

5100e359ae ▾

...

tensorflow / tensorflow / core / kernels / dequantize_op.cc

 **jpienaar** Rename to underlying type rather than alias ... ✕

 History

7 contributors



293 lines (265 sloc) | 11.8 KB

...

```

1  /* Copyright 2015 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
16 // See docs in ../ops/math_ops.cc.
17
18 #define EIGEN_USE_THREADS
19
20 #include "tensorflow/core/framework/op.h"
21 #include "tensorflow/core/framework/op_kernel.h"
22 #include "tensorflow/core/framework/type_traits.h"
23 #include "tensorflow/core/framework/types.h"
24 #include "tensorflow/core/kernels/meta_support.h"
25 #include "tensorflow/core/kernels/quantization_utils.h"
26 #include "tensorflow/core/lib/core/errors.h"
27 #include "tensorflow/core/platform/bfloat16.h"
28
29 namespace {
```

[illegible]

```

79     }
80
81     if (mode_string == "MIN_COMBINED") {
82         mode_ = QUANTIZE_MODE_MIN_COMBINED;
83     } else if (mode_string == "MIN_FIRST") {
84         mode_ = QUANTIZE_MODE_MIN_FIRST;
85     } else if (mode_string == "SCALED") {
86         mode_ = QUANTIZE_MODE_SCALED;
87     }
88     OP_REQUIRES_OK(ctx, ctx->GetAttr("narrow_range", &narrow_range_));
89     OP_REQUIRES_OK(ctx, ctx->GetAttr("axis", &axis_));
90 }

```

```

91
92 void Compute(OpKernelContext* ctx) override {
93     const Tensor& input = ctx->input(0);
94     const Tensor& input_min_tensor = ctx->input(1);
95     const Tensor& input_max_tensor = ctx->input(2);
96
97     int num_slices = 1;
98     if (axis_ > -1) {
99         num_slices = input.dim_size(axis_);
100     }
101     OP_REQUIRES(ctx, input_min_tensor.NumElements() == num_slices,
102                 errors::InvalidArgument(
103                     "input_min_tensor must have as many elements as input on "
104                     "the dequantization axis (",
105                     axis_, "), got ", input_min_tensor.NumElements(),
106                     ", expected ", num_slices));
107     OP_REQUIRES(ctx, input_max_tensor.NumElements() == num_slices,
108                 errors::InvalidArgument(
109                     "input_max_tensor must have as many elements as input on "
110                     "the dequantization axis (",
111                     axis_, "), got ", input_max_tensor.NumElements(),
112                     ", expected ", num_slices));
113
114     Tensor* output = nullptr;
115     OP_REQUIRES_OK(ctx, ctx->allocate_output(0, input.shape(), &output));
116     Tensor float_output =
117         need_cast_ ? tensorflow::Tensor(DT_FLOAT, input.shape()) : *output;
118     if (num_slices == 1) {
119         const float min_range = input_min_tensor.flat<float>()(0);
120         const float max_range = input_max_tensor.flat<float>()(0);
121         DequantizeTensor(ctx, input, min_range, max_range, &float_output);
122     } else {
123         OP_REQUIRES(ctx, mode_ != QUANTIZE_MODE_MIN_FIRST,
124                     errors::Unimplemented("MIN_FIRST mode is not implemented for "
125                                             "Dequantize with axis != -1."));
126
127         int64_t pre_dim = 1, post_dim = 1;

```

```

128     for (int i = 0; i < axis_; ++i) {
129         pre_dim *= float_output.dim_size(i);
130     }
131     for (int i = axis_ + 1; i < float_output.dims(); ++i) {
132         post_dim *= float_output.dim_size(i);
133     }
134     auto input_tensor = input.template bit_casted_shaped<T, 3>(
135         {pre_dim, num_slices, post_dim});
136     auto output_tensor =
137         float_output.flat_inner_outer_dims<float, 3>(axis_ - 1);
138     auto min_ranges = input_min_tensor.vec<float>();
139     auto max_ranges = input_max_tensor.vec<float>();
140     for (int i = 0; i < num_slices; ++i) {
141         DequantizeSlice(ctx->eigen_device<Device>(), ctx,
142             input_tensor.template chip<1>(i), min_ranges(i),
143             max_ranges(i), output_tensor.template chip<1>(i));
144     }
145 }
146 if (need_cast_) {
147     S* out_ptr = output->flat<S>().data();
148     float* in_ptr = float_output.flat<float>().data();
149     for (int64_t i = 0; i < float_output.NumElements(); ++i) {
150         out_ptr[i] = static_cast<S>(in_ptr[i]);
151     }
152 }
153 }
154
155 void DequantizeTensor(OpKernelContext* ctx, const Tensor& input,
156     const float min_range, const float max_range,
157     Tensor* output) {
158     const float half_range =
159         !std::is_signed<T>::value
160         ? 0.0f
161         : (static_cast<float>(std::numeric_limits<T>::max()) -
162             std::numeric_limits<T>::min() + 1) /
163             2.0f;
164
165     if (mode_ == QUANTIZE_MODE_MIN_COMBINED) {
166         const float scale_factor =
167             (max_range - min_range) /
168             (static_cast<float>(std::numeric_limits<T>::max()) -
169             std::numeric_limits<T>::min());
170
171         const auto& input_tensor = input.flat<T>();
172         output->flat<float>() =
173             ((input_tensor.template cast<float>() + half_range) * scale_factor) +
174             min_range;
175
176     } else if (mode_ == QUANTIZE_MODE_MIN_FIRST) {

```

```

177     if (meta::IsSupportedAndEnabled() && std::is_same<T, quint8>()) {
178         auto input_ui8_array = input.flat<quint8>();
179         meta::Dequantize(ctx, input_ui8_array.data(), input_ui8_array.size(),
180             min_range, max_range, output->flat<float>().data());
181     } else {
182         QuantizedTensorToFloatInPlaceUsingEigen<T>(
183             ctx->template eigen_device<Device>(), input, min_range, max_range,
184             output);
185     }
186 } else if (mode_ == QUANTIZE_MODE_SCALED) {
187     const int min_output_value =
188         std::numeric_limits<T>::min() + (narrow_range_ ? 1 : 0);
189     const float scale_factor =
190         std::numeric_limits<T>::min() == 0
191             ? (max_range / std::numeric_limits<T>::max())
192             : std::max(min_range / min_output_value,
193                 max_range / std::numeric_limits<T>::max());
194     const auto& input_tensor = input.flat<T>();
195     output->flat<float>() =
196         input_tensor.template cast<int>().template cast<float>() *
197         scale_factor;
198 }
199 }
200
201 template <typename ConstVec, typename Vec>
202 void DequantizeSlice(const Device& d, OpKernelContext* ctx,
203     const ConstVec& input, float min_range, float max_range,
204     Vec output) {
205     // TODO(pauldonnelly): Factor out the similar calculations in quantize,
206     // dequantize and quantize_and_dequantize ops.
207     const float half_range =
208         !std::is_signed<T>::value
209             ? 0.0f
210             : (static_cast<float>(std::numeric_limits<T>::max()) -
211                 std::numeric_limits<T>::min() + 1) /
212                 2.0f;
213
214     if (mode_ == QUANTIZE_MODE_MIN_COMBINED) {
215         const float scale_factor =
216             (max_range - min_range) /
217             (static_cast<float>(std::numeric_limits<T>::max()) -
218                 std::numeric_limits<T>::min());
219
220         output.device(d) =
221             ((input.template cast<float>() + half_range) * scale_factor) +
222             min_range;
223     } else if (mode_ == QUANTIZE_MODE_SCALED) {
224         const int min_output_value =
225             std::numeric_limits<T>::min() + (narrow_range_ ? 1 : 0);

```

```

226     const float scale_factor =
227         std::numeric_limits<T>::min() == 0
228         ? (max_range / std::numeric_limits<T>::max())
229         : std::max(min_range / min_output_value,
230                   max_range / std::numeric_limits<T>::max());
231     output.device(d) = input.template cast<float>() * scale_factor;
232 }
233 }
234
235 private:
236     int mode_;
237     int axis_;
238     bool narrow_range_;
239     bool need_cast_;
240 };
241
242 REGISTER_KERNEL_BUILDER(Name("Dequantize")
243                         .Device(DEVICE_CPU)
244                         .TypeConstraint<quint8>("T")
245                         .TypeConstraint<float>("dtype"),
246                         DequantizeOp<CPUDevice, quint8, float>);
247 REGISTER_KERNEL_BUILDER(Name("Dequantize")
248                         .Device(DEVICE_CPU)
249                         .TypeConstraint<qint8>("T")
250                         .TypeConstraint<float>("dtype"),
251                         DequantizeOp<CPUDevice, qint8, float>);
252 REGISTER_KERNEL_BUILDER(Name("Dequantize")
253                         .Device(DEVICE_CPU)
254                         .TypeConstraint<quint16>("T")
255                         .TypeConstraint<float>("dtype"),
256                         DequantizeOp<CPUDevice, quint16, float>);
257 REGISTER_KERNEL_BUILDER(Name("Dequantize")
258                         .Device(DEVICE_CPU)
259                         .TypeConstraint<qint16>("T")
260                         .TypeConstraint<float>("dtype"),
261                         DequantizeOp<CPUDevice, qint16, float>);
262 REGISTER_KERNEL_BUILDER(Name("Dequantize")
263                         .Device(DEVICE_CPU)
264                         .TypeConstraint<qint32>("T")
265                         .TypeConstraint<float>("dtype"),
266                         DequantizeOp<CPUDevice, qint32, float>);
267
268 REGISTER_KERNEL_BUILDER(Name("Dequantize")
269                         .Device(DEVICE_CPU)
270                         .TypeConstraint<quint8>("T")
271                         .TypeConstraint<bfloat16>("dtype"),
272                         DequantizeOp<CPUDevice, quint8, bfloat16>);
273 REGISTER_KERNEL_BUILDER(Name("Dequantize")
274                         .Device(DEVICE_CPU)

```

```

275         .TypeConstraint<qint8>("T")
276         .TypeConstraint<bfloat16>("dtype"),
277         DequantizeOp<CPUDevice, qint8, bfloat16>);
278 REGISTER_KERNEL_BUILDER(Name("Dequantize")
279         .Device(DEVICE_CPU)
280         .TypeConstraint<qint16>("T")
281         .TypeConstraint<bfloat16>("dtype"),
282         DequantizeOp<CPUDevice, quint16, bfloat16>);
283 REGISTER_KERNEL_BUILDER(Name("Dequantize")
284         .Device(DEVICE_CPU)
285         .TypeConstraint<qint16>("T")
286         .TypeConstraint<bfloat16>("dtype"),
287         DequantizeOp<CPUDevice, qint16, bfloat16>);
288 REGISTER_KERNEL_BUILDER(Name("Dequantize")
289         .Device(DEVICE_CPU)
290         .TypeConstraint<qint32>("T")
291         .TypeConstraint<bfloat16>("dtype"),
292         DequantizeOp<CPUDevice, qint32, bfloat16>);
293 } // namespace tensorflow

```