# DRP-AI Extension Pack (Pruning Tool) Version 1.1.0

## User's Manual

## Contents

## 1.  Overview

This section describes the operating environment and functions of the DRP-AI Extension Pack.

## 1.1  Product Configuration

**Table 1.1    Product Configuration**

| Item | Description |
|---|---|
| r20ut5188ej0200-drp-ai-extension-pack.pdf | This manual |
| drpai-extension-pack_ver1.1.0.tar.gz | DRP-AI Extension Pack (product covered by this manual) |

**Table 1.2    Configuration of Files in drpai-extension-pack_ver1.1.0.tar.gz**

| Configuration of Files | Description |
|---|---|
| 🗁 drpai-extension-pack_ver1.1.0.tar.gz | |
| 🗁 drpai_compaction_tool | API library of the functions and class listed in Table 4.1 |
| 🗁 samples | |
| 🗁 classification | |
| 🗁 pytorch_mobilenetv2 | Sample code using MobileNetV2 of the PyTorch version. See 3.8.1 for details. |
| 🗁 tensorflow_cnn | Sample code using CNN of the TensorFlow version. See 3.8.2 for details. |

## 1.2   Operating Environment

The operating environment and software to be installed for the DRP-AI Extension Pack in each case are shown in the following tables.

**Table 1.3    Operating Environment (When Using PyTorch)**

| Item | Software Name and Version Number, etc. |
|---|---|
| Operating environment | Ubuntu 20.04 LTS, 64-bit version |
| | CUDA 11.6 |
| Software to be installed | Python 3.8.10 |
| | torch==1.13.1 |
| | torchvision==0.14.1 |
| | torchstat==0.0.7 |
| | pandas==1.4.2 |
| | onnx==1.11.0 |

**Table 1.4    Operating Environment (When Using TensorFlow)**

| Item | Software Name and Version Number, etc. |
|---|---|
| Operating environment | Ubuntu 20.04 LTS, 64-bit version |
| | CUDA 11.2 |
| Software to be installed | Python 3.8.10 |
| | tensorflow==2.5.0 |
| | tensorflow-model-optimization==0.6.0 |
| | tf2onnx==1.14.0 |
| | onnx==1.11.0 |

## 1.3 Function

The DRP-AI Extension Pack provides a pruning function optimized for the DRP-AI. A general description of pruning is given under 1.4, Pruning, on the following page. This pruning function optimized for the DRP-AI can be used by using the DRP-AI Extension Pack in combination with the training code written with the use of PyTorch or TensorFlow.



**Figure 1-1     Deployment Flow**

## 1.4   Pruning

Nodes are interconnected in a neural network as shown in the figure below. Methods of reducing the number of parameters by removing weights between nodes or removing nodes are referred to as "pruning". A neural network to which pruning has not been applied is generally referred to as a dense neural network.

Applying pruning to a neural network leads to a slight deterioration in the accuracy of the model but can reduce the power required by hardware and accelerate the inference process.



**Figure 1-2     Schematic View of the Pruning of a Neural Network**

Note:   In the use of this product, we recommend pruning by at least 70% to improve the processing performance of the DRP-AI.

## 1.5   Relationship between Compressing AI Models and DRP-AI Performance

DRP-AI for RZ/V2H supports the feature of efficiently calculating the pruned AI model. Therefore, power efficiency is improved by using the pruned AI model.

The following graph provides an example of improvement in power efficiency when changing from an unpruned AI model to a pruned AI model. Compared to unpruned AI models, pruned AI models are significantly more power-efficient.



**Figure 1-3      DRP-AI Performance after Compressing AI Models**

Note: Quantization was applied to the AI models and measurements were performed.

Applying compressing processing such as pruning and quantization to AI models might generally lower the accuracy of models. Using the DRP-AI Extension Pack in pruning, however, allows ensuring the same or almost the same accuracy for the AI model as that before pruning by proceeding with retraining after pruning. The figure below shows the results of changes in accuracy with the use of the YOLOv2 models. Compared with an accuracy of 74.9% before compressing, that after compressing (pruning plus quantization) can reach 72.3%.



**Figure 1-4      Changes in the Accuracy of YOLOv2 Models after Compressing**

Note:   For details on calibration and quantization, see the DRP-AI_Quantizer User's Manual.

## 1.6  Two Pruning Modes

This product supports two pruning modes. One-shot pruning is characterized by a relatively short training time being associated with pruning. In gradual pruning, longer training times are associated with pruning than in one-shot pruning, but the accuracy may be improved. Table 1.5 shows a comparison of the two pruning modes. The pruning rate rises as training proceeds in gradual pruning. For details on how to set one-shot pruning or gradual pruning, see 4.2.2 or 4.3.2, depending on whether you are using PyTorch or TensorFlow, respectively. Note that the initial application of one-shot pruning is recommended.

**Table 1.5    Comparison of the Two Pruning Modes**

| One-shot pruning (recommended) | Pruning is applied only once as the initial step.  |
|---|---|
| Gradual pruning | Pruning is applied gradually.  |

## 1.7   Updates in Version 1.1.0

The updates in DRP-AI Extension Pack V1.1.0 are as follows.

✓   Supported PyTorch models with a Multi-Head Attention structure for a Transformer neural network.

For models using torch.nn.MultiheadAttention(), please note that there are usage considerations. For more information, please refer to 6. Usage Notes.

o   PyTorch torch.nn.MultiheadAttention() :
https://pytorch.org/docs/1.13/generated/torch.nn.MultiheadAttention.html#multiheadattention

## 2.  Setting Up the DRP-AI Extension Pack

This section describes how to set up the DRP-AI Extension Pack. Descriptions in this section are on the assumption that Python 3.8.10 has been set up on a PC running Ubuntu.

### 2.1   Installing the Library for Use by the DRP-AI Extension Pack

Install the following library on a PC running Ubuntu.

[When using PyTorch]

```
$ pip3 install torch==1.13.1+cu116 torchvision==0.14.1+cu116 \
                 --extra-index-url https://download.pytorch.org/whl/cu116
$ pip3 install torchstat==0.0.7 pandas==1.4.2 onnx==1.11.0
```

[When using TensorFlow]

```
$ pip3 install tensorflow==2.5.0 tensorflow-model-optimization==0.6.0 tf2onnx==1.14.0 onnx==1.11.0
```

### 2.2   Adding the Environment Variable

Register the working directory as an environment variable.

```
$ export WORK=/home/<Path to working directory>
```

Note:  Change <Path to working directory> to suit the environment of the PC you are using.

### 2.3   Decompressing the DRP-AI Extension Pack

Place drpai-extension-pack_ver*.tar.gz in the working directory and execute the following command.

```
$ cd $WORK
$ tar -xvf drpai-extension-pack_ver*.tar.gz
[When using TensorFlow, also execute the following command.]
$ drpai_compaction_tool/scripts/setup_tf.sh
```

### 2.4   Adding the Path to the DRP-AI Extension Pack

Execute the following command to add the path to the DRP-AI Extension Pack.

```
$ cd $WORK
$ export PYTHONPATH="$(pwd):$PYTHONPATH"
```

Note:  Once you have ended the terminal session, re-execute the command stated above when you intend to use the extension pack again.

Execute the following command. With output of the version number, the setup processing is completed.

```
$ python3 -c "import drpai_compaction_tool; print(drpai_compaction_tool.__version__)"
<DRP-AI Extension Pack version>
[When using TensorFlow, ensure that the following commands do not generate any errors.]
$ python3 -c "from drpai_compaction_tool.tensorflow import Pruner"
```

Note:  <DRP-AI Extension Pack version> depends on the version you are using.

## 3.　Using the DRP-AI Extension Pack

This section describes how to use the DRP-AI Extension Pack.

### 3.1　Flow of Using the DRP-AI Extension Pack

Use the DRP-AI Extension Pack in combination with the training code written with the use of PyTorch or TensorFlow. Figure 3-1 shows the flow of using the DRP-AI Extension Pack.
The flow consists of two steps. The first step is initial training. Initial training involves training of the AI model without pruning. Use the code for use in initial training and a dataset you have prepared.

The second step is pruning and then retraining. This includes retraining of the AI model by adding the DRP-AI Extension Pack to the code for use in initial training. For details on how to add the DRP-AI Extension Pack to the code for use in initial training, see 3.2, [PyTorch] Adding the DRP-AI Extension Pack. Check the accuracy of the AI model after one round of pruning then retraining has been completed. Repeat pruning then retraining with increasingly high pruning rates while confirming that the rates in use do not create problems in terms of accuracy.
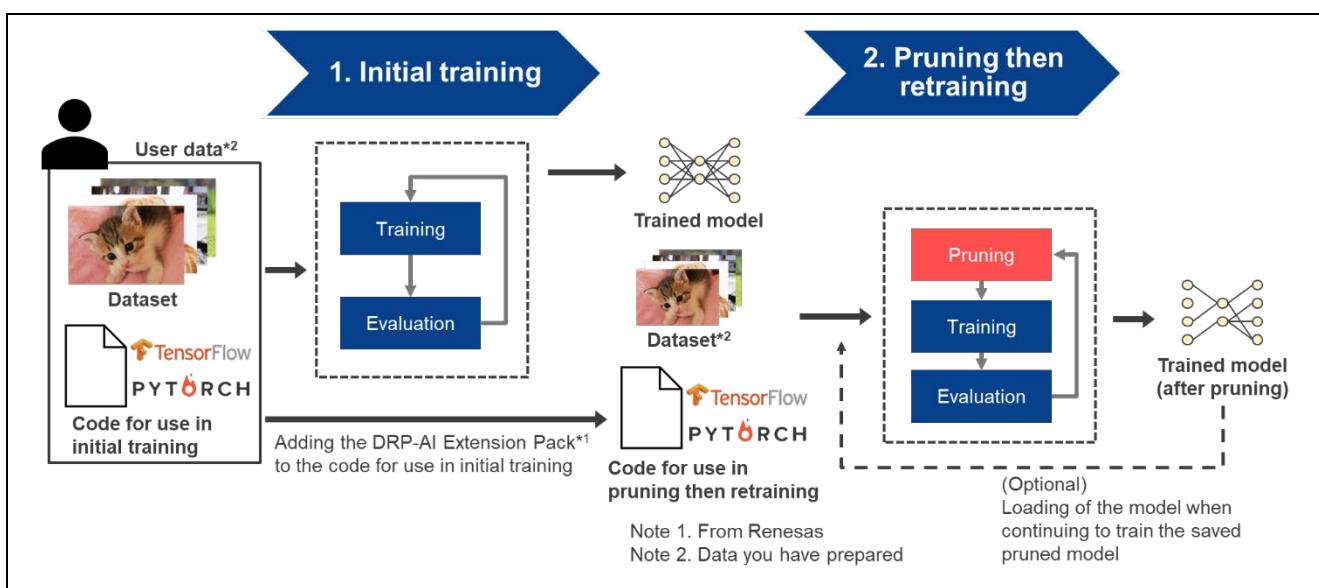


**Figure 3-1　　Flow of Using the DRP-AI Extension Pack**

Figure 3-2 consists of listings of the code written with PyTorch for use in initial training without and with addition of the DRP-AI Extension Pack. The code in the left column is that for initial training and the code in the right column is that for pruning then retraining. The green shading indicates the differences between the two listings, that is, the several lines that are added to make the DRP-AI Extension Pack usable. For details on the case of using PyTorch, see 3.2. For details on the case of using TensorFlow, see 3.5.

Left column (Code for Initial Training):

```
1  # Importing the library
2  import torch
3  from torch import nn
4  from torch.optim import SGD
5  from torch.utils.data import DataLoader
6  from torchvision import datasets
7  from torchvision.transforms import ToTensor


8
9  # Defining the neural network
10 class NeuralNetwork(nn.Module):
11     def __init__(self):
12         super(NeuralNetwork, self).__init__()
13         self.flatten = nn.Flatten()
14         self.linear_relu_stack = nn.Sequential(
15             nn.Linear(28*28, 512),
16             nn.ReLU(),
17             nn.Linear(512, 512),
18             nn.ReLU(),
19             nn.Linear(512, 10),
20             nn.ReLU()
21         )
22
23     def forward(self, x):
24         x = self.flatten(x)
25         logits = self.linear_relu_stack(x)
26         return logits
27 model = NeuralNetwork()


28
29 # Registering the model parameters with the optimizer
30 optimizer = SGD(model.parameters(), lr=1e-3)
31
32 # Defining the training data and loss function
33 training_data = datasets.FashionMNIST(
34     root="data",
35     train=True,
36     download=True,
37     transform=ToTensor(),
38 )
39 loss_fn = nn.CrossEntropyLoss()
40 batch_size = 64
41 max_epochs = 10
42



43 # Training
44 for epoch in range(max_epochs):
45     for batch_x, batch_y in DataLoader(training_data, batch_size):



46         # Compute prediction and loss
47         pred = model(batch_x)
48         loss = loss_fn(pred, batch_y)
49
50         # Backpropagation
51         optimizer.zero_grad()
52         loss.backward()
53         optimizer.step()
54
55 torch.save(model.state_dict(), "pretrained_model.pth")
```

Right column (Code for Pruning Then Retraining):

```
1  # Importing the library
2  import torch
3  from torch import nn
4  from torch.optim import SGD
5  from torch.utils.data import DataLoader
6  from torchvision import datasets
7  from torchvision.transforms import ToTensor
8  # 1. Importing the DRP-AI Extension Pack module
9  from drpai_compaction_tool.pytorch import make_pruning_layer_list, \
10                                           Pruner, \
11                                           get_model_info
12
13 # Defining the neural network
14 class NeuralNetwork(nn.Module):
15     def __init__(self):
16         super(NeuralNetwork, self).__init__()
17         self.flatten = nn.Flatten()
18         self.linear_relu_stack = nn.Sequential(
19             nn.Linear(28*28, 512),
20             nn.ReLU(),
21             nn.Linear(512, 512),
22             nn.ReLU(),
23             nn.Linear(512, 10),
24             nn.ReLU()
25         )
26
27     def forward(self, x):
28         x = self.flatten(x)
29         logits = self.linear_relu_stack(x)
30         return logits
31 model = NeuralNetwork()
32 # 2. Loading the trained model
33 model.load_state_dict(torch.load("pretrained_model.pth"))
34
35 # Registering the model parameters with the optimizer
36 optimizer = SGD(model.parameters(), lr=1e-3)
37
38 # Defining the training data and loss function
39 training_data = datasets.FashionMNIST(
40     root="data",
41     train=True,
42     download=True,
43     transform=ToTensor(),
44 )
45 loss_fn = nn.CrossEntropyLoss()
46 batch_size = 64
47 max_epochs = 10
48 # 3. Preparing for pruning the model
49 pruning_layer_list = make_pruning_layer_list(model, [(1,1,28,28)])
50 pruner = Pruner(model,
51                 pruning_layer_list,
52                 final_pr = 0.7)
53 print(get_model_info(model, [(1,1,28,28)]))
54 # Training
55 for epoch in range(max_epochs):
56     for batch_x, batch_y in DataLoader(training_data, batch_size):
57         # 4. Updating pruning parameters
58         pruner.update()
59
60         # Compute prediction and loss
61         pred = model(batch_x)
62         loss = loss_fn(pred, batch_y)
63
64         # Backpropagation
65         optimizer.zero_grad()
66         loss.backward()
67         optimizer.step()
68
69 # 5. Saving the pruned model
70 if pruner.is_finished:
71     torch.save(model.state_dict(), "oneshot_pruned_model.pth")
72     torch.onnx.export(model,
73                       training_data[0][0].unsqueeze(0),
74                       "oneshot_pruned_model.onnx",
75                       opset_version = 12)
```

**Figure 3-2    How to Add the DRP-AI Extension Pack to the Code for Use in Initial Training (Left: Code for Initial Training; Right: Code for Pruning Then Retraining)**

RENESAS

## 3.2   [PyTorch] Adding the DRP-AI Extension Pack

The steps involved in adding the DRP-AI Extension Pack to the code written with PyTorch for use in initial training and then proceeding with pruning and retraining are given below. Implement the five processes listed below in the code for initial training.

1. Importing the DRP-AI Extension Pack module
2. Loading the trained model
3. Preparing for pruning the model
4. Updating the pruning parameters
5. Saving the pruned model

The figure below shows a listing of the code for retraining, which is obtained by adding the DRP-AI Extension Pack to the code written with PyTorch for use in initial training. In this figure, red text indicates the statements added to the code for use in initial training.

```python
# Importing the library
import torch
from torch import nn
from torch.optim import SGD
from torch.utils.data import DataLoader
from torchvision import datasets
from torchvision.transforms import ToTensor
# 1. Importing the DRP-AI Extension Pack module
from drpai_compaction_tool.pytorch import make_pruning_layer_list, \
                                          Pruner, \
                                          get_model_info

# Defining the neural network
class NeuralNetwork(nn.Module):
    def __init__(self):
        super(NeuralNetwork, self).__init__()
        self.flatten = nn.Flatten()
        self.linear_relu_stack = nn.Sequential(
            nn.Linear(28*28, 512),
            nn.ReLU(),
            nn.Linear(512, 512),
            nn.ReLU(),
            nn.Linear(512, 10),
            nn.ReLU()
        )

    def forward(self, x):
        x = self.flatten(x)
        logits = self.linear_relu_stack(x)
        return logits
model = NeuralNetwork()
# 2. Loading the trained model
model.load_state_dict(torch.load("pretrained_model.pth"))

# Registering the model parameters with the optimizer
optimizer = SGD(model.parameters(), lr=1e-3)
```

```
37
38    # Defining the training data and loss function
39    training_data = datasets.FashionMNIST(
40        root="data",
41        train=True,
42        download=True,
43        transform=ToTensor(),
44    )
45    loss_fn = nn.CrossEntropyLoss()
46    batch_size = 64
47    max_epochs = 10
48    # 3. Preparing for pruning the model
49    pruning_layer_list = make_pruning_layer_list(model, [(1,1,28,28)])
50    pruner = Pruner(model,
51                    pruning_layer_list,
52                    final_pr=0.7)
53    print(get_model_info(model, [(1,1,28,28)]))
54    # Training
55    for epoch in range(max_epochs):
56        for batch_x, batch_y in DataLoader(training_data, batch_size):
57            # 4. Updating the pruning parameters
58            pruner.update()
59
60            # Compute prediction and loss
61            pred = model(batch_x)
62            loss = loss_fn(pred, batch_y)
63
64            # Backpropagation
65            optimizer.zero_grad()
66            loss.backward()
67            optimizer.step()
68
69    # 5. Saving the pruned model
70    if pruner.is_finished
71        torch.save(pruner.state_dict(), "pruned_model.pth")
72        torch.onnx.export(model,
73                        training_data[0][0].unsqueeze(0),
74                        'pruned_model.onnx',
75                        opset_version = 12)
```

**Figure 3-3    Training Code for Pruning Then Retraining**

### 3.2.1　[PyTorch] Importing the DRP-AI Extension Pack Module

Import the DRP-AI Extension Pack module to the code written with PyTorch for use in initial training.

```
1   # Importing the library
2   import torch
3   from torch import nn
4   from torch.optim import SGD
5   from torch.utils.data import DataLoader
6   from torchvision import datasets
7   from torchvision.transforms import ToTensor
8   # 1. Importing the DRP-AI Extension Pack module
9   from drpai_compaction_tool.pytorch import make_pruning_layer_list, \
10                                            Pruner \
11                                            get_model_info
```

**Figure 3-4　　Importing the DRP-AI Extension Pack Module**

### 3.2.2　[PyTorch] Loading the Trained Model

Define the model and load the trained model.

```
31  model = NeuralNetwork()
32  # 2. Loading the trained model
33  model.load_state_dict(torch.load("pretrained_model.pth"))
34
```

**Figure 3-5　　Loading the Trained Model**

### 3.2.3 [PyTorch] Preparing for Pruning the Model

Register the model parameters with the optimizer and then execute the API function for pruning. After the API function for pruning has been executed, confirming that pruning has been performed with the get_model_info() function is recommended.

```python
35  # Registering the model parameters with the optimizer
36  optimizer = SGD(model.parameters(), lr=1e-3)
37
38  # Defining the training data and loss function
39  training_data = datasets.FashionMNIST(
40      root="data",
41      train=True,
42      download=True,
43      transform=ToTensor(),
44  )
45  loss_fn = nn.CrossEntropyLoss()
46  batch_size = 64
47  max_epochs = 10
48  # 3. Preparing for pruning the model
49  pruning_layer_list = make_pruning_layer_list(model, [(1,1,28,28)])
50  pruner = Pruner(model,
51                  pruning_layer_list,
52                  final_pr=0.7)
53  print(get_model_info(model, [(1,1,28,28)]))
```

**Figure 3-6　　Pruning the Model**

Note: Execute the API function for pruning (Pruner) after registering the model parameters with the optimizer.

### 3.2.4 [PyTorch] Updating the Pruning Parameters

Update the pruning parameters during training. The API function in red text below (pruner.update()) must be called at the start of each iteration.

```python
54  # Training
55  for epoch in range(max_epochs):
56      for batch_x, batch_y in DataLoader(training_data, batch_size):
57          # 4. Updating the pruning parameters
58          pruner.update()
59
60          # Compute prediction and loss
61          pred = model(batch_x)
62          loss = loss_fn(pred, batch_y)
63
64          # Backpropagation
65          optimizer.zero_grad()
66          loss.backward()
67          optimizer.step()
```

**Figure 3-7　　Updating the Pruning Parameters**

### 3.2.5 [PyTorch] Saving the Pruned Model

After having confirmed the completion of pruning, use the PyTorch method to save the pruned model. If the model is to be exported to ONNX, specify 12 for opset_version.

```python
71  # 5. Saving the pruned model
72  if pruner.is_finished:
73      torch.save(pruner.state_dict(), "pruned_model.pth")
74      torch.onnx.export(model,
75                        training_data[0][0].unsqueeze(0),
76                        'pruned_model.onnx',
77                        opset_version = 12)
```

**Figure 3-8　　Saving the Pruned Model**

## 3.3    [PyTorch] Confirming the Result of Pruning

The steps involved in confirming the result of pruning are given below. The function (get_model_info) provided by the DRP-AI Extension Pack can be used to confirm how many parameters were pruned in which layers and the reductions in the number of multiply-and-accumulate calculations. For details on how to use the get_model_info function, see 4.2.5. Calling this function is possible both before and after pruning.

The figure below shows a listing of the sample code in one-shot pruning. In one-shot pruning, the confirmation of pruning being applied before training is recommended.

```
1    ・・・・(Omitted)
2    # 3. Preparing for pruning the model
3    pruning_layer_list = make_pruning_layer_list(model, [(1,1,28,28)])
4    pruner = Pruner(model,
5                    pruning_layer_list,
6                    final_pr = 0.7)
7    # Recommended: Confirming the result of pruning before training
8    print(get_model_info(model, [(1,1,28,28)]))
9
10   # Training
11   for epoch in range(max_epochs):
12   ・・・・(Omitted)
```

**Figure 3-9      Confirming the Result of Pruning: One-Shot Pruning**

The figure below shows a listing of the sample code in gradual pruning. In gradual pruning, the confirmation of pruning being applied during training is recommended.

```
1    ・・・・(Omitted)
2    # 3. Preparing for pruning the model
3    pruning_layer_list = make_pruning_layer_list(model, [(1,1,28,28)])
4    end_step = get_endstep(data_loader,
5                           max_epoch=max_epochs)
6    frequency = get_frequency(dataloader=data_loader)
7    pruner = Pruner(model,
8                    pruning_layer_list,
9                    final_pr = 0.7,
10                   end_step=end_step,
11                   frequency=frequency)
12
13   # Training
14   for epoch in range(max_epochs):
15       for i, (batch_x, batch_y) in enumerate(data_loader):
16           # 4. Updating the pruning parameters
17           pruner.update()
18           # Recommended: Confirming the result of pruning during training
19           if i % 100 == 0:
20               print(get_model_info(model, [(1,1,28,28)]))
21   ・・・・(Omitted)
```

**Figure 3-10      Confirming the Result of Pruning: Gradual Pruning**

The figure below shows the result of executing get_model_info. The meanings of the headings in the figure are as follows.

"module name": Layer name
"input shape": Input size to a layer
"output shape": Output size from a layer
"params": Number of parameters
"sparsity": Pruning rate
"Baseline MAC": Number of multiply-and-accumulate calculations before pruning
"Current MAC": Number of multiply-and-accumulate calculations after pruning

The result of pruning shown below indicates that pruning by about 70% was applied in the "linear_relu_stack.2" layer. It also indicates that pruning reduces the number of multiply-and-accumulate calculations from 262,144 before pruning to 78,848 after pruning.

| | module name | input shape | output shape | params | sparsity | Baseline MAC | Current MAC |
|---|---|---|---|---|---|---|---|
| 0 | linear_relu_stack.0 | 784 | 512 | 401920.0 | 0.000 | 401,408.0 | 401,408.0 |
| 1 | linear_relu_stack.2 | 512 | 512 | 262656.0 | **0.699** | **262,144.0** | **78,848.0** |
| 2 | linear_relu_stack.4 | 512 | 10 | 5130.0 | 0.000 | 5,120.0 | 5,120.0 |
| total | | | | 669706.0 | | 668,672.0 | 485,376.0 |

**Figure 3-11    Result of Executing get_model_info**

## 3.4   [PyTorch] Training or Inference with a Saved Pruned Model

The steps involved in loading a saved pruned model are given below. Refer to this section when continuing to train a saved pruned model or performing inference with a pruned model. There are 2 methods to load the saved pruned model. The method by using load_pruned_state_dict() is recommended because it is easy to use.

### 3.4.1   [Recommend] How to load the pruned model with load_pruned_state_dict()

Load the saved pruned model through the steps listed below.

1.  Importing the DRP-AI Extension Pack module
2.  Loading the pruned model

Calling the get_model_info function to check the pruning rate after loading of the pruned model is recommended. For details on how to use the get_model_info function, see 4.2.5.

The figure below shows a list of the sample code.

```python
# Importing the library
import torch
from torch import nn
# 1. Importing the DRP-AI Extension Pack module
from drpai_compaction_tool.pytorch import load_pruned_state_dict, \
                                          get_model_info

# Defining the neural network
class NeuralNetwork(nn.Module):
    def __init__(self):
        super(NeuralNetwork, self).__init__()
        self.flatten = nn.Flatten()
        self.linear_relu_stack = nn.Sequential(
            nn.Linear(28*28, 512),
            nn.ReLU(),
            nn.Linear(512, 512),
            nn.ReLU(),
            nn.Linear(512, 10),
            nn.ReLU()
        )

    def forward(self, x):
        x = self.flatten(x)
        logits = self.linear_relu_stack(x)
        return logits
model = NeuralNetwork()


# 2. Loading the pruned model
load_pruned_state_dict(model, torch.load("pruned_model.pth"))
print(get_model_info(model, [(1, 1, 28, 28)]))
```

**Figure 3-12     [PyTorch] Loading the Pruned Model (Method1)**

### 3.4.2　How to load the pruned model with make_pruning_layer_list() and Pruner()

Loads the saved pruned model through the steps listed below.

1. Importing the DRP-AI Extension Pack module
2. Preparing for pruning the model with a pruning rate of 0.0
3. Loading the pruned model

Calling the get_model_info function to check the pruning rate after loading of the pruned model is recommended. For details on how to use the get_model_info function, see 4.2.5.

The figure below shows a list of the sample code.

```python
# Importing the library
import torch
from torch import nn
# 1. Importing the DRP-AI Extension Pack module
from drpai_compaction_tool.pytorch import make_pruning_layer_list, \
                                          Pruner, \
                                          get_model_info

# Defining the neural network
class NeuralNetwork(nn.Module):
    def __init__(self):
        super(NeuralNetwork, self).__init__()
        self.flatten = nn.Flatten()
        self.linear_relu_stack = nn.Sequential(
            nn.Linear(28*28, 512),
            nn.ReLU(),
            nn.Linear(512, 512),
            nn.ReLU(),
            nn.Linear(512, 10),
            nn.ReLU()
        )

    def forward(self, x):
        x = self.flatten(x)
        logits = self.linear_relu_stack(x)
        return logits
model = NeuralNetwork()

# 2. Preparing for pruning the model with a pruning rate of 0.0
pruning_layer_list = make_pruning_layer_list(model, [(1,1,28,28)])
pruner = Pruner(model, pruning_layer_list, final_pr=0.0)

# 3. Loading the pruned model
model.load_state_dict(torch.load("pruned_model.pth"), strict=True)
print(get_model_info(model, [(1, 1, 28, 28)]))
```

**Figure 3-13　　[PyTorch] Loading the Pruned Model (Method2)**

Note: When loading the weights with load_state_dict() function, set strict argument to "True". When this
argument is set to "False", weights may not be loaded correctly.

## 3.5  [TensorFlow] Adding the DRP-AI Extension Pack

The steps involved in adding the DRP-AI Extension Pack to the code written with TensorFlow for use in initial training and then proceeding with pruning and retraining are given below. Implement the five processes listed below in the code for initial training.

1. Importing the DRP-AI Extension Pack module
2. Loading the trained model
3. Preparing for pruning the model
4. Registering the callback function for pruning
5. Saving the pruned model

Figure 3-14 consists of listings of the code written with TensorFlow for use in initial training without and with addition of the DRP-AI Extension Pack. The code in the left column is that for initial training and the code in the right column is that for pruning then retraining. The green shading indicates the differences between the two listings, that is, the several lines that are added to make the DRP-AI Extension Pack usable.

```
1  # Importing the library                                       1  # Importing the library
2  import tensorflow as tf                                       2  import tensorflow as tf
3  import tf2onnx                                                3  import tf2onnx
4  import onnx                                                   4  import onnx
                                                                 5+ import tensorflow_model_optimization as tfmot
                                                                 6+ # 1. Importing the DRP-AI Extension Pack module
                                                                 7+ from drpai_compaction_tool.tensorflow import make_pruning_layer_list, \
                                                                 8+                                                   Pruner
5                                                                9
6  # Defining the neural network                                10  # Defining the neural network
7  def NeuralNetwork(input_shape=(32, 32, 3)):                  11  def NeuralNetwork(input_shape=(32, 32, 3)):
8      num_classes = 10                                         12      num_classes = 10
9      return tf.keras.Sequential([                             13      return tf.keras.Sequential([
10         tf.keras.layers.Flatten(input_shape=input_shape),    14         tf.keras.layers.Flatten(input_shape=input_shape),
11         tf.keras.layers.Dense(512, activation='relu', name="dense1"),  15    tf.keras.layers.Dense(512, activation='relu', name="dense1"),
12         tf.keras.layers.Dense(512, activation='relu', name="dense2"),  16    tf.keras.layers.Dense(512, activation='relu', name="dense2"),
13         tf.keras.layers.Dropout(0.4),                        17         tf.keras.layers.Dropout(0.4),
14         tf.keras.layers.Dense(num_classes, activation='softmax', name="dense3")  18  tf.keras.layers.Dense(num_classes, activation='softmax', name="dense3")
15     ])                                                       19     ])
16- model = NeuralNetwork()                                     20+ # 2. Loading the trained model
                                                                21+ model = tf.keras.models.load_model("pretrained_model.h5")
17                                                               22
18 # Defining the training data and loss function              23  # Defining the training data and loss function
19 (train_images, train_labels), (_, _) = tf.keras.datasets.cifar10.load_data()  24  (train_images, train_labels), (_, _) = tf.keras.datasets.cifar10.load_data()
20 train_images = train_images / 255.0                         25  train_images = train_images / 255.0
21 compile_args = {                                            26  compile_args = {
22     'optimizer': 'adam',                                    27      'optimizer': 'adam',
23     'loss': 'sparse_categorical_crossentropy',             28      'loss': 'sparse_categorical_crossentropy',
24     'metrics': ['accuracy'],                                29      'metrics': ['accuracy'],
25 }                                                           30  }
26 fit_args = {'epochs': 10,                                   31  fit_args = {'epochs': 10,
27            'batch_size': 64,                                32             'batch_size': 64,
28            'validation_split': 0.1}                         33             'validation_split': 0.1}
29                                                              34
                                                                35+ # 3. Preparing for pruning the model
                                                                36+ pruning_layer_list = make_pruning_layer_list(model)
                                                                37+ pruner = Pruner(model, pruning_layer_list, final_pr=0.7)
                                                                38+ model_for_pruning = pruner.get_pruning_model()
30 # Compiling the model.                                      39  # Compiling the model.
31- model.compile(**compile_args)                              40+ model_for_pruning.compile(**compile_args)
32                                                              41
33 # Training                                                  42  # Training
34- model.fit(train_images, train_labels, **fit_args)          43+ # 4. Registering the callback function for pruning
                                                                44+ callbacks = [
                                                                45+     # Update pruning parameters
                                                                46+     tfmot.sparsity.keras.UpdatePruningStep(),
                                                                47+     # Save pruning informations
                                                                48+     tfmot.sparsity.keras.PruningSummaries(log_dir="./log_dir"),
                                                                49+ ]
                                                                50+ model_for_pruning.fit(train_images, train_labels, **fit_args, callbacks=callbacks)
35                                                              51
36- # Saving                                                   52+ # 5. Saving the pruned model
37- model.save("pretrained_model.h5", include_optimizer=False)  53+ model_for_pruning.save("oneshot_pruned_model.h5", include_optimizer=True)
38 onnx_model, _ = tf2onnx.convert.from_keras(model, opset=12)  54  onnx_model, _ = tf2onnx.convert.from_keras(model, opset=12)
39- onnx.save(onnx_model, 'pretrained_model.onnx')              55+ onnx.save(onnx_model, 'oneshot_pruned_model.onnx')
```

**Figure 3-14      How to Add the DRP-AI Extension Pack to the Code for Use in Initial Training
(Left: Code for Initial Training; Right: Code for Pruning Then Retraining)**

### 3.5.1 [TensorFlow] Importing the DRP-AI Extension Pack Module

Import the DRP-AI Extension Pack module to the code written with TensorFlow for use in initial training.

```
1  # Importing the library
2  import tensorflow as tf
3  import tf2onnx
4  import onnx
5  import tensorflow_model_optimization as tfmot
6  # 1. Importing the DRP-AI Extension Pack module
7  from drpai_compaction_tool.tensorflow import make_pruning_layer_list, \
8                                               Pruner
```

**Figure 3-15    Importing the DRP-AI Extension Pack Module**

Note:  Import tensorflow_model_optimization.

### 3.5.2 [TensorFlow] Loading the Trained Model

Load the trained model according to the usage method of TensorFlow.

```
20  # 2. Loading the trained model
21  model = tf.keras.models.load_model("pretrained_model.h5")
```

**Figure 3-16    Loading the Trained Model**

### 3.5.3 [TensorFlow] Preparing for Pruning the Model

After having executed the API function for pruning, obtain the model to which pruning is to be applied by using the get_pruning_model() function. After that, execute compilation of the model.

```
35  # 3. Preparing for pruning the model
36  pruning_layer_list = make_pruning_layer_list(model)
37  pruner = Pruner(model, pruning_layer_list, final_pr=0.7)
38  model_for_pruning = pruner.get_pruning_model()
39  # Compiling the model.
40  model_for_pruning.compile(**compile_args)
```

**Figure 3-17    Preparing for Pruning the Model**

### 3.5.4  [TensorFlow] Registering the Callback Function for Pruning

Register the callback function for pruning when carrying out training. For details on UpdatePruningStep() and PruningSummaries(), see 3.6 and the official documents of TensorFlow.

```
42  # Training
43  # 4. Registering the callback function for pruning
44  callbacks = [
45      # Update pruning parameters
46      tfmot.sparsity.keras.UpdatePruningStep(),
47      # Save pruning informations
48      tfmot.sparsity.keras.PruningSummaries(log_dir="./log_dir"),
49  ]
50  model_for_pruning.fit(train_images,
                          train_labels,
                          **fit_args,
                          callbacks=callbacks)
```

**Figure 3-18     Registering the Callback Function for Pruning**

Note:   Only executing the step of preparing for pruning a model, which was described in 3.5.3, does not lead to actual pruning of the model. Make sure to always execute that step in combination with the callback function.

### 3.5.5  [TensorFlow] Saving the Pruned Model

Save the model according to the usage method of TensorFlow. If the model is to be exported to ONNX, specify 12 for opset.

```
52  # 4. Saving the pruned model
53  model_for_pruning.save("one-shot_pruned_model.h5", include_optimizer=True)
54  onnx_model, _ = tf2onnx.convert.from_keras(model, opset=12)
55  onnx.save(onnx_model, 'one-shot_pruned_model.onnx')
```

**Figure 3-19     Saving the Pruned Model**

## 3.6    [TensorFlow] Confirming the Result of Pruning

The steps involved in confirming the result of pruning are given below. The callback function (PruningSummaries()) provided by TensorFlow can be used to obtain the result of how many parameters were pruned in which layers. TensorBoard provided by TensorFlow can be used to display the obtained information in a way that allows confirming the result as shown below.
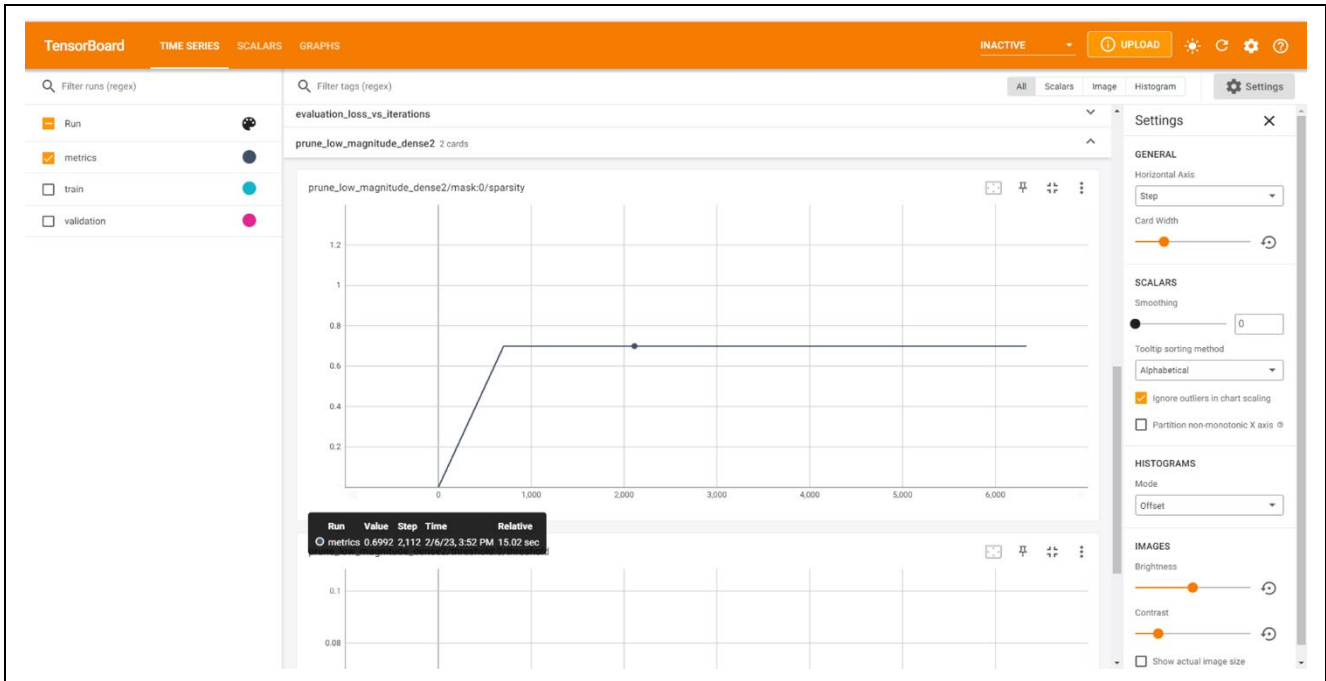


**Figure 3-20      Using TensorBoard to Confirm the Pruning Rate**

The changes in the pruning rate of the "prune_low_magnitude_dense2" layer are shown in the above figure. The horizontal axis indicates the number of steps (iterations) and the vertical axis indicates the pruning rate.

For example, PruningSummaries() may be set as follows:

```
tfmot.sparsity.keras.PruningSummaries(log_dir="./logdir")
```

The result of pruning can be confirmed by starting up TensorBoard as follows:

```
$ tensorboard –logdir ./logdir
```

For details, see the official documents of TensorFlow.

## 3.7 [TensorFlow] Training or Inference with a Saved Pruned Model

The steps involved in loading a saved pruned model are given below. Refer to this section when continuing to train a saved pruned model or performing inference with a pruned model.

Load the saved pruned model through the steps listed below.

1. Importing the DRP-AI Extension Pack module
2. Preparing for pruning the model with a pruning rate of 0.0
3. Loading the pruned model

The figure below shows a listing of the sample code.

```python
# Importing the library
import tensorflow as tf
import tensorflow_model_optimization as tfmot
# 1. Importing the DRP-AI Extension Pack module
from drpai_compaction_tool.tensorflow import make_pruning_layer_list, \
                                            Pruner


def print_sparsity(model):
    import numpy as np
    from tensorflow_model_optimization.python.core.sparsity.keras \
                                            import pruning_wrapper

    layer_info = {}
    for layer in model.layers:
        if not isinstance(layer, pruning_wrapper.PruneLowMagnitude):
            continue
        for weight, mask, threshold in layer.pruning_vars:
            np_mask = tf.keras.backend.get_value(mask)
            sparsity = 1.0 - np.count_nonzero(np_mask) / float(np_mask.size)
            layer_info[layer.name] = sparsity

    max_len = len(max(layer_info.keys(), key=lambda name: len(name)))
    for name, sparsity in layer_info.items():
        print(f'{name:{max_len+1}s} |  {sparsity:0.2f}')

# Defining the neural network
def NeuralNetwork(input_shape=(32, 32, 3)):
    num_classes = 10
    return tf.keras.Sequential([
        tf.keras.layers.Flatten(input_shape=input_shape),
        tf.keras.layers.Dense(512, activation='relu', name="dense1"),
        tf.keras.layers.Dense(512, activation='relu', name="dense2"),
        tf.keras.layers.Dropout(0.4),
        tf.keras.layers.Dense(num_classes,
                              activation='softmax', name="dense3")
    ])
model = NeuralNetwork()

# 2. Preparing for pruning the model with a pruning rate of 0.0
pruning_layer_list = make_pruning_layer_list(model)
```

```
40  pruner = Pruner(model, pruning_layer_list, final_pr=0.0)
41  model_for_pruning = pruner.get_pruning_model()
42  print_sparsity(model_for_pruning)
43
44  # 3. Loading the pruned model
45  model_for_pruning.load_weights("pruned_model.h5")
46  print_sparsity(model_for_pruning)
```

**Figure 3-21      [TensorFlow] Loading the Pruned Model**

## 3.8 Sample Code

This subsection describes how to execute the sample code and gives an outline of its operation. The accuracy of a model after pruning can be confirmed with the use of the sample code.

### 3.8.1 classification/pytorch_mobilenetv2

This sample code employs the MobileNetV2 architecture of PyTorch and is for use in initial training and pruning then retraining. The code is for use with the CIFAR-10 dataset. The following three files are provided. The method for adding the DRP-AI Extension Pack module can be confirmed by comparing train.py, code for use in initial training with the files with names of the form retrain*.py, containing the two variants of the code for use in pruning then retraining.

**Table 3.1     List of Provided Files**

| File Name | Description |
|---|---|
| train.py | MobileNetV2 sample code for initial training |
| retrain_with_oneshot_pruning.py | MobileNetV2 sample code for pruning then retraining (one-shot pruning) |
| retrain_with_gradual_pruning.py | MobileNetV2 sample code for pruning then retraining (gradual pruning) |

Pruning then retraining with the MobileNetV2 architecture can be performed by executing the following two steps.

```
# Step 1: Initial training

$ python3 train.py

# Step 2: Pruning then retraining

# In one-shot pruning

$ python3 retrain_with_oneshot_pruning.py

# In gradual pruning

$ python3 retrain_with_gradual_pruning.py
```

**Figure 3-22     Executing the MobileNetV2 Sample Code for Initial Training and Pruning Then Retraining**

After executing the sample code, the files listed in the table below will have been output.

**Table 3.2     List of Output Files**

| File Name | Description |
|---|---|
| pretrained_mobilenetv2.pth | Trained model file (pth format) |
| pretrained_mobilenetv2.onnx | Trained model file (ONNX format) |
| oneshot_pruned_mobilenetv2.pth | Model file after pruning in one-shot pruning mode then retraining (pth format) |
| oneshot_pruned_mobilenetv2.onnx | Model file after pruning in one-shot pruning mode then retraining (ONNX format) |
| gradual_pruned_mobilenetv2.pth | Model file after pruning in gradual pruning mode then retraining (pth format) |
| gradual_pruned_mobilenetv2.onnx | Model file after pruning in gradual pruning mode then retraining (ONNX format) |

Command-line options are listed in the table below.

**Table 3.3     List of Options of the MobileNetV2 Sample Code for Initial Training and Pruning Then Retraining**

| Option Argument | Description |
|---|---|
| -h, --help | Outputs a help message.<br>Example:<br>$ python3 train.py -h |
| --lr LR | Sets the learning rate.<br>Set a small learning rate for a case where the loss varies greatly.<br>Set a large learning rate for a case where the loss does not vary.<br>Example:<br>$ python3 train.py --lr 0.2 |
| --max_epochs MAX_EPOCHS | Specifies the maximum number of epochs.<br>If you want a greater accuracy, set a value greater than the default so that learning proceeds for a longer time.<br>Example:<br>$ python3 train.py --max_epochs 3 |
| --pretrained_weight | Specifies the name of a file (.pth format) for a model for initial training.<br>Note:  Can only be set for code for pruning then retraining.<br>Example:<br>$ python3 retrain_with_oneshot_pruning.py \<br>    –pretrained_weight ./pretrained_mobilenetv2.pth |
| --pruning_rate | Specifies the pruning rate.<br>Note:  Can only be set for code for pruning then retraining.<br>Example:<br>$ python3 retrain_with_oneshot_pruning.py \<br>    –pruning_rate 0.7 |

### 3.8.2 classification/tensorflow_cnn

This sample code employs the CNN model of TensorFlow and is for use in initial training and pruning then retraining. The code is for use with the CIFAR-10 dataset. The following three files are provided. The method for adding the DRP-AI Extension Pack module can be confirmed by comparing train.py, code for use in initial training with the files with names of the form retrain*.py, containing the two variants of the code for use in pruning then retraining.

**Table 3.4    List of Provided Files**

| File Name | Description |
|---|---|
| train.py | CNN sample code for initial training |
| retrain_with_oneshot_pruning.py | CNN sample code for pruning then retraining (one-shot pruning) |
| retrain_with_gradual_pruning.py | CNN sample code for pruning then retraining (gradual pruning) |

Pruning then retraining with the CNN model can be performed by executing the following two steps.

```
# Step 1: Initial training
$ python3 train.py
# Step 2: Pruning then retraining
# In one-shot pruning
$ python3 retrain_with_oneshot_pruning.py
# In gradual pruning
$ python3 retrain_with_gradual_pruning.py
```

**Figure 3-23    Executing the CNN Sample Code for Initial Training and Pruning Then Retraining**

After executing the sample code, the files listed in the table below will have been output.

**Table 3.5    List of Output Files**

| File Name | Description |
|---|---|
| pretrained_cnn.h5 | Trained model file (h5 format) |
| pretrained_cnn.onnx | Trained model file (ONNX format) |
| oneshot_pruned_cnn.h5 | Model file after pruning in one-shot pruning mode then retraining (h5 format) |
| oneshot_pruned_cnn.onnx | Model file after pruning in one-shot pruning mode then retraining (ONNX format) |
| gradual_pruned_cnn.h5 | Model file after pruning in gradual pruning mode then retraining (h5 format) |
| gradual_pruned_cnn.onnx | Model file after pruning in gradual pruning mode then retraining (ONNX format) |

Command-line options are listed in the table below.

**Table 3.6     List of Options of the CNN Sample Code for Initial Training and Pruning Then Retraining**

| Option Argument | Description |
|---|---|
| -h, --help | Outputs a help message.<br>Example:<br>$ python3 train.py -h |
| --lr LR | Sets the learning rate.<br>Set a small learning rate for a case where the loss varies greatly.<br>Set a large learning rate for a case where the loss does not vary.<br>Example:<br>$ python3 train.py --lr 0.2 |
| --max_epochs MAX_EPOCHS | Specifies the maximum number of epochs.<br>If you want a greater accuracy, set a value greater than the default so that learning proceeds for a longer time.<br>Example:<br>$ python3 train.py --max_epochs 3 |
| --pretrained_weight | Specifies the name of a file (.h5 format) of a model for initial training.<br>Note:  Can only be set for code for pruning then retraining.<br>Example:<br>$ python3 retrain_with_oneshot_pruning.py \<br>    –pretrained_weight ./pretrained_cnn.h5 |
| --pruning_rate | Specifies the pruning rate.<br>Note:  Can only be set for code for pruning then retraining.<br>Example:<br>$ python3 retrain_with_oneshot_pruning.py \<br>    –pruning_rate 0.7 |

## 4. Details on the DRP-AI Extension Pack API

This section describes the API functions and class provided by the DRP-AI Extension Pack.

### 4.1 List of DRP-AI Extension Pack API Functions and Class

The API functions and class provided by the DRP-AI Extension Pack are listed in the table below.

**Table 4.1     List of DRP-AI Extension Pack API Functions and Class**

| Module | Function/Class Name | Description | Section |
|---|---|---|---|
| drpai_compaction_tool. pytorch | make_pruning_layer_list | Sets layers to which pruning is not to be applied and creates the list of the target layers for pruning. | 4.2.1 |
| | Pruner | Applies pruning to the model. | 4.2.2 |
| | get_endstep | Gets the step at which pruning ends. | 4.2.3 |
| | get_frequency | Gets the frequency for updating of the pruning rate. | 4.2.4 |
| | get_model_info | Gets a list of information on the model, such as the numbers of parameters. | 4.2.5 |
| | deepcopy_model | Deep copy the model. (Deep copy means copies that are completely reproduced.) | 4.2.6 |
| | load_pruned_state_dict | Loads the model which weights is pruned | 4.2.7 |
| drpai_compaction_tool.t ensorflow | make_pruning_layer_list | Sets layers to which pruning is not to be applied and creates the list of the target layers for pruning. | 4.3.1 |
| | Pruner | Applies pruning to the model. | 4.3.2 |
| | get_endstep | Gets the step at which pruning ends. | 4.3.3 |
| | get_frequency | Gets the frequency for updating of the pruning rate. | 4.3.4 |

The API functions and class are described in terms of the following items on the following pages.

| | |
|---|---|
| [Overview] | Describes the API class or function in outline. |
| | |
| [Function/class name] | Function name or class name |
| | |
| [Calling format] | Describes the format for calling the function or the class as a function. |
| | |
| [Argument] | Describes the arguments. |
| | |
| [Returns] | Describes the return value. |
| | |
| [Feature] | Describes the function of the API. |
| | |
| [Remarks] | Describes points to note. |

## 4.2 [PyTorch]

### 4.2.1 make_pruning_layer_list

| [Overview] | Sets layers to which pruning is not to be applied and creates the list of the target layers for pruning. | |
|---|---|---|
| | | |
| [Function/class name] | make_pruning_layer_list | |
| | | |
| [Calling format] | make_pruning_layer_list(model:   torch.nn.Module,<br>                    input_size: List[Tuple[int]],<br>                    input_data: Union[List[Any], Mapping[Any]],<br>                    prune_last: bool = False,<br>                    prune_dwise: bool = False) -> List[str] | |
| | | |
| [Argument] | model: torch.nn.Module | PyTorch model |
| | input_size: List[Tuple[int]] | Shape of input data<br>Default: None<br>Set either input_size or input_data.<br>Input values in the order of [batch size, number of channels, height, width].<br>A value of at least 2 should be set as the batch size when using batch normalization. |
| | input_data:<br>Union[List[Any],<br>    Mapping[Any]] | Input data<br>Default: None<br>Set either input_size or input_data. |
| | prune_last: bool | When this argument is set to "true", the last layer is included among the targets for pruning.<br>Default: False |
| | prune_dwise: bool | When this argument is set to "true", the depthwise convolution layer is included among the targets for pruning.<br>Default: False |
| | | |
| [Returns] | pruning_layer_list: List[str] | List of the target layers for pruning |
| | | |
| [Feature] | Creates the list of the target layers for pruning based on the input model.<br>Pruning is to be applied to the layers defined with torch.nn.Conv2d or torch.nn.Linear. | |
| | | |

RENESAS

| [Remarks] | Pruning cannot be applied to the first layer because doing so significantly worsens the accuracy of the model. |
|---|---|
| | Pruning also cannot be applied to a layer for which the number of input channels is not a multiple of 32. |
| | Pruning is not applied to the last layer or depthwise convolution layer by default because doing so significantly worsens the accuracy of the model. |
| | For details on torch.nn.Module, torch.nn.Conv2d, and torch.nn.Linear, see the official documents of PyTorch. |
| | Usage example 1: |
| | >>>import torchvision |
| | >>> model = torchvision.models.resnet18(num_classes=1000) |
| | # Set a value of at least 2 as the batch size because batch normalization is used in ResNet18. |
| | >>> make_pruning_layer_list(model, \ |
| | input_size=[(2,3,224,224)]) |
| | Usage example 2: |
| | >>>import torchvision |
| | >>> model = torchvision.models.detection.ssd300_vgg16() |
| | # Inputting in the list format |
| | >>> make_pruning_layer_list(model, \ |
| | input_data=[[torch.rand(3, 300, 300)], \ |
| | [{'boxes': torch.tensor([[0, 0, 100, 100]]),\ |
| | 'labels': torch.tensor([0])}]]) |
| | # Inputting in the dict format |
| | >>> make_pruning_layer_list(model, |
| | input_data={ |
| | "images": [torch.rand(3, 300, 300)], \ |
| | "targets": [{'boxes': torch.tensor([[0, 0, 100, 100]]),\ |
| | 'labels': torch.tensor([0])}] \ |
| | } \ |
| | ) |

### 4.2.2   Pruner

| [Overview] | Controls the pruning parameters. | |
|---|---|---|
| | | |
| [Function/class name] | Pruner | |
| | | |
| [Calling format] | class Pruner(model: torch.nn.Module, | |
| | pruning_layer_list: List[str], | |
| | initial_pr: float, | |
| | final_pr: float, | |
| | begin_step: int, | |
| | end_step: int, | |
| | frequency: int, | |
| | ) | |
| | | |
| [Argument] | model: torch.nn.Module | PyTorch model |
| | pruning_layer_list: List[str] | List of the target layers for pruning |
| | initial_pr: float | Initial value of pruning rate<br>Default: 0.01<br>Input range: $0 \leq initial\_pr < 1$ |
| | final_pr: float | Final value of pruning rate<br>Default: 0.7<br>Input range: $0 \leq final\_pr < 1$ |
| | begin_step: int | Number of the step (iteration) where pruning starts<br>Default: 0<br>Input range: $0 \leq begin\_step$ |
| | end_step: int | Number of the step (iteration) where pruning ends<br>Default: -1<br>Input range: $-1 \leq end\_step$ |
| | frequency: int | Frequency for executing pruning (number of iterations)<br>Default: 100<br>Input range: $0 < frequency$ |
| | | |
| [Returns] | pruner: object | Object for setting up pruning |
| | | |
| [Feature] | The model is only pruned once when end_step is -1. (Any settings of initial_pr, begin_step, and frequency will be ignored in this case.)<br>When a value other than -1 is set, gradual pruning is applied to the model.<br>For details, refer to [Remarks] below. | |
| | | |

| [Remarks] | Call this API function after registering the parameters of the model to the optimizer. |
| --- | --- |
| | The setting of begin_step = end_step is prohibited.<br>When end_step = -1, the settings of initial_pr, begin_step, and frequency are ignored.<br>Pruning is carried out over the number of iterations set by [begin_step, end_step]. To complete pruning, training needs to have been performed for the number of iterations represented by (end_step − begin_step + 1).<br>initial_pr and final_pr must be values in the range [0.0, 1.0).<br><br>The use of the default values for initial_pr and begin_step is recommended.<br>For end_step, setting a value around 70% of the total number of iterations in training is recommended. (For example, when the total number of iterations was 100, set 70 iterations.) Note that get_endstep() can be used to set the value.<br>For frequency, setting the total number of iterations per epoch is recommended. Note that get_frequency() can be used to set the value.<br>For details on torch.nn.Module, see the official documents of PyTorch.<br><br>The setting of end_step determines the pruning mode.<br>The initial use of one-shot pruning is recommended. When this leads to an excessively great deterioration in the accuracy of the model, gradual pruning should be used.<br><br>**When end_step = -1 (one-shot pruning)**<br><br>Pruning is executed when this API function is called.<br><br>The settings of initial_pr, begin_step, and frequency are ignored.<br><br> |

**When end_step is other than -1 (gradual pruning)**

The pruning rate changes with the frequency.

Pruning is executed when the update() function (see 4.2.2.1) is called.



**Variables**

| Variable Name | Description |
|---|---|
| is_finished | The value of this variable becoming "true" indicates the completion of pruning. |

**Methods**

| Method Name | Description |
|---|---|
| update() | Updates the pruning parameters. |
| state_dict() | Returns the settings of Pruner in the dict format. |
| load_state_dict() | Loads the settings of Pruner. |

### 4.2.2.1　update

| [Overview] | Updates the pruning parameters. | |
|---|---|---|
| [Function/class name] | update | |
| [Calling format] | update() -> None | |
| [Argument] | — | — |
| [Returns] | — | — |
| [Feature] | Updates the pruning parameters. | |
| [Remarks] | This API function must always be called per iteration.<br>The timing for calling this API function is at the start (beginning) of each iteration.<br>To execute pruning during iterations of [begin_step, end_step], this API function must be called at least the number of times represented by (end_step − begin_step + 1).<br>**In one-shot pruning**<br>Updates the pruning parameters. The pruning rate of the model does not change.<br>**In gradual pruning**<br>Updates the pruning parameters and the pruning rate of the model. | |

### 4.2.2.2　state_dict

| [Overview] | Returns the settings of Pruner in the dict format. | |
|---|---|---|
| [Function/class name] | state_dict | |
| [Calling format] | state_dict() -> Dict[str, Any] | |
| [Argument] | — | — |
| [Returns] | — | Data in dict-format that includes the settings of Pruner |
| [Feature] | Returns the settings of Pruner in the dict format. | |
| [Remarks] | This API function is used when saving the settings of Pruner.<br>Usage example:<br>pruner = Pruner(model, pruning_layer_list)<br><br>torch.save(pruner.state_dict(),"pruner.pth") | |

### 4.2.2.3 load_state_dict

| | | |
|---|---|---|
| [Overview] | Loads the settings of Pruner. | |
| | | |
| [Function/class name] | load_state_dict | |
| | | |
| [Calling format] | load_state_dict(state_dict: Dict[str, Any]) -> None | |
| | | |
| [Argument] | state_dict: Dict[str, Any] | Data in dict-format that includes the settings of Pruner |
| | | |
| [Returns] | — | — |
| | | |
| [Feature] | Loads the settings of Pruner. | |
| | | |
| [Remarks] | This API function is used when loading the settings of Pruner.<br>Usage example:<br>pruner = Pruner(model, pruning_layer_list)<br><br>pruner.load_state_dict(torch.load("pruner.pth")) | |

RENESAS

### 4.2.3 get_endstep

| [Overview] | Gets the step at which pruning ends. | |
|---|---|---|
| | | |
| [Function/class name] | get_endstep | |
| | | |
| [Calling format] | get_endstep(dataloader: torch.utils.data.DataLoader,<br>max_iter: int,<br>max_epoch: int,<br>ratio: float) -> int | |
| | | |
| [Argument] | dataloader:<br>torch.utils.data.DataLoader | PyTorch data loader<br>Set the data loader for use in training. |
| | max_iter: int | Maximum number of iterations in training<br>Default: None<br>Input range: $0 < \max\_iter$<br><br>Set either max_iter or max_epoch.<br>Both cannot be set at the same time. |
| | max_epoch: int | Maximum number of epochs in training<br>Default: None<br>Input range: $0 < \max\_epoch$<br><br>Set either max_iter or max_epoch.<br>Both cannot be set at the same time. |
| | ratio: float | Ratio of the step where pruning ends to the maximum number of iterations<br>Default: 0.7<br>Input range: $0 < ratio$ |
| | | |
| [Returns] | end_step: int | Number of the step (iteration) where pruning ends |
| | | |
| [Feature] | Gets the step at which pruning ends. | |
| | | |
| [Remarks] | Set either max_iter or max_epoch.<br>Both cannot be set at the same time.<br>The use of the default value for ratio is recommended.<br>The step where pruning ends can be obtained from the following equation.<br><br>$$end\_step = Maximum\ iteration \times ratio$$<br><br>70% of the maximum number of iterations is returned by default. When training was performed for 100 iterations, this API function by default returns 70 iterations for pruning. | |

### 4.2.4　get_frequency

| [Overview] | Gets the frequency for updating of the pruning rate. | |
|---|---|---|
| [Function/class name] | get_frequency | |
| [Calling format] | get_frequency(dataloader: torch.utils.data.DataLoader,<br>　　　　　　ratio: float) -> int | |
| [Argument] | dataloader:<br>torch.utils.data.DataLoader | PyTorch data loader<br>Set the data loader for training. |
| | ratio: float | Ratio for controlling the frequency for updating of the pruning rate<br>Default: 1.0<br>Input range: $0 < ratio$ |
| [Returns] | frequency: int | Frequency (iteration) for updating of the pruning rate |
| [Feature] | Gets the frequency for updating of the pruning rate. | |
| [Remarks] | The use of the default value for ratio is recommended.<br>The total number of iterations per epoch is returned by default. When 1 epoch consists of 100 iterations, this API function returns 100 iterations.<br>When ratio is 1.0, the pruning rate is updated once every epoch. When ratio is 0.5, the pruning rate is updated twice every epoch. | |

### 4.2.5 get_model_info

| [Overview] | Gets a list of information on the model, such as the numbers of parameters. |
|---|---|
| | |
| [Function/class name] | get_model_info |
| | |

| [Calling format] | get_model_info (model: torch.nn.Module,<br>                input_size: List[Tuple[int]],<br>                input_data: Union[List[Any], Mapping[Any]],<br>        ) -> pandas.core.frame.DataFrame | | |
|---|---|---|
| | | |

| [Argument] | model: torch.nn.Module | PyTorch model |
|---|---|---|
| | input_size: List[Tuple[int]] | Shape of input data<br>Default: None<br>Set either input_size or input_data.<br>Input values in the order of [batch size, number of channels, height, width]. |
| | input_data:<br>Union[List[Any],<br>     Mapping[Any]] | Input data<br>Default: None<br>Set either input_size or input_data. |
| | | |

| [Returns] | model_info:<br>pandas.core.frame.DataFrame | List including the input shape, output shape, number of parameters, sparsity, number of multiply-and-accumulate calculations before pruning, and number of multiply-and-accumulate calculations after pruning |
|---|---|---|
| | | |

| [Feature] | Gets information on the convolution layers and fully connected layers of the model in terms of the input shape, output shape, number of parameters, sparsity, number of multiply-and-accumulate calculations before pruning, and number of multiply-and-accumulate calculations after pruning.<br>The values in terms of the items listed below are obtained for each layer. |
|---|---|

| module name | Layer name |
|---|---|
| input shape | Input shape |
| output shape | Output shape |
| params | Number of parameters |
| sparsity | Sparsity rate (pruning rate) |
| Baseline MAC | Number of multiply-and-accumulate calculations before pruning |
| Current MAC | Number of multiply-and-accumulate calculations after pruning |

| [Remarks] | This function is only applicable to the convolution layers and fully connected layers. |
|---|---|
| | The numbers of multiply-and-accumulate calculations are calculated on the assumption that each set of calculations in a multiply-and-accumulate operation is handled as a single bundle. |
| | For details on torch.nn.Module, see the official documents of PyTorch. |
| | For details on the pandas.core.frame.DataFrame class, see the official documents of pandas. |
| | Usage example: |
| | >>> import torch.nn as nn |
| | >>> import torch.nn.functional as F |
| | >>> class NeuralNetwork(nn.Module): |
| | >>>    def __init__(self): |
| | >>>      super(NeuralNetwork, self).__init__() |
| | >>>      self.conv1 = nn.Conv2d(3, 32, kernel_size=3, padding='same') |
| | >>>      self.conv2 = nn.Conv2d(32, 64, kernel_size=5, padding='same') |
| | >>>      self.fc1 = nn.Linear(64*32*32, 50) |
| | >>>      self.fc2 = nn.Linear(50, 10) |
| | >>>    def forward(self, x): |
| | >>>      x = F.relu(F.max_pool2d(self.conv1(x), 2)) |
| | >>>      x = F.relu(F.max_pool2d(self.conv2(x), 2)) |
| | >>>      x = x.view(-1, 64*32*32) |
| | >>>      x = F.relu(self.fc1(x)) |
| | >>>      x = self.fc2(x) |
| | >>>      return F.log_softmax(x, dim=1) |
| | >>> model = NeuralNetwork() |
| | >>> from drpai_compaction_tool.pytorch import get_model_info |
| | >>> print(get_model_info(model, [(1, 3, 128, 128)])) |

### 4.2.6  deepcopy_model

| | |
|---|---|
| [Overview] | Deep copy the model. (Deep copy means copies that are completely reproduced.) |
| [Function/Class Name] | deepcopy_model |
| [Calling format] | deepcopy_model(model: torch.nn.Module) -> torch.nn.Module |

| [Argument] | model: torch.nn.Module | PyTorch model |
|---|---|---|

| [Returns] | copied_model: torch.nn.Module | Deep copied PyTorch model |
|---|---|---|

| | |
|---|---|
| [Feature] | Deep copy the model. Deep copy means copies that are completely reproduced. |
| [Remarks] | deepcopy() cannot be executed for a pruned model.<br>In case of need to deep copy a pruned model, please use the deepcopy_model() function instead of using deepcopy().<br><br>For details on torch.nn.Module, see the official documents of PyTorch.<br>Usage example:<br><br>>>> import torch<br>>>> from collections import OrderedDict<br>>>> model = torch.nn.Sequential(OrderedDict([<br>    ("fc1", torch.nn.Linear(3,1024)),<br>    ("fc2", torch.nn.Linear(1024,10))<br>  ]))<br>>>> from drpai_compaction_tool.pytorch import Pruner, deepcopy_model<br>>>> _ = Pruner(model, ["fc1" "fc2"])<br>>>> copied_model = deepcopy_model(model) |

### 4.2.7  load_pruned_state_dict

| [Overview] | Loads the model which weights (stated_dict) is pruned | |
|---|---|---|
| | | |
| [Function/Class Name] | load_pruned_state_dict | |
| | | |
| [Calling format] | load_pruned_state_dict(model: torch.nn.Module,<br><br>pruned_state_dict: Dict,<br><br>strict: bool) -> None | |
| | | |
| [Argument] | model: torch.nn.Module | PyTorch model |
| | pruned_state_dict: Dict | Data in dict-format that includes the pruned weights. |
| | strict: bool | Whether to strictly enforce that the keys in pruned_state_dict match the keys returned by the model's state_dict() function.<br><br>Default: True<br><br>When this argument is set to "true", an error is returned if the keys do not match. |
| | | |
| [Returns] | - | - |
| | | |
| [Feature] | Loads the model which weights (stated_dict) is pruned.<br><br>The pruned weights (state_dict) means weight parameters stored in weight_orig / weight_mask format. | |
| | | |
| [Remarks] | This function performs the following behavior depending on each parameter. | |

This function performs the following behavior depending on each parameter.

| model | pruned_state_dict | behavior |
|---|---|---|
| sparse (pruned) | sparse (pruned) | This function loads the pruned_state_dict to model. |
| sparse (pruned) | dense | This function raises an error. |
| dense | sparse (pruned) | This function loads the pruned_state_dict to model.<br><br>Note: After the dense model to be input to this function, the weight parameters will change to weight_org/weight_mask format. |
| dense | dense | This function raises an error. |

For details on torch.nn.Module, see the official documents of PyTorch.

Usage example:

>>> import torch, torchvision

>>> pruned_model = torchvision.models.resnet18(num_classes=1000)

```
>>> from drpai_compaction_tool.pytorch import Pruner, load_pruned_state_dict

# This example prunes the layer called "layer1.0.conv1".

>>> _ = Pruner(pruned_model, ["layer1.0.conv1"])

>>> torch.save(pruned_model.state_dict(), "pruned_state_dict.pth")

>>> dense_model = torchvision.models.resnet18(num_classes=1000)

# Note: After the dense model to be input to this function, the weight parameters
will change to weight_org/weight_mask format.

>>> load_pruned_state_dict(dense_model, torch.load("pruned_state_dict.pth"))
```

## 4.3  [TensorFlow]

### 4.3.1  make_pruning_layer_list

| [Overview] | Sets layers to which pruning is not to be applied and creates the list of the target layers for pruning. | |
|---|---|---|
| | | |
| [Function/class name] | make_pruning_layer_list | |
| | | |
| [Calling format] | make_pruning_layer_list(model: tensorflow.python. <br><br> keras.engine.functional.Functional, <br><br> prune_last: bool = False, <br><br> prune_dwise: bool = False) -> List[str] | |
| | | |
| [Argument] | model: <br><br> tensorflow.python. <br><br>   keras.engine. <br><br>     functional.Functional | TensorFlow model <br> Only a functional model or sequential model can be input. |
| | prune_last: bool | When this argument is set to "true", the last layer is included among the targets for pruning. <br> Default: False |
| | prune_dwise: bool | When this argument is set to "true", the depthwise convolution layer is included among the targets for pruning. <br> Default: False |
| | | |
| [Returns] | pruning_layer_list: List[str] | List of the target layers for pruning |
| | | |
| [Feature] | Creates the list of the target layers for pruning based on the input model. <br> Pruning is to be applied to the layers defined with tensorflow.keras.layers.Conv2D or tensorflow.keras.layers.Dense. | |
| | | |
| [Remarks] | Pruning cannot be applied to the first layer because doing so significantly worsens the accuracy of the model. <br> Pruning also cannot be applied to a layer for which the number of input channels is not a multiple of 32. <br> Pruning is not applied to the last layer or depthwise convolution layer by default because doing so significantly worsens the accuracy of the model. <br> For details on tensorflow.keras.Model, tensorflow.keras.layers.Conv2D, and tensorflow.keras.layers.Dense, see the official documents of TensorFlow. | |

### 4.3.2 Pruner

| [Overview] | Controls the pruning parameters. | |
|---|---|---|
| | | |
| [Function/class name] | Pruner | |
| | | |
| [Calling format] | class Pruner (model: tensorflow.keras.Model,<br><br>                                pruning_layer_list: List[str],<br><br>                                initial_pr: float,<br><br>                                final_pr: float,<br><br>                                begin_step: int,<br><br>                                end_step: int,<br><br>                                frequency: int,<br><br>                                ) | |
| | | |
| [Argument] | model:<br>tensorflow.keras.Model | TensorFlow model |
| | pruning_layer_list: List[str] | List of the target layers for pruning |
| | initial_pr: float | Initial value of pruning rate<br>Default: 0.01<br>Input range: $0 \le initial\_pr < 1$ |
| | final_pr: float | Final value of pruning rate<br>Default: 0.7<br>Input range: $0 \le final\_pr < 1$ |
| | begin_step: int | Number of the step (iteration) where pruning starts<br>Default: 0<br>Input range: $0 \le begin\_step$ |
| | end_step: int | Number of the step (iteration) where pruning ends<br>Default: -1<br>Input range: $-1 \le end\_step$ |
| | frequency: int | Frequency for executing pruning (number of iterations)<br>Default: 100<br>Input range: $0 < frequency$ |
| | | |
| [Returns] | pruner: object | Object for setting up pruning |
| | | |
| [Feature] | The model is only pruned once when end_step is -1. (Any settings of initial_pr, begin_step, and frequency will be ignored in this case.)<br>When a value other than -1 is set, gradual pruning is applied to the model.<br>For details, refer to [Remarks] below. | |
| | | |

RENESAS

| [Remarks] | When end_step = -1, the settings of initial_pr, begin_step, and frequency are ignored.<br>The setting of begin_step = end_step is prohibited.<br>Pruning is carried out over the number of iterations set by [begin_step, end_step]. To complete pruning, training needs to have been performed for the number of iterations represented by (end_step − begin_step + 1).<br>initial_pr and final_pr must be values in the range [0.0, 1.0).<br><br>The use of the default values for initial_pr and begin_step is recommended.<br>For end_step, setting a value around 70% of the total number of iterations in training is recommended. (For example, when the total number of iterations was 100, set 70 iterations.) Note that get_endstep() can be used to set the value.<br>For frequency, setting the total number of iterations per epoch is recommended. Note that get_frequency() can be used to set the value.<br>For details on tensorflow.keras.Model, see the official documents of TensorFlow.<br><br>The setting of end_step determines the pruning mode.<br>The initial use of one-shot pruning is recommended. When this leads to an excessively great deterioration in the accuracy of the model, gradual pruning should be used. |
|---|---|

**When end_step = -1 (one-shot pruning)**

Pruning is executed only once at the beginning.

The settings of initial_pr, begin_step, and frequency are ignored.

RENESAS

**When end_step is other than -1 (gradual pruning)**

The pruning rate changes with the frequency.



**Variables**

| Variable Name | Description |
|---|---|
| — | — |

**Methods**

| Method Name | Description |
|---|---|
| get_pruning_model() | Gets the model for pruning. |

### 4.3.2.1 get_pruning_model

| | | |
|---|---|---|
| [Overview] | Gets the model for pruning. | |
| [Function/class name] | get_pruning_model | |
| [Calling format] | get_pruning_model() -> tensorflow.keras.Model | |
| [Argument] | — | — |
| [Returns] | pruning_model: tensorflow.keras.Model | Model for pruning |
| [Feature] | Gets the model for pruning. | |
| [Remarks] | — | |

### 4.3.3 get_endstep

| | |
|---|---|
| [Overview] | Gets the step at which pruning ends. |
| | |
| [Function/class name] | get_endstep |
| | |

| | |
|---|---|
| [Calling format] | get_endstep(num_data: int,<br><br>         batch_size: int,<br><br>         max_iter: int,<br><br>         max_epoch: int,<br><br>         ratio: float) -> int |
| | |

| | | |
|---|---|---|
| [Argument] | num_data: int | Number of training data |
| | batch_size: int | Batch size in training |
| | max_iter: int | Maximum number of iterations in training<br>Default: None<br>Input range: $0 < max\_iter$<br>Set either max_iter or max_epoch.<br>Both cannot be set at the same time. |
| | max_epoch: int | Maximum number of epochs in training<br>Default: None<br>Input range: $0 < max\_epoch$<br>Set either max_iter or max_epoch.<br>Both cannot be set at the same time. |
| | ratio: float | Ratio of the step where pruning ends to the maximum number of iterations<br>Default: 0.7<br>Input range: $0 < frequency$ |
| | | |

| | | |
|---|---|---|
| [Returns] | end_step: int | Number of the step (iteration) where pruning ends |
| | | |

| | |
|---|---|
| [Feature] | Gets the step at which pruning ends. |
| | |

| | |
|---|---|
| [Remarks] | Set either max_iter or max_epoch.<br>Both cannot be set at the same time.<br>The use of the default value for ratio is recommended.<br>The step where pruning ends can be obtained from the following equation.<br>$$end\_step = Maximum\ iteration \times ratio$$<br>70% of the maximum number of iterations is returned by default. When training was performed for 100 iterations, this API function by default returns 70 iterations for pruning. |

### 4.3.4 get_frequency

| [Overview] | Gets the frequency for updating of the pruning rate. | |
|---|---|---|
| [Function/class name] | get_frequency | |
| [Calling format] | get_frequency(num_data: int, <br><br> batch_size: int, <br><br> ratio: float) -> int | |
| [Argument] | num_data: int | Number of training data |
| | batch_size: int | Batch size in training |
| | ratio: float | Ratio for controlling the frequency for updating of the pruning rate <br> Default: 1.0 <br> Input range: $0 < ratio$ |
| [Returns] | frequency: int | Frequency (iteration) for updating of the pruning rate |
| [Feature] | Gets the frequency for updating of the pruning rate. | |
| [Remarks] | The use of the default value for ratio is recommended. <br> The total number of iterations per epoch is returned by default. When 1 epoch consists of 100 iterations, this API function returns 100 iterations. <br> When ratio is 1.0, the pruning rate is updated once every epoch. When ratio is 0.5, the pruning rate is updated twice every epoch. | |

## 5.  Recommendations during Application of Pruning

This section gives recommendations on how to suppress deterioration of the accuracy of the model due to the application of pruning. Attempt the recommended measures if unacceptably low accuracy is encountered after pruning then retraining.

- Perform pruning then retraining with 70% pruning rate to check the accuracy and processing performance. After that, change pruning rate depending on the accuracy and processing performance, and perform pruning then retraining again to check the accuracy and processing performance.
- Pruning should initially be performed as one-shot pruning. If the resulting accuracy is too low, try gradual pruning.
- Do not apply pruning to the first and last layers.
- Do not apply pruning to the depthwise convolution layer.
- Use the same parameters in training, such as the learning rate, optimizer, epoch, as those that were set for initial training.

Note:   The recommended measures are not guaranteed to always suppress unacceptable deterioration of the accuracy.

## 6.  Usage Notes

- deepcopy() cannot be executed for a pruned model.
  To copy a pruned model, please use deepcopy_model() function. For more details about this function, see 4.2.6.
- When loading a saved pruned model, see 3.4 and 3.7.
- Do not use early stopping with gradual pruning. Doing so may cause training to be terminated when pruning has not yet been completed. For details, see the TensorFlow sample code.
- Do not use Exponential Moving Average (EMA) with gradual pruning. It may cause incorrect pruning result. After pruning then retraining, you can confirm whether the pruning rate is correct by using get_model_info().
- When using the PruningSummaries() callback function of TensorFlow, do not also use the TensorBoard() callback function of TensorFlow. Since the TensorBoard() callback function is initialized in the

PruningSummaries() callback function, use of both callback functions is judged to represent a double definition and an error will occur. For details, see the TensorFlow sample code.

- Notes on using torch.nn.MultiheadAttention() in models
  - ✓ Considerations for pruning

    The Linear layer indicated by the red frame below will not be pruned. If pruning is desired, replace torch.nn.MultiheadAttention() with torch.nn.Linear() and re-implement.
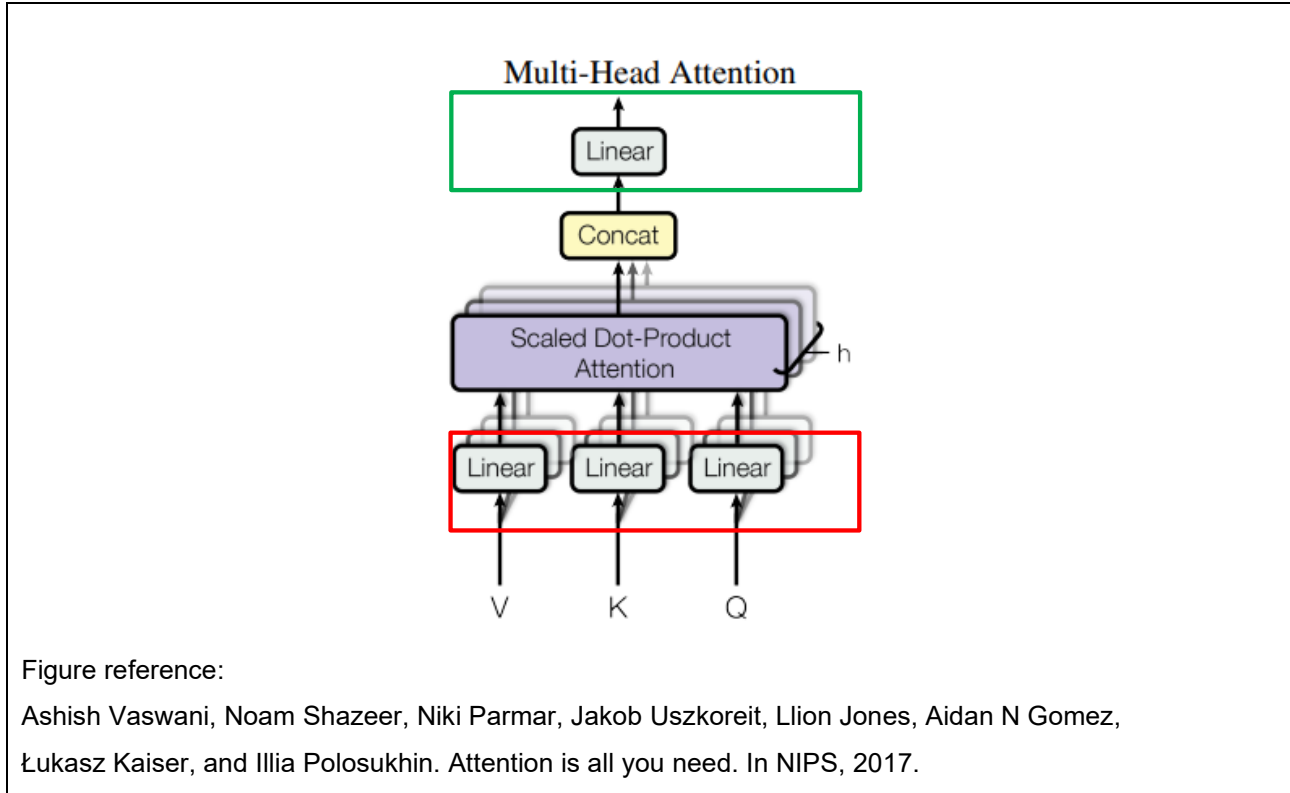


Figure reference:

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,

Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In NIPS, 2017.

**Figure 6-1 Structure of Multi-Head Attention**

- ✓ Considerations for using get_model_info()

  The input Linear layer enclosed in the red frame on Figure 6-1 does not display layer information. Additionally, the layer enclosed in the green frame (*.out_proj layer) displays "-1" for the items input_shape, output_shape, Baseline MAC, and Current MAC. The params and sparsity are displayed. This does not affect pruning behavior.



**Figure 6-2  Example output of get_model_info()**

✓ Considerations for exporting to ONNX

When exporting a model that uses torch.nn.MultiheadAttention() to ONNX, please call the following code snippet highlighted in red before exporting to ONNX.

```python
# Define a sample model using MultiheadAttention()
class SelfAttentionLike(torch.nn.Module):
    def __init__(self, embed_dim = 32, num_heads = 2):
        super(SelfAttentionLike, self).__init__()
        self._embed_dim = embed_dim
        self._num_heads = num_heads
        self.linear1 = torch.nn.Linear(self._embed_dim,
                                       self._embed_dim)
        self.multihead_attention = \
            torch.nn.MultiheadAttention(self._embed_dim,
                                        self._num_heads)
        self.linear2 = torch.nn.Linear(self._embed_dim,
                                       self._embed_dim)

    def forward(self, x):
        x = self.linear1(x)
        attn_output, _ = self.multihead_attention(query=x,
                                                  key=x,
                                                  value=x)
        output = self.linear2(attn_output)
        return output

# Create a model
embed_dim = 32
seq_len = 5
batchsize = 10
input_size = [(batchsize, seq_len, embed_dim)]
model = SelfAttentionLike(embed_dim=embed_dim)

# Perform a pruning to a model
pruning_layer_list = make_pruning_layer_list(model,
                                             input_size=input_size)
pruner = Pruner(model,
                pruning_layer_list,
                final_pr=0.7)
print(get_model_info(model,
                     input_size=input_size))

# Perform a postprocess to export a model as ONNX
import torch.nn.utils.prune as prune
for name, module in model.named_modules():
    if isinstance(module, torch.nn.MultiheadAttention):
        prune.remove(module.out_proj, 'weight')
```

```
45   # Export a model as ONNX
46   img = torch.randn(input_size[0])
47   torch.onnx.export(model,
48                     img,
49                     "exported_model.onnx",
50                     input_names=['input'],
51                     output_names=["output"])
```

- Resolving PT_GMI_002 errors when using get_model_info()

  When using get_model_info() to get the pruning rate, the following error may occur:

[PT_GMI_002] Failed to run. See above stack traces for more details. If that does not help, please contact Renesas.

  If this error occurs, please use the print_sparsity() function instead of get_model_info() as follows to verify the pruning rate:

```
1    def calc_sparsity(module):
2        if hasattr(module, 'weight_mask'):
3            m = module.weight_mask
4            sparsity = 1.0 - m.count_nonzero().item() / m.nelement()
5        else:
6            w = module.weight
7            sparsity = 1.0 - w.count_nonzero().item() / w.nelement()
8        return sparsity
9
10   def print_sparsity(model):
11       for name, module in model.named_modules():
12           if not isinstance(module, (torch.nn.Conv2d, torch.nn.Linear)):
13               continue
14           sparsity = calc_sparsity(module)
15           print(f'{name} = {sparsity:.04f}')
16
17   # Use the print_sparsity() function instead of get_model_info().
18   # Before
19   # print(get_model_info(model, [(1,3,224,224)]))
20   # After
21   print_sparsity(model)
22
```

**Revision History**

| Rev. | Date | Description | |
| --- | --- | --- | --- |
| | | Page | Summary |
| 1.00 | Dec.05.23 | — | First edition issued |
| 2.00 | Oct.08.24 | 10, 54 | Added updates in DRP-AI Extension Pack V1.1.0 |
| | | | Added usage notes on using torch.nn.MultiheadAttention() |

# Notice

1. Descriptions of circuits, software and other related information in this document are provided only to illustrate the operation of semiconductor products and application examples. You are fully responsible for the incorporation or any other use of the circuits, software, and information in the design of your product or system. Renesas Electronics disclaims any and all liability for any losses and damages incurred by you or third parties arising from the use of these circuits, software, or information.

2. Renesas Electronics hereby expressly disclaims any warranties against and liability for infringement or any other claims involving patents, copyrights, or other intellectual property rights of third parties, by or arising from the use of Renesas Electronics products or technical information described in this document, including but not limited to, the product data, drawings, charts, programs, algorithms, and application examples.

3. No license, express, implied or otherwise, is granted hereby under any patents, copyrights or other intellectual property rights of Renesas Electronics or others.

4. You shall be responsible for determining what licenses are required from any third parties, and obtaining such licenses for the lawful import, export, manufacture, sales, utilization, distribution or other disposal of any products incorporating Renesas Electronics products, if required.

5. You shall not alter, modify, copy, or reverse engineer any Renesas Electronics product, whether in whole or in part. Renesas Electronics disclaims any and all liability for any losses or damages incurred by you or third parties arising from such alteration, modification, copying or reverse engineering.

6. Renesas Electronics products are classified according to the following two quality grades: "Standard" and "High Quality". The intended applications for each Renesas Electronics product depends on the product's quality grade, as indicated below.

    "Standard": Computers; office equipment; communications equipment; test and measurement equipment; audio and visual equipment; home electronic appliances; machine tools; personal electronic equipment; industrial robots; etc.

    "High Quality": Transportation equipment (automobiles, trains, ships, etc.); traffic control (traffic lights); large-scale communication equipment; key financial terminal systems; safety control equipment; etc.

    Unless expressly designated as a high reliability product or a product for harsh environments in a Renesas Electronics data sheet or other Renesas Electronics document, Renesas Electronics products are not intended or authorized for use in products or systems that may pose a direct threat to human life or bodily injury (artificial life support devices or systems; surgical implantations; etc.), or may cause serious property damage (space system; undersea repeaters; nuclear power control systems; aircraft control systems; key plant systems; military equipment; etc.). Renesas Electronics disclaims any and all liability for any damages or losses incurred by you or any third parties arising from the use of any Renesas Electronics product that is inconsistent with any Renesas Electronics data sheet, user's manual or other Renesas Electronics document.

7. No semiconductor product is absolutely secure. Notwithstanding any security measures or features that may be implemented in Renesas Electronics hardware or software products, Renesas Electronics shall have absolutely no liability arising out of any vulnerability or security breach, including but not limited to any unauthorized access to or use of a Renesas Electronics product or a system that uses a Renesas Electronics product. RENESAS ELECTRONICS DOES NOT WARRANT OR GUARANTEE THAT RENESAS ELECTRONICS PRODUCTS, OR ANY SYSTEMS CREATED USING RENESAS ELECTRONICS PRODUCTS WILL BE INVULNERABLE OR FREE FROM CORRUPTION, ATTACK, VIRUSES, INTERFERENCE, HACKING, DATA LOSS OR THEFT, OR OTHER SECURITY INTRUSION ("Vulnerability Issues"). RENESAS ELECTRONICS DISCLAIMS ANY AND ALL RESPONSIBILITY OR LIABILITY ARISING FROM OR RELATED TO ANY VULNERABILITY ISSUES. FURTHERMORE, TO THE EXTENT PERMITTED BY APPLICABLE LAW, RENESAS ELECTRONICS DISCLAIMS ANY AND ALL WARRANTIES, EXPRESS OR IMPLIED, WITH RESPECT TO THIS DOCUMENT AND ANY RELATED OR ACCOMPANYING SOFTWARE OR HARDWARE, INCLUDING BUT NOT LIMITED TO THE IMPLIED WARRANTIES OF MERCHANTABILITY, OR FITNESS FOR A PARTICULAR PURPOSE.

8. When using Renesas Electronics products, refer to the latest product information (data sheets, user's manuals, application notes, "General Notes for Handling and Using Semiconductor Devices" in the reliability handbook, etc.), and ensure that usage conditions are within the ranges specified by Renesas Electronics with respect to maximum ratings, operating power supply voltage range, heat dissipation characteristics, installation, etc. Renesas Electronics disclaims any and all liability for any malfunctions, failure or accident arising out of the use of Renesas Electronics products outside of such specified ranges.

9. Although Renesas Electronics endeavors to improve the quality and reliability of Renesas Electronics products, semiconductor products have specific characteristics, such as the occurrence of failure at a certain rate and malfunctions under certain use conditions. Unless designated as a high reliability product or a product for harsh environments in a Renesas Electronics data sheet or other Renesas Electronics document, Renesas Electronics products are not subject to radiation resistance design. You are responsible for implementing safety measures to guard against the possibility of bodily injury, injury or damage caused by fire, and/or danger to the public in the event of a failure or malfunction of Renesas Electronics products, such as safety design for hardware and software, including but not limited to redundancy, fire control and malfunction prevention, appropriate treatment for aging degradation or any other appropriate measures. Because the evaluation of microcomputer software alone is very difficult and impractical, you are responsible for evaluating the safety of the final products or systems manufactured by you.

10. Please contact a Renesas Electronics sales office for details as to environmental matters such as the environmental compatibility of each Renesas Electronics product. You are responsible for carefully and sufficiently investigating applicable laws and regulations that regulate the inclusion or use of controlled substances, including without limitation, the EU RoHS Directive, and using Renesas Electronics products in compliance with all these applicable laws and regulations. Renesas Electronics disclaims any and all liability for damages or losses occurring as a result of your noncompliance with applicable laws and regulations.

11. Renesas Electronics products and technologies shall not be used for or incorporated into any products or systems whose manufacture, use, or sale is prohibited under any applicable domestic or foreign laws or regulations. You shall comply with any applicable export control laws and regulations promulgated and administered by the governments of any countries asserting jurisdiction over the parties or transactions.

12. It is the responsibility of the buyer or distributor of Renesas Electronics products, or any other party who distributes, disposes of, or otherwise sells or transfers the product to a third party, to notify such third party in advance of the contents and conditions set forth in this document.

13. This document shall not be reprinted, reproduced or duplicated in any form, in whole or in part, without prior written consent of Renesas Electronics.

14. Please contact a Renesas Electronics sales office if you have any questions regarding the information contained in this document or Renesas Electronics products.

(Note1) "Renesas Electronics" as used in this document means Renesas Electronics Corporation and also includes its directly or indirectly controlled subsidiaries.

(Note2) "Renesas Electronics product(s)" means any product developed or manufactured by or for Renesas Electronics.

(Rev.5.0-1   October 2020)

## Corporate Headquarters

TOYOSU FORESIA, 3-2-24 Toyosu,
Koto-ku, Tokyo 135-0061, Japan

www.renesas.com

## Trademarks

Renesas and the Renesas logo are trademarks of Renesas Electronics Corporation. All trademarks and registered trademarks are the property of their respective owners.

## Contact information

For further information on a product, technology, the most up-to-date version of a document, or your nearest sales office, please visit:
www.renesas.com/contact/.