

▼ Copyright 2020 Google LLC.

Licensed under the Apache License, Version 2.0 (the "License");



## ▼ BigTransfer (BiT): A step-by-step tutorial for state-of-the-art vision

This colab demonstrates how to:

1. Load BiT models in PyTorch
2. Make predictions using BiT pre-trained on ImageNet
3. Fine-tune BiT on 5-shot CIFAR10 and get amazing results!

It is good to get an understanding or quickly try things. However, to run longer training runs, we recommend using the commandline scripts at [http://github.com/google-research/big\\_transfer](http://github.com/google-research/big_transfer)

```
from functools import partial
from collections import OrderedDict
```

```
%config InlineBackend.figure_format = 'retina'
```

```
import numpy as np
```

```
import matplotlib.pyplot as plt
from mpl_toolkits.axes_grid1 import ImageGrid
```

```
import torch
import torch.nn as nn
import torch.nn.functional as F
from torch.autograd import Variable
```

```
import torchvision as tv
```

```
torch.cuda.empty_cache()
```

```
device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
print(device)
```

```
cuda:0
```

## ▼ Reading weight data from the Cloud bucket

```
import requests
import io

def get_weights(bit_variant):
    response = requests.get(f'https://storage.googleapis.com/bit\_models/{bit\_variant}.npy')
    response.raise_for_status()
    return np.load(io.BytesIO(response.content))

weights = get_weights('BiT-S-R50x3') # You could use other variants, such as R101x3 c
```

## ▼ Defining the architecture and loading weights

```
class StdConv2d(nn.Conv2d):
    def forward(self, x):
        w = self.weight
        v, m = torch.var_mean(w, dim=[1, 2, 3], keepdim=True, unbiased=False)
        w = (w - m) / torch.sqrt(v + 1e-10)
        return F.conv2d(x, w, self.bias, self.stride, self.padding, self.dilation, self.groups)

def conv3x3(cin, cout, stride=1, groups=1, bias=False):
    return StdConv2d(cin, cout, kernel_size=3, stride=stride, padding=1, bias=bias, groups=groups)

def conv1x1(cin, cout, stride=1, bias=False):
    return StdConv2d(cin, cout, kernel_size=1, stride=stride, padding=0, bias=bias)

def tf2th(conv_weights):
    """Possibly convert HWIO to OIHW"""
    if conv_weights.ndim == 4:
        conv_weights = np.transpose(conv_weights, [3, 2, 0, 1])
    return torch.from_numpy(conv_weights)

class PreActBottleneck(nn.Module):
    """
    Follows the implementation of "Identity Mappings in Deep Residual Networks" here:
    https://github.com/KaimingHe/resnet-1k-layers/blob/master/resnet-pre-act.lua

    Except it puts the stride on 3x3 conv when available.
    """
    def __init__(self, cin, cout=None, cmid=None, stride=1):
```

```

super().__init__()
cout = cout or cin
cmid = cmid or cout//4

self.gn1 = nn.GroupNorm(32, cin)
self.conv1 = conv1x1(cin, cmid)
self.gn2 = nn.GroupNorm(32, cmid)
self.conv2 = conv3x3(cmid, cmid, stride) # Original ResNetv2 has it on conv1!!
self.gn3 = nn.GroupNorm(32, cmid)
self.conv3 = conv1x1(cmid, cout)
self.relu = nn.ReLU(inplace=True)

if (stride != 1 or cin != cout):
    # Projection also with pre-activation according to paper.
    self.downsample = conv1x1(cin, cout, stride)

def forward(self, x):
    # Conv'ed branch
    out = self.relu(self.gn1(x))

    # Residual branch
    residual = x
    if hasattr(self, 'downsample'):
        residual = self.downsample(out)

    # The first block has already applied pre-act before splitting, see Appendix.
    out = self.conv1(out)
    out = self.conv2(self.relu(self.gn2(out)))
    out = self.conv3(self.relu(self.gn3(out)))

    return out + residual

def load_from(self, weights, prefix=''):
    with torch.no_grad():
        self.conv1.weight.copy_(tf2th(weights[prefix + 'a/standardized_conv2d/kernel']))
        self.conv2.weight.copy_(tf2th(weights[prefix + 'b/standardized_conv2d/kernel']))
        self.conv3.weight.copy_(tf2th(weights[prefix + 'c/standardized_conv2d/kernel']))
        self.gn1.weight.copy_(tf2th(weights[prefix + 'a/group_norm/gamma']))
        self.gn2.weight.copy_(tf2th(weights[prefix + 'b/group_norm/gamma']))
        self.gn3.weight.copy_(tf2th(weights[prefix + 'c/group_norm/gamma']))
        self.gn1.bias.copy_(tf2th(weights[prefix + 'a/group_norm/beta']))
        self.gn2.bias.copy_(tf2th(weights[prefix + 'b/group_norm/beta']))
        self.gn3.bias.copy_(tf2th(weights[prefix + 'c/group_norm/beta']))
        if hasattr(self, 'downsample'):
            self.downsample.weight.copy_(tf2th(weights[prefix + 'a/proj/standardized_conv2d/kernel']))
    return self

class ResNetV2(nn.Module):
    BLOCK_UNITS = {
        'r50': [3, 4, 6, 3],
        'r101': [3, 4, 23, 3],
    }

```

```

    'r152': [3, 8, 36, 3],
}

def __init__(self, block_units, width_factor, head_size=21843, zero_head=False):
    super().__init__()
    wf = width_factor # shortcut 'cause we'll use it a lot.

    self.root = nn.Sequential(OrderedDict([
        ('conv', StdConv2d(3, 64*wf, kernel_size=7, stride=2, padding=3, bias=False)),
        ('padp', nn.ConstantPad2d(1, 0)),
        ('pool', nn.MaxPool2d(kernel_size=3, stride=2, padding=0)),
        # The following is subtly not the same!
        #('pool', nn.MaxPool2d(kernel_size=3, stride=2, padding=1)),
    ]))

    self.body = nn.Sequential(OrderedDict([
        ('block1', nn.Sequential(OrderedDict(
            [('unit01', PreActBottleneck(cin= 64*wf, cout=256*wf, cmid=64*wf))] +
            [(f'unit{i:02d}', PreActBottleneck(cin=256*wf, cout=256*wf, cmid=64*wf)) 1
        ))),
        ('block2', nn.Sequential(OrderedDict(
            [('unit01', PreActBottleneck(cin=256*wf, cout=512*wf, cmid=128*wf, stride=
            [(f'unit{i:02d}', PreActBottleneck(cin=512*wf, cout=512*wf, cmid=128*wf))
        ))),
        ('block3', nn.Sequential(OrderedDict(
            [('unit01', PreActBottleneck(cin= 512*wf, cout=1024*wf, cmid=256*wf, stric
            [(f'unit{i:02d}', PreActBottleneck(cin=1024*wf, cout=1024*wf, cmid=256*wf)
        ))),
        ('block4', nn.Sequential(OrderedDict(
            [('unit01', PreActBottleneck(cin=1024*wf, cout=2048*wf, cmid=512*wf, stric
            [(f'unit{i:02d}', PreActBottleneck(cin=2048*wf, cout=2048*wf, cmid=512*wf)
        ))),
    ]))

    self.zero_head = zero_head
    self.head = nn.Sequential(OrderedDict([
        ('gn', nn.GroupNorm(32, 2048*wf)),
        ('relu', nn.ReLU(inplace=True)),
        ('avg', nn.AdaptiveAvgPool2d(output_size=1)),
        ('conv', nn.Conv2d(2048*wf, head_size, kernel_size=1, bias=True)),
    ]))

def forward(self, x):
    x = self.head(self.body(self.root(x)))
    assert x.shape[-2:] == (1, 1) # We should have no spatial shape left.
    return x[... ,0,0]

def load_from(self, weights, prefix='resnet/'):
    with torch.no_grad():
        self.root.conv.weight.copy_(tf2th(weights[f'{prefix}root_block/standardized_conv
        self.head.gn.weight.copy_(tf2th(weights[f'{prefix}group_norm/gamma' ]))
        self.head.gn.bias.copy_(tf2th(weights[f'{prefix}group_norm/beta' ]))

```

```

self.head.gn.bias.copy_(tf2th(weights[f'{prefix}group_norm/beta']))
if self.zero_head:
    nn.init.zeros_(self.head.conv.weight)
    nn.init.zeros_(self.head.conv.bias)
else:
    self.head.conv.weight.copy_(tf2th(weights[f'{prefix}head/conv2d/kernel']))
    self.head.conv.bias.copy_(tf2th(weights[f'{prefix}head/conv2d/bias']))

    for bname, block in self.body.named_children():
        for uname, unit in block.named_children():
            unit.load_from(weights, prefix=f'{prefix}{bname}/{uname}/')
return self

```

## ▼ Boilerplate

```

from IPython.display import HTML, display

def progress(value, max=100):
    return HTML("""
        <progress
            value='{value}'
            max='{max}',
            style='width: 100%'
        >
            {value}
        </progress>
    """).format(value=value, max=max)

def stairs(s, v, *svs):
    """ Implements a typical "stairs" schedule for learning-rates.
    Best explained by example:
    stairs(s, 0.1, 10, 0.01, 20, 0.001)
    will return 0.1 if s<10, 0.01 if 10<=s<20, and 0.001 if 20<=s
    """
    for s0, v0 in zip(svs[::2], svs[1::2]):
        if s < s0:
            break
        v = v0
    return v

def rampup(s, peak_s, peak_lr):
    if s < peak_s: # Warmup
        return s/peak_s * peak_lr
    else:
        return peak_lr

def schedule(s):
    step_lr = stairs(s, 3e-3, 200, 3e-4, 300, 3e-5, 400, 3e-6, 500, None)
    return rampup(s, 100, step_lr)

```

```
return rampup(s, 100, step_lr)
```

## ▼ CIFAR-10 Example

```
import PIL
```

```
preprocess_train = tv.transforms.Compose([
    tv.transforms.Resize((160, 160), interpolation=PIL.Image.BILINEAR), # It's the de
    tv.transforms.RandomCrop((128, 128)),
    tv.transforms.RandomHorizontalFlip(),
    tv.transforms.ToTensor(),
    tv.transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5)) # Get data into [-1, 1]
])
```

```
preprocess_eval = tv.transforms.Compose([
    tv.transforms.Resize((128, 128), interpolation=PIL.Image.BILINEAR),
    tv.transforms.ToTensor(),
    tv.transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5))
])
```

```
trainset = tv.datasets.CIFAR10(root='./data', train=True, download=True, transform=pre
testset = tv.datasets.CIFAR10(root='./data', train=False, download=True, transform=pre
```

```
/usr/local/lib/python3.7/dist-packages/torchvision/transforms/transforms.py:258:
"Argument interpolation should be of type InterpolationMode instead of int. "
Downloading https://www.cs.toronto.edu/~kriz/cifar-10-python.tar.gz to ./data/ci
170499072/? [00:13<00:00, 12682716.75it/s]
```

```
Extracting ./data/cifar-10-python.tar.gz to ./data
Files already downloaded and verified
```

```
import random
```

```
n = random.randint(0, 10)
```

## ▼ Find indices to create a N-shot CIFAR10 variant

```
preprocess_tiny = tv.transforms.Compose([tv.transforms.CenterCrop((2, 2)), tv.transfoi
trainset_tiny = tv.datasets.CIFAR10(root='./data', train=True, download=True, transfoi
loader = torch.utils.data.DataLoader(trainset_tiny, batch_size=50000, shuffle=False, r
images, labels = iter(loader).next()
```

```
Files already downloaded and verified
```

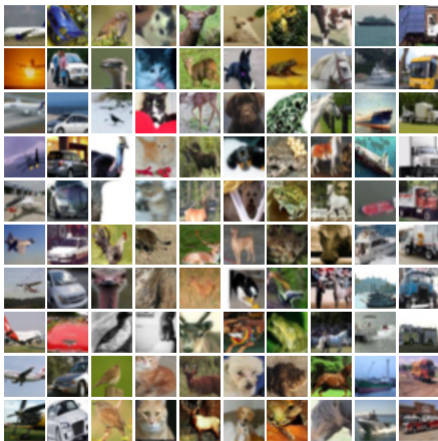
```
indices = {cls: np.random.choice(np.where(labels.numpy() == cls)[0], n, replace=False)}

print(indices)

{0: array([10148, 4651, 38344, 17965, 13225, 9616, 46343, 42010, 33933,
          16529]), 1: array([10711, 6269, 46693, 36550, 30714, 3062, 26375, 5841,
          21003]), 2: array([12675, 42098, 7664, 28238, 21102, 34654, 29730, 803,
          15765]), 3: array([18127, 41866, 33955, 27766, 28488, 2982, 45467, 39135,
          35245]), 4: array([46319, 37004, 24942, 5919, 44977, 37875, 2106, 25803,
          13909]), 5: array([ 928, 18280, 26010, 3837, 12371, 43330, 10542, 42862,
          11155]), 6: array([47464, 19586, 36304, 15048, 13511, 23683, 20148, 7994,
          4026]), 7: array([20738, 17469, 37882, 6939, 45459, 33019, 12838, 47522,
          12394]), 8: array([22760, 24629, 18045, 39139, 35171, 38809, 46735, 4502,
          21134]), 9: array([26374, 45735, 17704, 2917, 27694, 1590, 35280, 47896,
          41348])}
```

```
fig = plt.figure(figsize=(10, 4))
ig = ImageGrid(fig, 111, (n, 10))
for c, cls in enumerate(indices):
    for r, i in enumerate(indices[cls]):
        img, _ = trainset[i]
        ax = ig.axes_column[c][r]
        ax.imshow((img.numpy().transpose([1, 2, 0]) * 127.5 + 127.5).astype(np.uint8))
        ax.set_axis_off()
fig.suptitle('The whole 10-shot CIFAR10 dataset');
```

The whole 10-shot CIFAR10 dataset



```
train_n_shot = torch.utils.data.Subset(trainset, indices=[i for v in indices.values()
len(train_n_shot)
```

100

## ▼ Fine-tune BiT-M on this 10-shot CIFAR10 variant

**NOTE:** In this very low data regime, the performance heavily depends on how "representative" the 5 examples you got are of the class. As shown in the paper, variance is very large, I'm getting anywhere between 78%-85% depending on luck.

Another point is that here I use `batch_size=512` for consistency with the paper. But actually, a much smaller `batch_size=50` works just as well and is about 10x faster!

```
model = ResNetV2(ResNetV2.BLOCK_UNITS['r50'], width_factor=3, head_size=10, zero_head=
model.load_from(weights)
model.to(device);
```

```
# Yes, we still use 512 batch-size! Maybe something else is even better, who knows.
# loader_train = torch.utils.data.DataLoader(train_5shot, batch_size=512, shuffle=True)
```

```
# NOTE: This is necessary when the batch-size is larger than the dataset.
sampler = torch.utils.data.RandomSampler(train_10shot, replacement=True, num_samples=2
loader_train = torch.utils.data.DataLoader(train_10shot, batch_size=50, num_workers=2,
```

```
crit = nn.CrossEntropyLoss()
opti = torch.optim.SGD(model.parameters(), lr=0.003, momentum=0.9)
model.train();
```

```
S = 500
def schedule(s):
    step_lr = stairs(s, 3e-3, 200, 3e-4, 300, 3e-5, 400, 3e-6, S, None)
    return rampup(s, 100, step_lr)
```

```
pb_train = display(progress(0, S), display_id=True)
pb_test = display(progress(0, 100), display_id=True)
losses = [[]]
accus_train = [[]]
accus_test = []
```

```
steps_per_iter = 512 // loader_train.batch_size
```

```
while len(losses) < S:
    for x, t in loader_train:
        x, t = x.to(device), t.to(device)

        logits = model(x)
        loss = crit(logits, t) / steps_per_iter
        loss.backward()
        losses[-1].append(loss.item())
```



```

with torch.no_grad():
    accus_train[-1].extend(torch.max(logits, dim=1)[1].cpu().numpy() == t.cpu().num

if len(losses[-1]) == steps_per_iter:
    losses[-1] = sum(losses[-1])
    losses.append([])
    accus_train[-1] = np.mean(accus_train[-1])
    accus_train.append([])


# Update learning-rate according to schedule, and stop if necessary
lr = schedule(len(losses) - 1)
for param_group in opti.param_groups:
    param_group['lr'] = lr

opti.step()
opti.zero_grad()

pb_train.update(progress(len(losses) - 1, S))
print(f'\r[Step {len(losses) - 1}] loss={losses[-2]:.2e} '
      f'train accu={accus_train[-2]:.2%} '
      f'test accu={accus_test[-1] if accus_test else 0:.2%} '
      f'(lr={lr:g})', end='', flush=True)

if len(losses) % 25 == 0:
    accus_test.append(eval_cifar10(model, progressbar=pb_test))
    model.train()

```



```

[Step 305] loss=6.14e-04 train accu=100.00% test accu=82.22% (lr=3e-05)

```

```

fig, (ax1, ax2, ax3) = plt.subplots(1, 3, figsize=(15, 4))
ax1.plot(losses[:-1])
ax1.set_yscale('log')
ax1.set_title('loss')
ax2.plot(accus_train[:-1])
ax2.set_title('training accuracy')
ax3.plot(np.arange(25, 501, 25), accus_test)
ax3.set_title('test accuracy');

```

