Closed book exam, no books, papers, notes, phones etc. allowed. The exam has 8 coding questions (labs), 8 open questions (lectures) and 14 multiple choice questions (papers). **Answer on the separate answer sheets.** Explain all answers, i.e. explicitly show intermediate steps to clarify if needed for coding / calculations / motivations / etc. Good luck!

We scan and crop answer boxes; please do not write outside answer boxes.

# 1   Lab assignments (38pts)

1. **Question** (4pts) Implement the forward pass of a shallow network with an input layer, one hidden layer, followed by a Sigmoid activation function. The forward pass equations are given below:

$$\mathbf{h} = \mathbf{x}\mathbf{W_1} + \mathbf{b_1}$$
$$\mathbf{y} = \mathbf{h}\mathbf{W_2} + \mathbf{b_2}$$
$$\mathbf{y} = \text{Sigmoid}(\mathbf{y})$$

Reminder:

$$\text{Sigmoid}(\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{x})}$$

Here $\mathbf{x}$ has dimensions [batch_size, input_channels] and $\mathbf{h}$ has dimensions [batch_size, hidden_dims] . Use the code template provided below. You may use elementary operations from the PyTorch library only, i.e. no predefined layers from torch.nn. Sub-questions:

A) Placeholder for layer weight and bias (2pts)

B) Forward pass (2pts)

```
import torch
class ShallowNet(object):
    """

    Args:
        in_feat: number of input features
        hidden_dims: number of hidden neurons
        out_feat: number of output features
    """

    def __init__(self, in_feat, out_feat):
        super(Linear, self).__init__()
        #################################
        # A) Your code goes here.
        #################################

        self.init_params() # Init. parameters

    def init_params(self):
        self.weight_1 = torch.randn_like(
            self.weight_1)
        self.bias_1 = torch.rand_like(
            self.bias_1)
        self.weight_2 = torch.randn_like(
            self.weight_2)
        self.bias_2 = torch.rand_like(
            self.bias_2)

    def forward(self, x):
        #################################
        # B) Your code goes here.
        #################################
        return y
```

2. **Question** (4pts) Implement the backward passes for the LeakyReLU and Tanh non-linearities. The non-linearities are defined as follows:

$$\text{LeakyReLU}(\mathbf{x}) = \max(0.01\mathbf{x}, \mathbf{x})$$
$$\text{Tanh}(\mathbf{x}) = \frac{\exp(\mathbf{x}) - \exp(-\mathbf{x})}{\exp(\mathbf{x}) + \exp(-\mathbf{x})}$$

You may use elementary operations from PyTorch library only, i.e. no predefined layers from torch.nn. Sub-questions:

A) LeakyReLU backward (2pts)

B) Tanh backward (2pts)

```
import torch
class LeakyReLU(object):

    def __init__(self):
        super(LeakyReLU, self).__init__()
        self.cache = None

    def forward(self, x):
        y = torch.clamp(x, min=0.01*x)
        self.cache = y
        return y

    def backward(self, dupstream):
        dupstream = dupstream.clone()
        #################################
        # A) Your code goes here.
        #################################
        return dx

class Tanh(object):

    def __init__(self):
        super(Tanh, self).__init__()
        self.cache = None

    def forward(self, x):
        y = (torch.exp(x) - torch.exp(-x))
            /(torch.exp(x) + torch.exp(-x))
        self.cache = y
        return y

    def backward(self, dupstream):
        #################################
        # B) Your code goes here
        #################################
        return dx
```

3. **Question** (6pts) Implement the following pooling methods for 2D inputs: You may use elementary operations from PyTorch library only, i.e. no predefined layers from torch.nn. Sub-questions:

A) Max Pooling (3pts)

B) Average Pooling (3pts)

```
1  import torch
2  class MaxPool2d(object):
3
4      def __init__(self, kernel_size, stride=1,
        padding=0):
5          self.kernel_size = kernel_size
6          self.stride = stride
7          self.padding = padding
8
9      def forward(self, x):
10         """
11         Args:
12             x: input tensor with shape of (N, C,
        H, W)
13         Returns:
14             y: output tensor with shape of (N, C
        , H', W') where
15                 H' = 1 + (H + 2 * padding -
        kernel_size) / stride
16                 W' = 1 + (W + 2 * padding -
        kernel_size) / stride
17         """
18         # Pad the input
19         x_padded = torch.nn.functional.pad(x, [
        self.padding] * 4)
20         # Unpack the needed dimensions
21         N, C, H, W = x.shape
22         KS = self.kernel_size
23         # Calculate output height and width
24         Hp = 1 + (H + 2 * self.padding - KS) //
        self.stride
25         Wp = 1 + (W + 2 * self.padding - KS) //
        self.stride
26         # Create an empty output to fill in.
27         # We combine first and second dim to
        speed up as we need no loop for each
28         # channel.
29         y = torch.empty((N*C, Hp, Wp), dtype=x.
        dtype, device=x.device)
30
31         ##################################
32         # A)    Your Code Here          #
33         ##################################
34
35         # Reshape output to seperate sample dim
        from channel dim since we
36         # combined them
37         y = y.reshape(N, C, Hp, Wp)
38
39         # Cache padded input to use in backward
        pass
40         self.cache = x
41
42         return y
43
44  class AveragePool2d(object):
45
46      def __init__(self, kernel_size, stride=1,
        padding=0):
47          self.kernel_size = kernel_size
48          self.stride = stride
49          self.padding = padding
50
51      def forward(self, x):
52          """
53          Args:
54              x: input tensor with shape of (N, C,
        H, W)
55
56          Returns:
57              y: output tensor with shape of (N, C
        , H', W') where
58                  H' = 1 + (H + 2 * padding -
```

```
        kernel_size) / stride
59                  W' = 1 + (W + 2 * padding -
        kernel_size) / stride
60          """
61
62
63          # Pad the input
64          x_padded = torch.nn.functional.pad(x, [
        self.padding] * 4)
65
66          # Unpack the needed dimensions
67          N, C, H, W = x.shape
68          KS = self.kernel_size
69
70          # Calculate output height and width
71          Hp = 1 + (H + 2 * self.padding - KS) //
        self.stride
72          Wp = 1 + (W + 2 * self.padding - KS) //
        self.stride
73
74          # Create an empty output to fill in.
75          # We combine first and second dim to
        speed up as we need no loop for each
76          # channel.
77          y = torch.empty((N*C, Hp, Wp), dtype=x.
        dtype, device=x.device)
78
79          ##################################
80          # B)    Your Code Here          #
81          ##################################
82
83          # Reshape output to seperate sample dim
        from channel dim since we
84          # combined them
85          y = y.reshape(N, C, Hp, Wp)
86
87          # Cache padded input to use in backward
        pass
88          self.cache = x
89
90          return y
```

4. **Question** (4pts) Implement the gradient update of the Adam optimizer, given by:

$$v_i = \rho_1 v_{i-1} + (1 - \rho_1) \nabla_\theta$$

[handwritten: rho ↓, prev. ↓, ← grad]

$$\hat{v}_i = \frac{v_i}{1 - \rho_1^i}$$

$$r_i = \rho_2 r_{i-1} + (1 - \rho_2) \nabla_\theta^2$$

$$\hat{r}_i = \frac{r_i}{1 - \rho_2^i}$$

$$\theta' = \theta - \epsilon \frac{\hat{v}_i}{\sqrt{\hat{r}_i + \delta}}$$

[handwritten: ← delta]

```
1  def adam(X, rhos, learning_rate, prev_values,
   index, Grad=Grad_f):
2      """
3      Adam optimization step.
4      Args:
5          X: Current value of objective function.
6          rhos: Optimization hyperparameter - see
   formula above.
7          learning_rate: Optimization step size.
8          prev_value: Momentum parameter from
   previous iteration.
9          index: Optimization step counter.
10         Grad: Gradient of quadratic function.
11     """
12
```

```
13    delta = 1e-5              # Tiny amount
      to prevent division by zero
14    gradient = Grad(*X)       # Gradient of
      current values
15    rho_v, rho_r = rhos       # Rho values
      for momentum & rmsProp part of Adam
16    v_prev, r_prev = prev_values # Adam
      parameters from previous iterations
17
18    v = r = 0                 # Adam
      paramters for momentum & rmsProp
19    v_bc = r_bc = 0           # Bias
      corrected adam parameters
20
21    ####################################
22    #   TODO: Create gradient update   #
23    #with Adam: update v, v_bc, r, r_bc #
24    #and X.                            #
25    ####################################
26    #        Your Code Here            #
27    ####################################
28    return X,(v,r)
```

5. **Question** (4pts) Implement the following regularization methods:

   - L2 Regularization

   - Early Stopping

```
1  def train_wd(train_loader, net, optimizer,
      criterion, wd):
2      """
3      Args:
4          train_loader: Data loader.
5          net: Neural network model.
6          optimizer: Optimizer (e.g. SGD).
7          criterion: Loss function
8          wd: Weight decay (L2 penalty)
9      """
10
11     avg_loss = 0
12     correct = 0
13     total = 0
14
15     for i, data in enumerate(train_loader):
16         inputs, labels = data
17         optimizer.zero_grad()
18         outputs = net(inputs)
19         loss = criterion(outputs, labels)
20         ####################################
21         #        A) Your Code Here         #
22         ####################################
23         loss.backward()
24         optimizer.step()
25         avg_loss += loss
26         _, predicted = torch.max(outputs.data,
      1)
27         total += labels.size(0)
28         correct += (predicted == labels).sum().
      item()
29
30     return avg_loss/len(train_loader), 100 *
      correct / total
31
32
33
34
35
36
37     #########EARLY STOPPING#########
38 net = FCNet()
39 criterion = nn.CrossEntropyLoss()
```

```
40  optimizer = optim.SGD(net.parameters(), lr=5e-1,
       weight_decay=3e-3)
41
42  # Set the number of epochs to for training
43  epochs = 100
44
45  # Patience - how many epochs to keep training
       after accuracy has not improved
46  patience = 0
47
48  # Initialize early stopping variables
49  val_acc_best = 0
50  patience_cnt = 0
51
52  for epoch in tqdm(range(epochs)):  # loop over
       the dataset multiple times
53      train_loss, train_acc = train(train_loader,
       net,optimizer,criterion)
54      val_loss, val_acc = test(val_loader,net,
       criterion)
55
56      writer.add_scalars("Loss", {'Train':
       train_loss, 'Val': val_loss}, epoch)
57      writer.add_scalars('Accuracy', {'Train':
       train_acc, 'Val': val_acc}, epoch)
58
59      ####################################
60      #         B) Your Code Here        #
61      ####################################
```

6. **Question** (6pts) Implement the GRU. [Hint]: think about all the architectural hyperparameters and especially sizes of the input, hidden states, weights and outputs.

The update rule is given as:

$$r_t = \sigma(W_{xr}x_t + b_{xr} + W_{hr}h(t-1) + b_{hr})$$
$$z_t = \sigma(W_{xz}x_t + b_{xz} + W_{hz}h(t-1) + b_{hz})$$
$$n_t = \tanh(W_{xn}x_t + b_{xn} + r_t \odot (W_{hn}h(t-1) + b_{hn}))$$
$$h_t = (1 - z_t) \odot n_t + z_t \odot h(t-1)$$

Sub-questions:

A) Parameter initialization (2pts)

B) Forward pass implementation (4pts)

```
1  class GRU(nn.Module):
2      def __init__(self, input_size, hidden_size):
3          super(GRU, self).__init__()
4          self.hidden_size = hidden_size
5          self.weight_xh = None
6          self.weight_hh = None
7          self.bias_xh = None
8          self.bias_hh = None
9
10         ####################################
11         # YOUR CODE HERE
12         ####################################
13
14         # Initialize parameters
15         self.reset_params()
16
17     def reset_params(self):
18         std = 1.0 / math.sqrt(self.hidden_size)
19         self.weight_xh.data.uniform_(-std, std)
20         self.weight_hh.data.uniform_(-std, std)
21         self.bias_xh.data.uniform_(-std, std)
22         self.bias_hh.data.uniform_(-std, std)
23
```

```
24    def forward(self, x):
25        """
26        Args:
27            x: input with shape (N, T, D)
28            where N is number of samples,
29            T is number of timestep and
30            D is input size which must be equal
31            to self.input_size.
32
33        Returns:
34            y: output with a shape of (N, T, H)
35            where H is hidden size
36        """
37
38        # Transpose input for efficient
        vectorized calculation. After transposing
        the input will have shape(T, N, D).
39        x = x.transpose(0, 1)
40        T, N, H = x.shape[0], x.shape[1], self.
        hidden_size
41        h0 = torch.zeros(N, H, device=x.device)
42
43        # Define a list to store outputs. We
        will then stack them.
44        y = []
45
46        ##################################
47        # TODO: Implement GRU forward pass   #
48        ##################################
49
50        # Stack the outputs. After this
        operation, output will have shape of
51        # (T, N, H)
52        y = torch.stack(y)
53
54        # Switch time and batch dimension, (T, N
        , H) -> (N, T, H)
55        y = y.transpose(0, 1)
56        return y
```
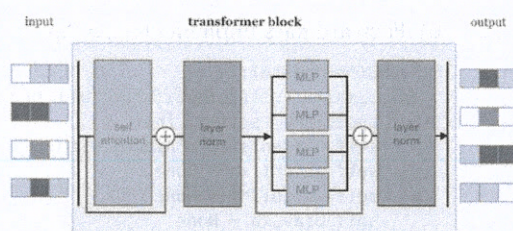
7. **Question** (4pts) The transformer block consists of a self-attention layer, followed by layer norm, a MLP applied on each vector individually and another layer norm. Note the residual connections in the self-attention and MLP layer. The architecture is given below:



Now implement the transformer block as a PyTorch layer. All different components are already defined in init function - your task is to connect them in the proper way in the forward method. Ima Sub-questions:

A) Perform the forward pass of a transformer block as depicted in the image.

```
1  import torch.nn as nn
2  class TransformerBlock(nn.Module):
3      def __init__(self, k, heads):
4          """
5          Args:
```

```
6            k: embedding dimension
7            heads: number of heads (k mod heads
   must be 0)
8          """
9          super(TransformerBlock, self).__init__()
10
11         self.att = MultiHeadAttention(k, heads=
   heads)
12         self.norm1 = nn.LayerNorm(k)
13         self.ff = nn.Sequential(
14             nn.Linear(k, 4 * k),
15             nn.ReLU(),
16             nn.Linear(4 * k, k))
17         self.norm2 = nn.LayerNorm(k)
18
19     def forward(self, x):
20         """
21         Args:
22             x: input with shape of (b, k)
23         Returns:
24             y: output with shape of (b, k)
25         """
26         ############################
27         # TODO: Perform the forward  #
28         # pass of a transformer block#
29         # as depicted in the image.  #
30         ############################
31         return y
```

8. **Question** (6pts) Implement a Variational Autoencoder (VAE) for the following code block.

- KL Loss (2pt)

- Reparametrization (2pt)

- Forward function: creating mean and variance (2pt)

```
1  #encoder
2  class VarEncoder(nn.Module):
3      def __init__(self, latent_dims, s_img, hdim)
   :
4          super(VarEncoder, self).__init__()
5
6          #layers for g1
7          self.linear1_1 = nn.Linear(s_img*s_img,
   hdim[0])
8          self.linear2_1 = nn.Linear(hdim[0], hdim
   [1])
9          self.linear3_1 = nn.Linear(hdim[1],
   latent_dims)
10
11         #layers for g2
12         self.linear1_2 = nn.Linear(s_img*s_img,
   hdim[0])
13         self.linear2_2 = nn.Linear(hdim[0], hdim
   [1])
14         self.linear3_2 = nn.Linear(hdim[1],
   latent_dims)
15
16         self.relu    = nn.ReLU()
17
18         #distribution setup
19         self.N = torch.distributions.Normal(0,
   1)
20         self.N.loc = self.N.loc.to(try_gpu()) #
   hack to get sampling on the GPU
21         self.N.scale = self.N.scale.to(try_gpu()
   )
22         self.kl = 0
23
24         ##################################
```
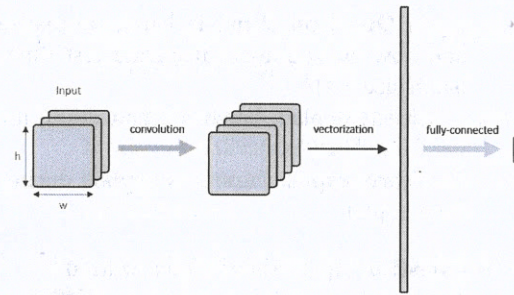
```
25    #    TODO: Define function for:   #
26    #  A) the Kullback-Leibner        #
27    #     loss "kull_leib"            #
28    #  B) the Reparameterization trick #
29    ###################################
30    #         Your code here          #
31    ###################################
32
33    def forward(self, x):
34
35
36        ###################################
37        # TODO: Create mean and variance   #
38        ###################################
39        # C) Forward Function your code here
40        ###################################
41
42        #reparameterize to find z
43        z = self.reparameterize(mu, sig)
44
45        #loss between N(0,I) and learned
       distribution
46        self.kl = self.kull_leib(mu, sig)
47
48        return z
49
50  #decoder: same as before
51
52  #autoencoder
53  class VarAutoencoder(nn.Module):
54      def __init__(self, latent_dims, s_img, hdim
       = [100, 50]):
55          super(VarAutoencoder, self).__init__()
56
57          self.encoder = VarEncoder(latent_dims,
       s_img, hdim)
58          self.decoder = Decoder(latent_dims,
       s_img, hdim)
59
60      def forward(self, x):
61
62          z = self.encoder(x)
63          y = self.decoder(z)
64
65          return y
```

## 2   Lectures (38pt)

1. **Question** (4pts) Consider a 2-layered network. The first layer is convolutional with a kernel size of $5 \times 5$, with padding = 1, stride = 2. The second layer is a fully connected layer. The number of input channels is 3, the number of hidden channels is 7. The number of output features is 2. The network uses bias terms. The size of an input image is: $h \times w = 12 \times 12$. The image below clarifies the network further. **What is the total number of learnable parameters?** Motivate your answer with a detailed step-by-step approach. (Don't forget the bias terms).



2. **Question** (6pts) Calculate MLE (Maximum Likelihood Estimator) for a Possion distribution.

Consider that we are given a Poisson probability distribution given by:

$$P_{model}(X = x) = \frac{\lambda^x e^{-\lambda}}{x!}$$

We draw $m$ samples from the probability distribution namely : $\{x_1, x_2 \ldots x_m\}$. The objective is to find the maximum likelihood estimate for the parameter $\lambda$ as a function of the data samples.

Sub-questions:
1) (1pt) Write the likelihood function for the given probability distribution.
2) (2pt) Write the log-likelihood function by using the logarithm operator on the function obtained in the previous step. Make sure to simplify the terms as much as possible.
3) (2pt) Calculate the derivative of the natural log likelihood function with respect to $\lambda$
4) (1pt) Set the derivative equal to zero and solve for $\lambda$

Hint: use the log-likelihood, defined as

$$MLE(X) = \arg\max_{\theta} \prod_{i=1}^{m} \log P_{model}(x_i, \theta)$$

where $X = \{x_1 \ldots x_m\}$

3. **Question** (6pts)

Calculate the receptive field of a feature/pixel in the output of the architecture given by the table below. Provide a step-by-step explanation.

**Note**: the receptive field here refers to the number of pixels in the input image that a particular feature ("pixel") in the output of Conv4 is looking at, i.e. the answer should be a single integer.

| Layer | Kernel size | Stride |
|-------|-------------|--------|
| Conv1 | 3           | 1      |
| Pool1 | 2           | 2      |
| Conv2 | 3           | 1      |
| Pool2 | 2           | 2      |
| Conv3 | 3           | 1      |
| Conv4 | 3           | 1      |

4. (4pts) **Question** What is batch normalization and how do you apply it during test time (e.g. batch size = 1)?
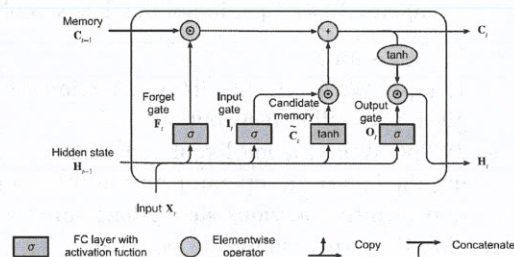   1. Clearly explain batch normalization, including formulas. (2pts)
   2. Clearly explain how it is applied during test time. (2pts)

5. **Question** (4pts) About regularization.
   1. Draw the learning curve for overfitting. (1pt)
   2. Name and explain one type of parameter norm regularization (include its equation). (1pt)
   3. How is dropout performed? Explain both training and evaluation. (2pt)

6. **Question** (6pts) Long Short Term Memory (LSTM).
   Using equations, explain the working of the LSTM, i.e. current hidden state based on the current memory cell state $(C_t)$, previous hidden state $(h_{t-1})$ and the current input $(X_t)$.



   1. Equation of the forget gate $(f_t)$. (1pt)
   2. Equation of the input gate $(i_t)$. (1pt)
   3. Equation of the candidate cell gate $(c_t)$. (2pt)
   4. Equation of current hidden state $(h_t)$. (2pt)

7. **Question**(4pts) There are two types of positional information used in self-attention: embeddings or encodings.

   1. Why do we need to explicitly include positional information in self-attention? (2pt)
   2. What is the difference between positional embeddings and positional encodings? (2pt)

8. **Question** (4pts) Contractive auto-encoder.

   1. What is an auto-encoders and what is it used for? (2pt)

   A contractive auto-encoder is given by

   $$l(g(f(x)), x) + \lambda \left\| \frac{\partial f(x)}{\partial x} \right\|_F^2$$

   2. Explain the equation (terms, operators), and its goal. (2pt)

## 3　Papers (14pts)

A 1pt multiple-choice question per paper. Select the single best fitting answer per question.

1. Paper: *A Step Toward Quantifying Independently Reproducible Machine Learning Research.* What statement matches the paper best:

   A) Releasing code is very important to reproduce ML research.

   B) Looking at the released code by the authors tends to improve reproduction quality.

   C) Paper reproduction rates have not significantly changed over the past 35 years

   D) More equations tends to improve reproduction quality.

2. Paper: *Troubling Trends in Machine Learning Scholarship.* The paper mentions "Failure to identify the sources of empirical gains". What exactly is meant by that?

   A) Gains come from using additional data instead of the proposed method

   B) Gains are not tested statistical for significance

   C) Gains come from hyper-parameter tuning instead of the proposed method

   D) Gains are tuned on the test set

3. Paper: *Do ImageNet Classifiers Generalize to ImageNet?* What is their main motivation?

   A) That repeated re-use of the test set leads to overfitting

   B) That the data collection procedure is not reproducible

   C) That the 1,000 ImageNet classes are too specialized

   D) That the concept of generalization is too narrowly defined

4. Paper: *Scaling down deep learning.* What do they scale down?

   A) The number of layers

   B) The dataset

   C) The weights in a layer

   D) The training epochs

5. Paper: *Highway and Residual Networks learn Unrolled Iterative Estimation.* Whats the difference between highway vs residual?

   A) Highway simplifies the residual

   B) Highway is a gated variant of Residual

C) Residual is a different name for Highway

D) Residual is specially adapted for images

6. Paper: *ResNet strikes back: An improved training procedure in timm.* What does the author hope to achieve?

A) A better ResNet baseline

B) Better understanding of the ResNet

C) An improved ResNet architecture

D) Much faster training times

7. Paper: *Deep Image Prior.* What statement fits best with the paper:

A) Neural network optimizers are a form of prior knowledge

B) ConvNets are useful for images, even without training

C) Optimization is less important than the network architecture

D) ConvNets can learn to denoise images

8. Paper: *Approximating CNNs with Bag-of-local-Features models works surprisingly well on ImageNet.* What statement best describes the paper?

A) Its based on shallow networks

B) The receptive field is artificially enlarged

C) The class evidence used for explainability is locally bagged

D) The receptive field is artificially reduced

9. Paper: *Group Normalization.* The statistics are computed over groups of pixels. The grouping is done by:

A) Grouping inside a single featuremap

B) Grouping activations

C) Grouping parameters

D) Grouping residuals

10. Paper: *Torch.manual_seed(3407) is all you need: On the influence of random seeds in deep learning architectures for computer vision.* What is the author surprised about?

A) That with not very large variance its easy to find an outlier that performs much better or much worse than the average.

B) That even though the variance is small, that the differences are still statistically significant.

C) That the random seed plays such a big role in finding the best results.

D) That the random seed influences so many things (initialization, batch elements, gradient steps, etc.) yet does account for so little differences.

11. Paper: *Attention Is All You Need.* Select the best answer.

A) For each token, the query, key and value parameters are identical

B) The connections between the losses are gated, similar to a GRU

C) One token is best seen as a set.

D) Without positional embeddings, it would exploit the absolute token position, similar to a MLP.

12. Paper: *Perceiver IO: A General Architecture for Structured Inputs & Outputs.*

A) Uses standard self-attention, and it is outperformed by specialized solutions

B) Uses a special form of self-attention, and it is outperformed by specialized solutions

C) Uses standard self-attention and it outperforms specialized solutions

D) Uses a special form of self-attention and it outperforms specialized solutions

13. Paper: *Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks.* What does the "cycle" refer to?

A) Cyclic featuremap padding

B) Training data cycles

C) Cycling between two annotations of an image

D) Cycling between domains

14. Paper: *GANORCON: Are Generative Models Useful for Few-shot Segmentation?* So, are GANs useful for few-show segmentation?

A) Yes, more useful than contrastive learning.

B) Yes, but not as useful as contrastive learning

C) No, but contrastive learning is

D) No, and contrastive learning also is not