_dim_has_nonzeroes.end(), false);
120     std::vector<int> inner_compressed_dim(num_expanded_dims);
121     int most_recent_compressed_dim = -1;
122     std::vector<int> num_segments_of_next_compressed_dim(num_expanded_dims);
123     int segment_count = 1;
124     for (int i = num_expanded_dims - 1; i >= 0; --i) {
125         inner_compressed_dim[i] = most_recent_compressed_dim;
126         if (format_[i] == kTfLiteDimSparseCSR) {
127             most_recent_compressed_dim = i;

```

```

128     num_segments_of_next_compressed_dim[i] = segment_count;
129     segment_count = 1;
130 } else {
131     num_segments_of_next_compressed_dim[i] = -1;
132     segment_count *= expanded_shape[traversal_order_[i]];
133 }
134 }
135
136 dim_metadata_.resize(num_expanded_dims * 2);
137 std::vector<int> dst_sparse_dims;
138 dst_sparse_dims.reserve(num_expanded_dims);
139 for (int i = 0; i < num_expanded_dims; ++i) {
140     dim_metadata_[i * 2].clear();
141     dim_metadata_[i * 2 + 1].clear();
142     if (format_[i] == kTfLiteDimDense) {
143         // If dimension is dense, just store the shape.
144         dim_metadata_[i * 2].push_back(expanded_shape[traversal_order_[i]]);
145     } else {
146         dim_metadata_[i * 2].push_back(0); // Segment array always begins with 0.
147         dst_sparse_dims.push_back(i);      // Add dimension to the sparse list.
148     }
149 }
150
151 // This algorithm assumes that the block size is small enough for all the
152 // elements to fit in cache, so the strided accesses from different traversal
153 // order and the write-first-erase-later strategy shouldn't be too slow
154 int dst_dim_idx = num_expanded_dims;
155 std::vector<int> coordinate(num_expanded_dims, 0);
156 int dense_tensor_idx = 0;
157 while (dst_dim_idx >= 0) {
158     if (dst_dim_idx == num_expanded_dims) {
159         // We have a complete coordinate. Add the element to the value array if it
160         // is not zero, or if the last dimension is dense.
161         if (!IsZero(src_data[dense_tensor_idx])) {
162             data_.push_back(src_data[dense_tensor_idx]);
163             // Mark all sparse dimensions that their current indices have nonzeros.
164             for (auto dst_dim : dst_sparse_dims) {
165                 if (!dst_dim_has_nonzeroes[dst_dim]) {
166                     // Only add the index to the indices array if the current nonzero
167                     // is the first nonzero of the block.
168                     dim_metadata_[2 * dst_dim + 1].push_back(coordinate[dst_dim]);
169                     dst_dim_has_nonzeroes[dst_dim] = true;
170                 }
171             }
172         } else if (format_[num_expanded_dims - 1] == kTfLiteDimDense) {
173             data_.push_back(src_data[dense_tensor_idx]);
174         }
175         --dst_dim_idx;
176     } else {

```

```

177     int original_dim_idx = traversal_order_[dst_dim_idx];
178     int dim_size = expanded_shape[original_dim_idx];
179     if (dst_dim_has_nonzeroes[dst_dim_idx]) {
180         // If the previous block has nonzeroes, reset the flag to false since
181         // we have just moved to a new block.
182         dst_dim_has_nonzeroes[dst_dim_idx] = false;
183     } else if (format_[dst_dim_idx] == kTfLiteDimSparseCSR) {
184         // This block is empty. Delete unnecessary values if compressed.
185         int next_compressed_dim = inner_compressed_dim[dst_dim_idx];
186         int erase_offset = dim_metadata_[2 * dst_dim_idx + 1].size() *
187             num_segments_of_next_compressed_dim[dst_dim_idx];
188         if (next_compressed_dim >= 0) {
189             auto& segments = dim_metadata_[2 * inner_compressed_dim[dst_dim_idx]];
190             segments.erase(segments.begin() + 1 + erase_offset, segments.end());
191         } else {
192             data_.erase(data_.begin() + erase_offset, data_.end());
193         }
194     }
195     if (++coordinate[dst_dim_idx] < dim_size) {
196         // The current dst_dim_idx is valid (not out of bound).
197         dense_tensor_idx += dst_ordered_offset[dst_dim_idx];
198         ++dst_dim_idx;
199     } else {
200         // dst_dim_idx has reached its dim size. Update segment array and go
201         // back to incrementing the previous dimension (dst_dim_idx - 1).
202         if (format_[dst_dim_idx] == kTfLiteDimSparseCSR) {
203             dim_metadata_[2 * dst_dim_idx].push_back(
204                 dim_metadata_[2 * dst_dim_idx + 1].size());
205         }
206         coordinate[dst_dim_idx] = -1;
207         dense_tensor_idx -= dst_ordered_offset[dst_dim_idx] * dim_size;
208         --dst_dim_idx;
209     }
210 }
211 }
212
213 return kTfLiteOk;
214 }
215
216 template <typename T>
217 FormatConverter<T>::FormatConverter(
218     const std::vector<int>& shape, const std::vector<int>& traversal_order,
219     const std::vector<TfLiteDimensionType>& format,
220     const std::vector<int>& dense_size,
221     const std::vector<std::vector<int>>& segments,
222     const std::vector<std::vector<int>>& indices,
223     const std::vector<int>& block_map) {
224     InitSparseToDenseConverter(shape, traversal_order, format, dense_size,
225         segments, indices, block_map);

```

```

226 }
227
228 template <typename T>
229 FormatConverter<T>::FormatConverter(const std::vector<int>& shape,
230                                   const TfLiteSparsity& sparsity) {
231     auto traversal_order = TfLiteIntArrayToVector(sparsity.traversal_order);
232     auto block_map = TfLiteIntArrayToVector(sparsity.block_map);
233
234     std::vector<TfLiteDimensionType> format(sparsity.dim_metadata_size);
235     std::vector<int> dense_size(sparsity.dim_metadata_size);
236     std::vector<std::vector<int>> segments(sparsity.dim_metadata_size);
237     std::vector<std::vector<int>> indices(sparsity.dim_metadata_size);
238     for (int i = 0; i < sparsity.dim_metadata_size; i++) {
239         format[i] = sparsity.dim_metadata[i].format;
240         dense_size[i] = sparsity.dim_metadata[i].dense_size;
241         segments[i] =
242             TfLiteIntArrayToVector(sparsity.dim_metadata[i].array_segments);
243         indices[i] = TfLiteIntArrayToVector(sparsity.dim_metadata[i].array_indices);
244     }
245
246     InitSparseToDenseConverter(shape, std::move(traversal_order),
247                               std::move(format), std::move(dense_size),
248                               std::move(segments), std::move(indices),
249                               std::move(block_map));
250 }

```

...

```

251
252 template <typename T>
253 void FormatConverter<T>::InitSparseToDenseConverter(
254     std::vector<int> shape, std::vector<int> traversal_order,
255     std::vector<TfLiteDimensionType> format, std::vector<int> dense_size,
256     std::vector<std::vector<int>> segments,
257     std::vector<std::vector<int>> indices, std::vector<int> block_map) {
258     dense_shape_ = std::move(shape);
259     traversal_order_ = std::move(traversal_order);
260     block_map_ = std::move(block_map);
261     format_ = std::move(format);
262
263     dense_size_ = 1;
264     for (int i = 0; i < dense_shape_.size(); i++) {
265         dense_size_ *= dense_shape_[i];
266     }
267
268     dim_metadata_.resize(2 * format_.size());
269     for (int i = 0; i < format_.size(); i++) {
270         if (format_[i] == kTfLiteDimDense) {
271             dim_metadata_[2 * i] = {dense_size[i]};
272         } else {
273             dim_metadata_[2 * i] = std::move(segments[i]);
274             dim_metadata_[2 * i + 1] = std::move(indices[i]);

```

```

275     }
276 }
277
278 int original_rank = dense_shape_.size();
279 int block_dim = 0;
280
281 blocked_shape_.resize(original_rank);
282 block_size_.resize(block_map_.size());
283 for (int i = 0; i < original_rank; i++) {
284     if (block_dim < block_map_.size() && block_map_[block_dim] == i) {
285         int orig_dim = traversal_order_[original_rank + block_dim];
286         block_size_[block_dim] = dense_size[orig_dim];
287         blocked_shape_[i] = dense_shape_[i] / dense_size[orig_dim];
288         block_dim++;
289     } else {
290         blocked_shape_[i] = dense_shape_[i];
291     }
292 }
293 }
294
295 template <typename T>
296 void FormatConverter<T>::Populate(const T* src_data, std::vector<int> indices,
297                                 int level, int prev_idx, int* src_data_ptr,
298                                 T* dest_data) {
299     if (level == indices.size()) {
300         int orig_rank = dense_shape_.size();
301         std::vector<int> orig_idx;
302         orig_idx.resize(orig_rank);
303         int i = 0;
304         for (; i < orig_idx.size(); i++) {
305             int orig_dim = traversal_order_[i];
306             orig_idx[orig_dim] = indices[i];
307         }
308
309         for (; i < indices.size(); i++) {
310             const int block_idx = traversal_order_[i] - orig_rank;
311             const int orig_dim = block_map_[block_idx];
312             orig_idx[orig_dim] =
313                 orig_idx[orig_dim] * block_size_[block_idx] + indices[i];
314         }
315
316         dest_data[GetFlattenedIndex(orig_idx, dense_shape_)] =
317             src_data[*src_data_ptr];
318
319         *src_data_ptr = *src_data_ptr + 1;
320         return;
321     }
322
323     const int metadata_idx = 2 * level;

```

```

324     const int shape_of_level = dim_metadata_[metadata_idx][0];
325     if (format_[level] == kTfLiteDimDense) {
326         for (int i = 0; i < shape_of_level; i++) {
327             indices[level] = i;
328             Populate(src_data, indices, level + 1, prev_idx * shape_of_level + i,
329                     src_data_ptr, dest_data);
330         }
331     } else {
332         const auto& array_segments = dim_metadata_[metadata_idx];
333         const auto& array_indices = dim_metadata_[metadata_idx + 1];
334         for (int i = array_segments[prev_idx]; i < array_segments[prev_idx + 1];
335             i++) {
336             indices[level] = array_indices[i];
337             Populate(src_data, indices, level + 1, i, src_data_ptr, dest_data);
338         }
339     }
340 }
341
342 template <typename T>
343 TfLiteStatus FormatConverter<T>::SparseToDense(const T* src_data) {
344     data_.resize(dense_size_);
345     std::fill(data_.begin(), data_.end(), T(0));
346
347     int total_rank = traversal_order_.size();
348     int src_data_ptr = 0;
349     std::vector<int> indices(total_rank);
350     Populate(src_data, indices, 0, 0, &src_data_ptr, data_.data());
351
352     return kTfLiteOk;
353 }
354
355 template <typename T>
356 TfLiteStatus FormatConverter<T>::SparseToDense(const T* src_data,
357                                                const size_t dest_size,
358                                                T* dest_data,
359                                                TfLiteContext* context) {
360     if (dest_size != dense_size_) {
361         TF_LITE_MAYBE_KERNEL_LOG(
362             context, "unexpected buffer size for densified data, expected %lld.\n",
363             dense_size_);
364         return kTfLiteError;
365     }
366
367     // For types like Eigen::half, we cannot do a simple memset() with 0 values.
368     for (auto i = 0; i < dest_size; i++) {
369         dest_data[i] = T(0);
370     }
371
372     const int total_rank = traversal_order_.size();

```



```
373     int src_data_ptr = 0;
374     std::vector<int> indices(total_rank);
375     Populate(src_data, indices, 0, 0, &src_data_ptr, dest_data);
376
377     return kTfLiteOk;
378 }
379
380 template <typename T>
381 bool FormatConverter<T>::IsZero(const T val) {
382     return (val == static_cast<T>(0));
383 }
384
385 template class FormatConverter<int32_t>;
386 template class FormatConverter<int8_t>;
387 template class FormatConverter<float>;
388 template class FormatConverter<Eigen::half>;
389
390 } // namespace sparsity
391 } // namespace internal
392 } // namespace tfLite
```