

```

loss = F.cross_entropy(preds.view(-1,preds.size(-1)), labels.view(-1))
acc = (preds.argmax(dim=-1) == labels).float().mean()

# Logging
self.log("f'{mode}_loss", loss)
self.log("f'{mode}_acc", acc)
return loss, acc

def training_step(self, batch, batch_idx):
    loss, _ = self._calculate_loss(batch, mode="train")
    return loss

def validation_step(self, batch, batch_idx):
    _ = self._calculate_loss(batch, mode="val")

def test_step(self, batch, batch_idx):
    _ = self._calculate_loss(batch, mode="test")

```

Finally, we can create a training function similar to the one we have seen in Tutorial 5 for PyTorch Lightning. We create a `pl.Trainer` object, running for N epochs, logging in TensorBoard, and saving our best model based on the validation. Afterward, we test our models on the test set. An additional parameter we pass to the trainer here is `gradient_clip_val`. This clips the norm of the gradients for all parameters before taking an optimizer step and prevents the model from diverging if we obtain very high gradients at, for instance, sharp loss surfaces (see many good blog posts on gradient clipping, like [DeepAI glossary](#)). For Transformers, gradient clipping can help to further stabilize the training during the first few iterations, and also afterward. In plain PyTorch, you can apply gradient clipping via `torch.nn.utils.clip_grad_norm_(...)` (see [documentation](#)). The clip value is usually between 0.5 and 10, depending on how harsh you want to clip large gradients. After having explained this, let's implement the training function:

In [18]:

```

def train_reverse(**kwargs):
    # Create a PyTorch Lightning trainer with the generation callback
    root_dir = os.path.join(CHECKPOINT_PATH, "ReverseTask")
    os.makedirs(root_dir, exist_ok=True)
    trainer = pl.Trainer(default_root_dir=root_dir,
                         callbacks=[ModelCheckpoint(save_weights_only=True, mode="max", monitor="val_acc")],
                         accelerator="gpu" if str(device).startswith("cuda") else "cpu",
                         devices=1,
                         max_epochs=10,
                         gradient_clip_val=5)
    trainer.logger._default_hp_metric = None # Optional Logging argument that we don't need

# Check whether pretrained model exists. If yes, Load it and skip training

```

```

pretrained_filename = os.path.join(CHECKPOINT_PATH, "ReverseTask.ckpt")

if os.path.isfile(pretrained_filename):
    print("Found pretrained model, loading...")
    model = ReversePredictor.load_from_checkpoint(pretrained_filename)

else:
    model = ReversePredictor(max_iters=trainer.max_epochs*len(train_loader), **kwargs)

    # Test best model on validation and test set
    val_result = trainer.test(model, val_loader, verbose=False)
    test_result = trainer.test(model, test_loader, verbose=False)
    result = {"test_acc": test_result[0]["test_acc"], "val_acc": val_result[0]["test_acc"]}

    model = model.to(device)
    return model, result

```

Finally, we can train the model. In this setup, we will use a single encoder block and a single head in the Multi-Head Attention. This is chosen because of the simplicity of the task, and in this case, the attention can actually be interpreted as an "explanation" of the predictions (compared to the other papers above dealing with deep Transformers).

```

In [19]: reverse_model, reverse_result = train_reverse(input_dim=train_loader.dataset.num_categories,
                                                    model_dim=32,
                                                    num_heads=1,
                                                    num_classes=train_loader.dataset.num_categories,
                                                    num_layers=1,
                                                    dropout=0.0,
                                                    lr=5e-4,
                                                    warmup=50)

```

```

GPU available: True, used: True
TPU available: False, using: 0 TPU cores
LOCAL_RANK: 0 - CUDA_VISIBLE_DEVICES: [0]
Found pretrained model, loading...
/home/phillip/anaconda3/envs/nlp1/lib/python3.7/site-packages/pytorch_lightning/utilities/distributed.py:45: UserWarning: The dataloader, test dataloader 0, does not have many workers which may be a bottleneck. Consider increasing the value of the `num_workers` argument (try 16 which is the number of cpus on this machine) in the `DataLoader` init to improve performance.
warnings.warn(*args, **kwargs)
HBox(children=(FloatProgress(value=1.0, bar_style='info', description='Testing', layout=Layout(flex='2'), max=...,
HBox(children=(FloatProgress(value=1.0, bar_style='info', description='Testing', layout=Layout(flex='2'), max=...

```

The warning of PyTorch Lightning regarding the number of workers can be ignored for now. As the data set is so simple and the `__getitem__` finishes a neglectable time, we don't need subprocesses to provide us the data (in fact, more workers can slow down the training as we have communication overhead among processes/threads). First, let's print the results:

```
In [20]: print(f"Val accuracy: {100.0 * reverse_result['val_acc']}%")
print(f"Test accuracy: {100.0 * reverse_result['test_acc']}%")

Val accuracy: 100.00%
Test accuracy: 100.00%
```

As we would have expected, the Transformer can correctly solve the task. However, how does the attention in the Multi-Head Attention block looks like for an arbitrary input? Let's try to visualize it below.

```
In [21]: data_input, labels = next(iter(val_loader))
inp_data = F.one_hot(data_input, num_classes=reverse_model.hparams.num_classes).float()
inp_data = inp_data.to(device)
attention_maps = reverse_model.get_attention_maps(inp_data)
```

The object `attention_maps` is a list of length N where N is the number of layers. Each element is a tensor of shape [Batch, Heads, SeqLen, SeqLen], which we can verify below.

```
In [22]: attention_maps[0].shape
Out[22]: torch.Size([128, 1, 16, 16])
```

Next, we will write a plotting function that takes as input the sequences, attention maps, and an index indicating for which batch element we want to visualize the attention map. We will create a plot where over rows, we have different layers, while over columns, we show the different heads. Remember that the softmax has been applied for each row separately.

```
In [23]: def plot_attention_maps(input_data, attn_maps, idx=0):
    if input_data is not None:
        input_data = input_data[idx].detach().cpu().numpy()
    else:
        input_data = np.arange(attn_maps[0][idx].shape[-1])
    attn_maps = [m[idx].detach().cpu().numpy() for m in attn_maps]

    num_heads = attn_maps[0].shape[0]
    num_layers = len(attn_maps)
```

```

seq_len = input_data.shape[0]
fig_size = 4 if num_heads == 1 else 3
fig, ax = plt.subplots(num_layers, num_heads, figsize=(num_heads*fig_size, num_layers*fig_size))

if num_layers == 1:
    ax = [ax]

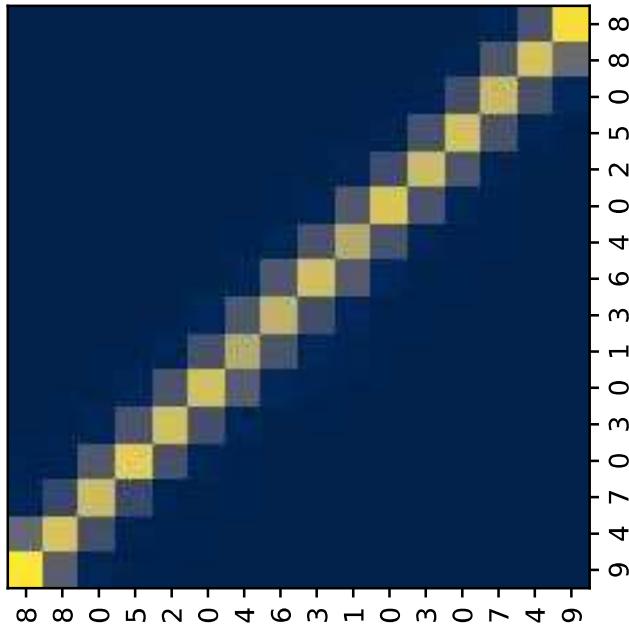
if num_heads == 1:
    ax = [[a] for a in ax]
for row in range(num_layers):
    for column in range(num_heads):
        ax[row][column].imshow(attn_maps[row][column], origin='lower', vmin=0)
        ax[row][column].set_xticks(list(range(seq_len)))
        ax[row][column].set_yticks(list(range(seq_len)))
        ax[row][column].set_yticklabels(input_data.tolist())
        ax[row][column].set_xticklabels(input_data.tolist())
        ax[row][column].set_title(f"Layer {row+1}, Head {column+1}")
fig.subplots_adjust(hspace=0.5)
plt.show()

```

Finally, we can plot the attention map of our trained Transformer on the reverse task:

```
In [24]: plot_attention_maps(data_input, attention_maps, idx=0)
```

Layer 1, Head 1



The model has learned to attend to the token that is on the flipped index of itself. Hence, it actually does what we intended it to do. We see that it however also pays some attention to values close to the flipped index. This is because the model doesn't need the perfect, hard attention to solve this problem, but is fine with this approximate, noisy attention map. The close-by indices are caused by the similarity of the positional encoding, which we also intended with the positional encoding.

Set Anomaly Detection

Besides sequences, sets are another data structure that is relevant for many applications. In contrast to sequences, elements are unordered in a set. RNNs can only be applied on sets by assuming an order in the data, which however biases the model towards a non-existing order in the data. [Vinyals et al. \(2015\)](#) and other papers have shown that the assumed order can have a significant impact on the model's performance, and hence, we should try to not use RNNs on sets. Ideally, our model should be permutation-equivariant/invariant such that the output is the same no matter how we sort the elements in a set.

Transformers offer the perfect architecture for this as the Multi-Head Attention is permutation-equivariant, and thus, outputs the same values no matter in what order we enter the inputs (inputs and outputs are permuted equally). The task we are looking at for sets is *Set Anomaly Detection* which means that we try to find the element(s) in a set that does not fit the others. In the research community, the common application of

anomaly detection is performed on a set of images, where $N - 1$ images belong to the same category/have the same high-level features while one belongs to another category. Note that category does not necessarily have to relate to a class in a standard classification problem, but could be the combination of multiple features. For instance, on a face dataset, this could be people with glasses, male, beard, etc. An example of distinguishing different animals can be seen below. The first four images show foxes, while the last represents a different animal. We want to recognize that the last image shows a different animal, but it is not relevant which class of animal it is.



In this tutorial, we will use the CIFAR100 dataset. CIFAR100 has 600 images for 100 classes each with a resolution of 32x32, similar to CIFAR10. The larger amount of classes requires the model to attend to specific features in the images instead of coarse features as in CIFAR10, therefore making the task harder. We will show the model a set of 9 images of one class, and 1 image from another class. The task is to find the image that is from a different class than the other images. Using the raw images directly as input to the Transformer is not a good idea, because it is not translation invariant as a CNN, and would need to learn to detect image features from high-dimensional input first of all. Instead, we will use a pre-trained ResNet34 model from the torchvision package to obtain high-level, low-dimensional features of the images. The ResNet model has been pre-trained on the [ImageNet](#) dataset which contains 1 million images of 1k classes and varying resolutions. However, during training and testing, the images are usually scaled to a resolution of 224x224, and hence we rescale our CIFAR images to this resolution as well. Below, we will load the dataset, and prepare the data for being processed by the ResNet model.

```
In [25]: # Import statistics
DATA_MEANS = np.array([0.485, 0.456, 0.406])
DATA_STD = np.array([0.229, 0.224, 0.225])
# As torch tensors for later preprocessing
TORCH_DATA_MEANS = torch.from_numpy(DATA_MEANS).view(1,3,1,1)
TORCH_DATA_STD = torch.from_numpy(DATA_STD).view(1,3,1,1)

# Resize to 224x224, and normalize to ImageNet statistic
transform = transforms.Compose([transforms.Resize((224,224)),
                               transforms.ToTensor(),
                               transforms.Normalize(DATA_MEANS, DATA_STD)])
])

# Loading the training dataset.
train_set = CIFAR100(root=DATASET_PATH, train=True, transform=transform, download=True)

# Loading the test set
test_set = CIFAR100(root=DATASET_PATH, train=False, transform=transform, download=True)
```

Files already downloaded and verified

Next, we want to run the pre-trained ResNet model on the images, and extract the features before the classification layer. These are the most high-level features, and should sufficiently describe the images. CIFAR100 has some similarity to ImageNet, and thus we are not retraining the ResNet model in any form. However, if you would want to get the best performance and have a very large dataset, it would be better to add the ResNet to the computation graph during training and finetune its parameters as well. As we don't have a large enough dataset and want to train our model efficiently, we will extract the features beforehand. Let's load and prepare the model below.

```
In [26]:  
import os  
os.environ["TORCH_HOME"] = CHECKPOINT_PATH  
pretrained_model = torchvision.models.resnet34(weights='IMAGENET1K_V1')  
# Remove classification layer  
# In some models, it is called "fc", others have "classifier"  
# Setting both to an empty sequential represents an identity map of the final features.  
pretrained_model.fc = nn.Sequential()  
pretrained_model.classifier = nn.Sequential()  
# To GPU  
pretrained_model = pretrained_model.to(device)  
  
# Only eval, no gradient required  
pretrained_model.eval()  
for p in pretrained_model.parameters():  
    p.requires_grad = False
```

We will now write a extraction function for the features below. This cell requires access to a GPU, as the model is rather deep and the images relatively large. The GPUs on GoogleColab are sufficient, but running this cell can take 2-3 minutes. Once it is run, the features are exported on disk so they don't have to be recalculated every time you run the notebook. However, this requires >150MB free disk space. So it is recommended to run this only on a local computer if you have enough free disk and a GPU (GoogleColab is fine for this). If you do not have a GPU, you can download the features from the [GoogleDrive folder](#).

```
In [27]:  
@torch.no_grad()  
def extract_features(dataset, save_file):  
    if not os.path.isfile(save_file):  
        data_loader = data.DataLoader(dataset, batch_size=128, shuffle=False, drop_last=False, num_workers=4)  
        extracted_features = []  
        for imgs, _ in tqdm(data_loader):  
            imgs = imgs.to(device)  
            feats = pretrained_model(imgs)
```

```

    extracted_features.append(feats)
    extracted_features = torch.cat(extracted_features, dim=0)
    extracted_features = extracted_features.detach().cpu()
    torch.save(extracted_features, save_file)

else:
    extracted_features = torch.load(save_file)

return extracted_features

```

Let's verify the feature shapes below. The training should have 50k elements, and the test 10k images. The feature dimension is 512 for the ResNet34. If you experiment with other models, you likely see a different feature dimension.

```

In [28]: print("Train:", train_set_feats.shape)
         print("Test:", test_feats.shape)

```

```

Train: torch.Size([50000, 512])
Test: torch.Size([10000, 512])

```

As usual, we want to create a validation set to detect when we should stop training. In this case, we will split the training set into 90% training, 10% validation. However, the difficulty is here that we need to ensure that the validation set has the same number of images for all 100 labels. Otherwise, we have a class imbalance which is not good for creating the image sets. Hence, we take 10% of the images for each class, and move them into the validation set. The code below does exactly this.

```

## Split train into train+val
# Get Labels from train set
labels = train_set.targets

# Get indices of images per class
labels = torch.LongTensor(labels)
num_labels = labels.max() + 1 # [classes, num_imgs_per_class]

# Determine number of validation images per class
sorted_indices = torch.argsort(labels).reshape(num_labels, -1) # [classes, num_imgs_per_class]
num_val_exmps = sorted_indices.shape[1] // 10

# Get image indices for validation and training

```

In [29]:

```

val_indices = sorted_indices[:, :num_val_exmps].reshape(-1)
train_indices = sorted_indices[:, num_val_exmps:].reshape(-1)

# Group corresponding image features and labels
train_feats, train_labels = train_set_feats[train_indices], labels[train_indices]
val_feats, val_labels = train_set_feats[val_indices], labels[val_indices]

```

Now we can prepare a dataset class for the set anomaly task. We define an epoch to be the sequence in which each image has been exactly once as an "anomaly". Hence, the length of the dataset is the number of images in it. For the training set, each time we access an item with `__getitem__`, we sample a random, different class than the image at the corresponding index `idx` has. In a second step, we sample $N - 1$ images of this sampled class. The set of 10 images is finally returned. The randomness in the `__getitem__` allows us to see a slightly different set during each iteration. However, we can't use the same strategy for the test set as we want the test dataset to be the same every time we iterate over it. Hence, we sample the sets in the `__init__` method, and return those in `__getitem__`. The code below implements exactly this dynamic.

In [30]: `class SetAnomalyDataset(data.Dataset):`

```

def __init__(self, img_feats, labels, set_size=10, train=True):
    """
    Inputs:
        img_feats - Tensor of shape [num_imgs, img_dim]. Represents the high-level features.
        labels - Tensor of shape [num_imgs], containing the class labels for the images
        set_size - Number of elements in a set. N-1 are sampled from one class, and one from another one.
        train - If True, a new set will be sampled every time __getitem__ is called.
    """

```

```

super().__init__()
self.img_feats = img_feats
self.labels = labels
self.set_size = set_size # The set size is here the size of correct images
self.train = train

```

```

# Tensors with indices of the images per class
self.num_labels = labels.max() + 1
self.img_idx_by_label = torch.argsort(self.labels).reshape(self.num_labels, -1)

if not train:
    self.test_sets = self._create_test_sets()

def _create_test_sets(self):

```

```

# Pre-generates the sets for each image for the test set
test_sets = []
num_imgs = self.img_feats.shape[0]
np.random.seed(42)
test_sets = [self.sample_img_set(self.labels[idx]) for idx in range(num_imgs)]
test_sets = torch.stack(test_sets, dim=0)
return test_sets

def sample_img_set(self, anomaly_label):
    """
    Samples a new set of images, given the label of the anomaly.
    The sampled images come from a different class than anomaly_label
    """

    # Sample class from 0,...,num_classes-1 while skipping anomaly_label as class
    set_label = np.random.randint(self.num_labels-1)
    if set_label >= anomaly_label:
        set_label += 1

    # Sample images from the class determined above
    img_indices = np.random.choice(self.img_idx_by_label.shape[1], size=self.set_size, replace=False)
    img_indices = self.img_idx_by_label[set_label, img_indices]
    return img_indices

def __len__(self):
    return self.img_feats.shape[0]

def __getitem__(self, idx):
    anomaly = self.img_feats[idx]
    if self.train: # If train => sample
        img_indices = self.sample_img_set(self.labels[idx])
    else: # If test => use pre-generated ones
        img_indices = self.test_sets[idx]

    # Concatenate images. The anomaly is always the last image for simplicity
    img_set = torch.cat([self.img_feats[img_indices], anomaly[None], dim=0])
    indices = torch.cat([img_indices, torch.LongTensor([idx]), dim=0])
    label = img_set.shape[0]-1

    # We return the indices of the images for visualization purpose. "Label" is the index of the anomaly
    return img_set, indices, label

```

Next, we can setup our datasets and data loaders below. Here, we will use a set size of 10, i.e. 9 images from one category + 1 anomaly. Feel free to change it if you want to experiment with the sizes.

```
In [31]:  
SET_SIZE = 10  
test_labels = torch.LongTensor(test_set.targets)
```

```
train_anom_dataset = SetAnomalyDataset(train_feats, train_labels, set_size=SET_SIZE, train=True)  
val_anom_dataset = SetAnomalyDataset(val_feats, val_labels, set_size=SET_SIZE, train=False)  
test_anom_dataset = SetAnomalyDataset(test_feats, test_labels, set_size=SET_SIZE, train=False)  
  
train_anom_loader = data.DataLoader(train_anom_dataset, batch_size=64, shuffle=True, drop_last=True, num_workers=4, pin_memory=True)  
val_anom_loader = data.DataLoader(val_anom_dataset, batch_size=64, shuffle=False, drop_last=False, num_workers=4)  
test_anom_loader = data.DataLoader(test_anom_dataset, batch_size=64, shuffle=False, drop_last=False, num_workers=4)
```

To understand the dataset a little better, we can plot below a few sets from the test dataset. Each row shows a different input set, where the first 9 are from the same class.

```
In [32]:  
def visualize_exmp(indices, orig_dataset):  
    images = [orig_dataset[idx][0] for idx in indices.reshape(-1)]  
    images = torch.stack(images, dim=0)  
    images = images * TORCH_DATA_STD + TORCH_DATA_MEANS  
  
    img_grid = torchvision.utils.make_grid(images, nrow=SET_SIZE, normalize=True, pad_value=0.5, padding=16)  
  
    plt.figure(figsize=(12,8))  
    plt.title("Anomaly examples on CIFAR100")  
    plt.imshow(img_grid)  
    plt.axis('off')  
    plt.show()  
    plt.close()  
  
    indices, _ = next(iter(test_anom_loader))  
    visualize_exmp(indices[:4], test_set)
```

Anomaly examples on CIFAR100



We can already see that for some sets the task might be easier than for others. Difficulties can especially arise if the anomaly is in a different, but yet visually similar class (e.g. train vs bus, flour vs worm, etc.).

After having prepared the data, we can look closer at the model. Here, we have a classification of the whole set. For the prediction to be permutation-equivariant, we will output one logit for each image. Over these logits, we apply a softmax and train the anomaly image to have the highest score/probability. This is a bit different than a standard classification layer as the softmax is applied over images, not over output classes in the classical sense. However, if we swap two images in their position, we effectively swap their position in the output softmax. Hence, the prediction is equivariant with respect to the input. We implement this idea below in the subclass of the Transformer Lightning module.

```
In [33]: class AnomalyPredictor(TransformerPredictor):
```

```
def calculate_loss(self, batch, mode="train"):  
    img_sets, labels = batch  
    preds = self.forward(img_sets, add_position  
    preds = preds.squeeze(dim=-1) # Shape: [Batch,
```

```

loss = F.cross_entropy(preds, labels) # Softmax/CE over set dimension
acc = (preds.argmax(dim=-1) == labels).float().mean()
self.log(f"{'{mode}_loss'", loss)
self.log(f"{'{mode}_acc'", acc, on_step=False, on_epoch=True)
return loss, acc

def training_step(self, batch, batch_idx):
    loss, _ = self._calculate_loss(batch, mode="train")
    return loss

def validation_step(self, batch, batch_idx):
    _ = self._calculate_loss(batch, mode="val")

def test_step(self, batch, batch_idx):
    _ = self._calculate_loss(batch, mode="test")

```

Finally, we write our train function below. It has the exact same structure as the reverse task one, hence not much of an explanation is needed here.

```

In [34]: def train_anomaly(**kwargs):
    # Create a PyTorch Lightning trainer with the generation callback
    root_dir = os.path.join(CHECKPOINT_PATH, "SetAnomalyTask")
    os.makedirs(root_dir, exist_ok=True)
    trainer = pl.Trainer(default_root_dir=root_dir,
                          callbacks=[ModelCheckpoint(save_weights_only=True, mode="max",
                                                      monitor="val_acc")],
                          accelerator="gpu" if str(device).startswith("cuda") else "cpu",
                          devices=1,
                          max_epochs=100,
                          gradient_clip_val=2)
    trainer.logger._default_hp_metric = None # Optional Logging argument that we don't need

# Check whether pretrained model exists. If yes, Load it and skip training
pretrained_filename = os.path.join(CHECKPOINT_PATH, "SetAnomalyTask.ckpt")
if os.path.isfile(pretrained_filename):
    print("Found pretrained model, loading...")
    model = AnomalyPredictor.load_from_checkpoint(pretrained_filename)
else:
    model = AnomalyPredictor(max_iters=trainer.max_epochs*len(train_anom_loader), **kwargs)
    trainer.fit(model, train_anom_loader, val_anom_loader)
    model = AnomalyPredictor.load_from_checkpoint(trainer.checkpoint_callback.best_model_path)

# Test best model on validation and test set

```