

# FYS-3033/8033 Home Exam - Spring 2026

Anonymous Full Paper  
Submission 79

## 1 Problem 1

### 1.1 Task 1a)

We can interpret equation 1 in the exam sheet as the sum of two expected values:

$$V(G, D) = \mathbb{E}_{x \sim p_{data}(x)}[\log D(x)] + \mathbb{E}_{z \sim p_z(z)}[\log(1 - D(g(z)))] \quad (1)$$

Then assuming that we have reached the optimal solution, where the generator-distribution is equal to the data distribution we can write:

$$\begin{aligned} V(G, D) &= \mathbb{E}_{x \sim p_{data}(x)}[\log D(x)] \\ &+ \mathbb{E}_{x \sim p_g(x)}[\log(1 - D(x))] \\ &= \int_x [p_{data}(x) \log D(x) \\ &+ p_g(x) \log(1 - D(x))] dx \end{aligned} \quad (2)$$

We will now maximize the integrand, which we denote as  $y$ , with respect to the discriminator. We take the derivative w.r.t.  $D$  and set equal to zero:

$$\frac{\partial y}{\partial D} = \frac{p_{data}(x)}{D(x)} - \frac{p_g(x)}{1 - D(x)} = 0 \quad (4)$$

$$\implies 0 = (1 - D(x))p_{data}(x) - p_g(x)D(x) \quad (5)$$

$$D^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)} \quad (6)$$

With our assumption that  $p_g = p_{data}$  we get the final result that the optimal discriminator value is  $D(x) = \frac{1}{2}$ . Inserting this value into 4 gives:

$$V(G, D) = -\log 2 \left( \int_x p_{data}(x) dx + \int_x p_g(x) dx \right) \quad (7)$$

Note that both of the integrals have to be equal to one, since we are integrating PDFs over their entire range. Thus, the optimal loss is  $V(G, D) = -2 \log 2 = -\log 4$ . Note that the optima is what we one gets when minimizing the Jensen-Shannon divergence between  $p_{data}$  and  $p_g$ .

### 1.2 Task 1b)

In order to train our VAE we need to have a way of describing the latent variable  $z$  as a deterministic function (so that we can differentiate and perform backpropagation), instead of a stochastic sampling.

This is known as the reparametrization trick and is given by,

$$z = \mu_{z|x} + \sigma_{z|x} \epsilon \quad (8)$$

where  $\epsilon \sim \mathcal{N}(0, 1)$ . So instead of sampling  $z \sim q_\phi(z|x)$ , we now have a deterministic function which lets the gradients flow through  $\mu_{z|x}$  and  $\sigma_{z|x}$ . Note that we can rewrite 8 as:

$$\epsilon = \frac{z - \mu_{z|x}}{\sigma_{z|x}} \quad (9)$$

Since  $\epsilon \sim \mathcal{N}(0, 1)$ , this implies that  $z \sim \mathcal{N}(\mu_{z|x}, \sigma_{z|x}^2)$ , because the right hand side of the equation is normalizing  $z$ .

### 1.3 Task 1c)

The Kullback-Leibler divergence is given by: [1]

$$KL(p, q) = \int_x p(x) \log \frac{p(x)}{q(x)} dx \quad (10)$$

We start by computing  $p(x)/q(x)$ :

$$\frac{p(x)}{q(x)} = \frac{\sigma_2}{\sigma_1} \exp \left[ -\frac{1}{2} \left( \frac{(x - \mu_1)^2}{\sigma_1^2} - \frac{(x - \mu_2)^2}{\sigma_2^2} \right) \right] \quad (11)$$

where we used the the density functions for the normal distributions [2]. Taking the logarithm on both sides, and then applying  $\mathbb{E}_{p(x)}[*]$  on both sides gives:

$$KL(p, q) = \log \frac{\sigma_2}{\sigma_1} - \frac{\mathbb{E}[(x - \mu_1)^2]}{2\sigma_1^2} + \frac{\mathbb{E}[(x - \mu_2)^2]}{2\sigma_2^2} \quad (12)$$

Now note that  $\mathbb{E}[(x - \mu_1)^2] = \sigma_1^2$  by definition. We can rewrite,

$$\mathbb{E}[(x - \mu_2)^2] = \mathbb{E}[(x - \mu_1 + \mu_1 - \mu_2)^2] \quad (13)$$

$$= \mathbb{E}[(x - \mu_1)^2 + (\mu_1 - \mu_2)^2] \quad (14)$$

$$= \sigma_1^2 + (\mu_1 - \mu_2)^2 \quad (15)$$

Inserting this into 12 gives the final expression:

$$KL(p, q) = \log \frac{\sigma_2}{\sigma_1} - \frac{1}{2} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} \quad (16)$$

The goal of VAEs is to learn some underlying prior data distribution  $p(x)$ , so that we can generate new samples from the distribution. To do this we generate samples from a latent representation instead of directly from the other pixel values or tokens:

$$p_\theta(x) = \int p_\theta(z) p_\theta(x|z) dz \quad (17)$$

069 where  $p_\theta(z)$  is a self-chosen prior distribution and  
 070  $p_\theta(x|z)$  is the "code" that takes us from the latent-  
 071 space to the "data-space". The integral is in-  
 072 tractable, so we have to estimate it instead of com-  
 073 puting it directly. We do this by replacing  $p_\theta(z|x)$   
 074 with some simple distribution  $q_\phi(z|x)$  (in Bayes' for-  
 075 mula). Note that we now have two distributions  
 076 that we can choose:  $p_\theta(z)$  and  $q_\phi(z|x)$ . In order to  
 077 find the optimal  $q_\phi$ , one can optimize the ELBO: [3]

$$078 \mathcal{L} = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - KL(q_\phi(z|x)||p_\theta(z)) \quad (18)$$

079 By making both the distributions Gaussian, we can use  
 080 16 as a differentiable loss function to push the  
 081 mean and standard deviation of  $q_\phi$  closer to the  
 082  $(\mu_z(x), \sigma_z(x))$  (since the first part of the equation is  
 083 "set").

## 084 2 Problem 2

### 085 2.1 Task 2a)

086 The BERT model [4] has a self-attention mechanism,  
 087 which means it can capture relationships between  
 088 tokens across the entire input. Unlike for instance  
 089 RNNs, the BERT model processes all the tokens  
 090 in the input sequence simultaneously, which lets us  
 091 parallelize the computation (and take advantage of  
 092 the provided GPUs). When inputting a sequence  
 093 into the BERT-model, the sequence will always start  
 094 with the special [CLS] token. The purpose of this  
 095 token, is that through every layer of self-attention  
 096 the [CLS] token will attend to all the other tokens  
 097 and opposite (bidirectional). Thus, the [CLS] token  
 098 will essentially contain a compressed summary of  
 099 the input sequence, which can then be used for clas-  
 100 sification. This fits our 4-class classification problem  
 101 perfectly. Although the original paper [4] specifically  
 102 states that **pre-trained** BERT-models are able to  
 103 obtain state-of-the-art performance on language in-  
 104 ference tasks, our dataset should be simple and large  
 105 enough to obtain good results without pre-training  
 106 and thus it is a good starting point.

### 107 2.2 Task 2b)

108 For this task we will use the pre-built *BertForSe-*  
 109 *quenceClassification* architecture [5]. Using the fol-  
 110 lowing code, we initialize the model with random  
 111 weights, instead of the pre-trained ones:

```
1 from transformers import BertConfig, BertForSequenceClassification,
  ↳ BertTokenizer
2 import torch
3 from torch.utils.data import DataLoader, Dataset
4
5 #Using BertTokenizer for vocabulary
6 tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
7 #initializing random weights for Bert model
8 config = BertConfig(num_labels = 4)
9 model = BertForSequenceClassification(config)
```

113 Note that we also use the pre-built BertTokenizer  
 114 as our "vocabulary". We tokenize the data using  
 115 the following function:

```
1 def tokenize_data(texts, labels, tokenizer):
2     data = []
3     #Looping over all samples and tokenizing with BertTokenizer
4     for txt, label in zip(texts, labels):
5         #max length 128 should be enough for news headlines
6         encoded = tokenizer(txt, truncation=True, max_length=128,
  ↳ padding="max_length", return_tensors="pt")
7         data.append({"input_ids": encoded["input_ids"].squeeze(),
8                     "attention_mask": encoded["attention_mask"].squeeze(),
9                     "labels": torch.tensor(label)})
10    return data
```

116

Note that we use a max length of 128 tokens,  
 which should be sufficient for this particular case.  
 Next we do a 80/20 train-validation-split, before  
 creating batches (size 32) using the PyTorch Data-  
 Loader. As our optimizer we use AdamW [6] with  
 a learning rate  $\gamma = 10^{-5}$ , which is good for BERT  
 models due to the way it applies weight decay separ-  
 atly from the adaptive updates. We train the model  
 for a total of 5 epochs. The training loop is in the  
 following code,

```
1 #Training loop
2 epochs = []
3 train_loss = []
4 val_acc = []
5 val_loss = []
6 for epoch in tqdm(range(0, n_epochs)):
7     model.train()
8     curr_loss = 0
9
10    for batch in train_loader:
11        #Resetting gradients for batch
12        optimizer.zero_grad()
13        #Retrieving data from batch
14        input_ids = batch["input_ids"].to(device)
15        attention_mask = batch["attention_mask"].to(device)
16        labels = batch["labels"].to(device)
17
18        out = model(input_ids=input_ids, attention_mask=attention_mask,
  ↳ labels=labels)
19        #Computing loss, applying backpropagation and taking a GD-step
20        batch_loss = out.loss
21        batch_loss.backward()
22        optimizer.step()
23
24        curr_loss += batch_loss.item()
25
26    train_loss.append(curr_loss / len(train_loader))
27    epochs.append(epoch+1)
28
29    #Testing on validation set
30    model.eval()
31    curr_val_loss = 0
32    correct_preds = 0
33    total_samples = 0
34
35    with torch.no_grad():
36        for batch in val_loader:
37            input_ids = batch["input_ids"].to(device)
38            attention_mask = batch["attention_mask"].to(device)
39            labels = batch["labels"].to(device)
40            out = model(input_ids=input_ids,
  ↳ attention_mask=attention_mask, labels=labels)
41            batch_loss = out.loss
42            curr_val_loss += batch_loss.item()
43
44            #Predictions: argmax of the logits of 4-class output vector =>
  ↳ highest probability
45            preds = out.logits.argmax(dim=1)
46            correct_preds_batch = sum(preds==labels).item()
47            correct_preds += correct_preds_batch
48            total_samples += batch_size
49
50    accuracy = correct_preds / total_samples
51    val_acc.append(accuracy)
52    val_loss.append(curr_val_loss / len(val_loader))
```

127

In figure 1 we see how the training and validation  
 loss/accuracy evolves with the number of epochs.  
 The performance on the validation set is actually  
 the best after only two epochs, but the training  
 loss gets significantly better after more epochs. A  
 validation accuracy of **91.5%** is satisfactory for our  
 initial model without pre-training.

### 2.3 Task 2c)

In order to use pre-trained weights from Huggingface,  
 the only change we have to make before training is  
 the following line:

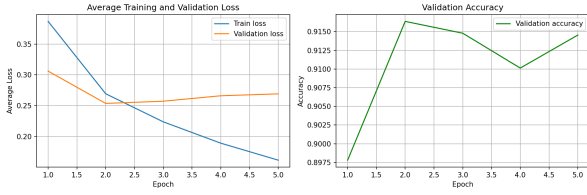


Figure 1. Left: Training and validation loss plotted as a function of epochs. Right: Validation accuracy plotted as a function of epochs.

```

1 model = BertForSequenceClassification.from_pretrained(
2     "bert-base-uncased",
3     num_labels=4
4 )
    
```

139

140 which means we are **not** overwriting the pre-  
 141 trained weights with random weights. Since we  
 142 are now finetuning weights that we know are good  
 143 for general NLP tasks we expect at least a slight  
 144 increase in performance. As we can see in figure 2  
 145 we do in fact get a significant boost in the validation  
 146 accuracy, which now peaks at **94.7%**. Once again  
 147 we note that the model is starting to overfit after  
 148 3 epochs, so going forward we will only be using 3  
 149 epochs to avoid this and to also save computation-  
 150 time. The total time for training and inference was  
 151 5 hours and 51 minutes. Note that we used the same  
 152 hyperparameters as for the model we trained from  
 scratch.

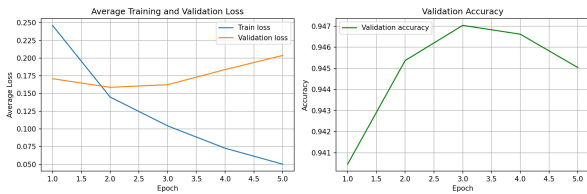


Figure 2. Training and validation loss/accuracy as a function of epochs, with pre-trained weights.

153

## 154 2.4 Task 2d)

155 To extract the representation from the 4 last trans-  
 156 former layers, we start by enabling hidden state  
 157 outputs in our model. We then pass the entire  
 158 validation set through the model and save the repre-  
 159 sentations from these 4 layers. We try two different  
 160 dimensionality reduction methods: (1) UMAP [7]  
 161 and (2) Principal Component Analysis (PCA) [8].  
 162 The results from the UMAP reduction is shown in  
 163 figure 3, while the PCA reduction is shown in 4. In  
 164 the UMAP representation we clearly see that there  
 165 is already minimal overlap between the classes in the  
 166 3rd last layer, which suggests that we can perhaps  
 167 drop the last two layers. The PCA representation  
 168 on the other hand shows a clear improvement in  
 169 separability between each layer, so based only on  
 170 that it would be hard to justify dropping any of the  
 171 layers.

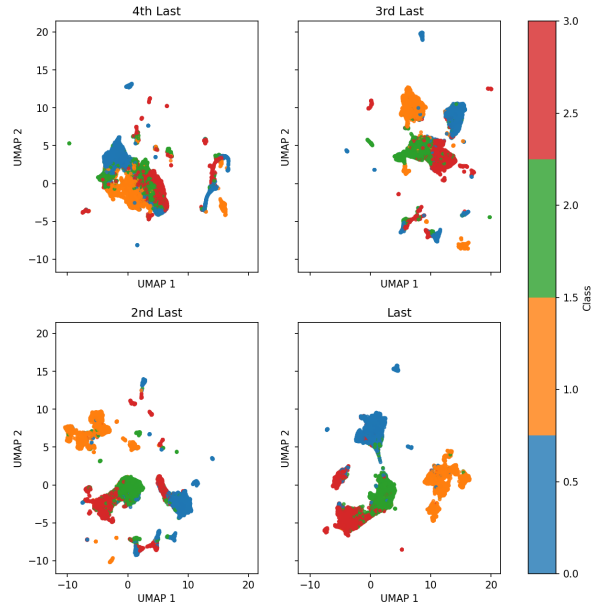


Figure 3. UMAP representation of the 4 last trans-  
 former layers.

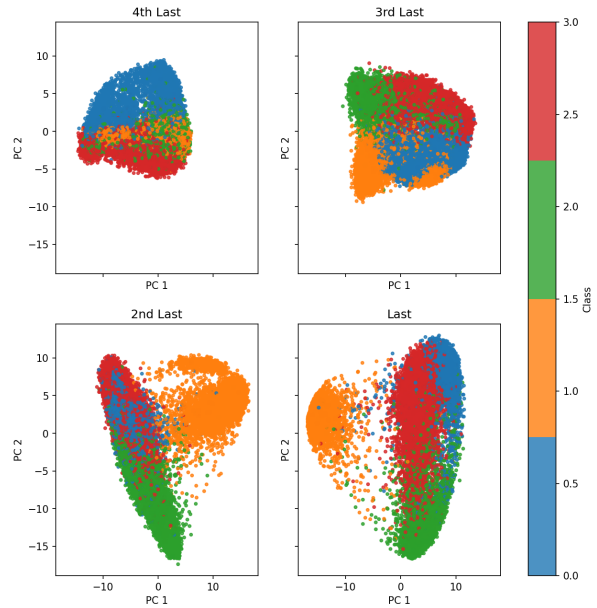


Figure 4. PCA representation of the 4 last trans-  
 former layers.

## 172 2.5 Task 2e)

173 Based on the UMAP representations from the pre-  
 174 vious task, we will try to remove the last two layers  
 175 of the transformer. According to [9], the top-layer  
 176 dropping strategy with two layers dropped generally  
 177 only slightly reduces performance for most of the  
 178 NLP tasks tested in the paper. We used to following  
 179 code to setup our slim version of the BERT model:

```

1 #Setting up a "slim" BERT model with 10 layers
2 config = BertConfig.from_pretrained("bert-base-uncased")
3 config_num_hidden_layers = 10
4 config_num_labels = 4
5 #10 layer Encoder with pretrained weights
    
```

180

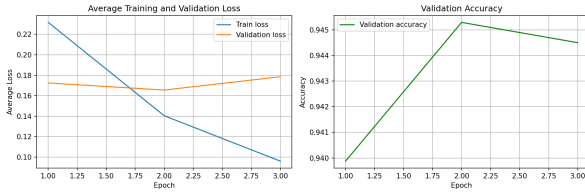
```

6 encoder_slim = BertModel.from_pretrained("bert-base-uncased",
  ↪ config=config)
7 #Wrapping encoder in BertForSequenceClassification model
8 model = BertForSequenceClassification(config=config)
9 model.bert = encoder_slim

```

181

182 In figure 5 we see the results from the training  
 183 and inference. We note that the validation accuracy  
 184 peaked at **94.55%** after two epochs, which is only a  
 185 drop of 0.15% compared to the full model. The total  
 186 training and inference time was also significantly  
 187 reduced to about **34 minutes** per epoch, compared  
 188 to 70 minutes per epoch for the full model. Note  
 189 that due to GPU constraints we had to reduce the  
 190 batch size from 32 to 16 for this particular run, but  
 191 all the other hyperparameters are the same as for  
 192 the regular BERT model. Our results show that we  
 193 can cut the training time in (roughly) half, while  
 194 barely dropping in performance.



**Figure 5.** Training and validation loss/accuracy as a function of epochs for "slim" model with last two layers dropped.

## 195 2.6 Task 2f)

196 Monte Carlo Dropout [10] approximates the stan-  
 197 dard dropout method as Bayesian inference, which  
 198 can then be used to model uncertainty. During test-  
 199 ing one performs  $T$  stochastic forward passes for  
 200 the same input  $\mathbf{x}^*$ , each with a separately sampled  
 201 dropout mask. We can then think of the  $T$  outputs  
 202 of the network,  $\{\hat{\mathbf{y}}^{(t)}(\mathbf{x}^*)\}_{t=1}^T$ , as Monte Carlo sam-  
 203 ples from the approximate predictive distribution,

$$204 \quad q(\mathbf{y}^*|\mathbf{x}^*) = \int p(\mathbf{y}^*|\mathbf{x}^*, \omega)q(\omega)d\omega. \quad (19)$$

205  $q(\omega)$  is the dropout-induced approximate posterior  
 206 of the weights. According to [10] we can approximate  
 207 the variance (uncertainty) using the formula:

$$208 \quad \text{Var}(\mathbf{y}^*) \approx \tau^{-1}\mathbf{I}_D + \frac{1}{T} \sum_t \hat{\mathbf{y}}^{(t)}\hat{\mathbf{y}}^{(t)T} - E[\mathbf{y}^*]E[\mathbf{y}^*]^T, \quad (20)$$

209 where  $\tau$  is the precision hyper-parameter determined  
 210 by,

$$211 \quad \tau = \frac{pl^2}{2N\lambda}. \quad (21)$$

212  $\lambda$  is weight decay,  $p$  is dropout-probability,  $N$  is  
 213 dataset size and  $l$  is the prior length-scale. We can  
 214 also use the following approximation to compute  
 215  $E[\mathbf{y}^*]$ ,

$$216 \quad E[\mathbf{y}^*] \approx \frac{1}{T} \sum_{i=1}^T \hat{\mathbf{y}}^*(\mathbf{x}^*, \mathbf{W}_1^t, \dots, \mathbf{W}_L^t). \quad (22)$$

Note that the expected values and variance in the  
 previous equations are with respect to  $q(\mathbf{y}^*|\mathbf{x}^*)$ . In  
 our BERT model, dropout is already being used.  
 Thus, by keeping dropout active during testing and  
 doing  $T$  forward passes per sample, we can obtain  
 Monte Carlo samples of the predictive distribution  
 and determine the uncertainty of all the validation  
 predictions, by for example using entropy as our  
 measure.

In practice, a system could for example send the top  
 $X\%$  highest entropy predictions to human review.  
 Depending on what one chooses  $X$  to be, this would  
 significantly increase the overall accuracy of the  
 system, while still saving a lot of time for employees.  
 This would reduce the risk of completely wrong  
 predictions on ambiguous samples.

## 217 2.7 Task 2g)

We chose to use Monte Carlo Dropout for uncer-  
 tainty modeling, because it is easy to implement  
 due to the fact that Dropout is already being used  
 during training. This lets us re-enable dropout dur-  
 ing inference to get the uncertainty estimates. Our  
 approach is doing  $T = 20$  forward passes with differ-  
 ent dropout masks and computing the probability of  
 each class. We then compute the predictive entropy,

$$242 \quad H[y|\mathbf{x}, \mathcal{D}] = - \sum_c p_c \log p_c \quad (23)$$

243 where  $p_c$  is the probability of class  $c$  which we get  
 244 from the estimated probability density:

$$245 \quad p(y|\mathbf{x}, \mathcal{D}) \approx \frac{1}{T} \sum_{t=1}^T p(y|\mathbf{x}, \omega_t) \quad (24)$$

246 where  $\omega_t$  is the sampled set of weights (dropout  
 247 mask). A higher entropy value implies that the  
 248 model is more uncertain about its classification of  
 249 the sample. Our implementation is provided in the  
 250 code below:

```

1 #Helper function to enable dropout
2 def enable_dropout(model):
3     if isinstance(model, torch.nn.Dropout):
4         model.train()
5
6 def mc_dropout_unc(model, dataloader, device, T=30):
7     model.to(device)
8     model.eval()
9     #Enabling dropout during inference
10    model.apply(enable_dropout)
11
12    sample_entropy = []
13
14    with torch.no_grad():
15        for batch in tqdm(dataloader):
16            input_ids = batch["input_ids"].to(device)
17            attention_mask = batch["attention_mask"].to(device)
18            preds_T = []
19
20            #Doing T forward passes through network
21            for i in range(0, T):
22                output = model(input_ids, attention_mask)
23                pred = F.softmax(output.logits, dim=-1)
24                preds_T.append(pred.unsqueeze(0))
25
26            #Converting to tensor
27            preds_T = torch.cat(preds_T, dim=0)
28            mean_preds = preds_T.mean(dim=0)
29            #Computing per sample predictive entropy in batch
30            entropy = -(mean_preds*(mean_preds+1e-12).log()).sum(dim=-1)
31            sample_entropy.append(entropy.cpu())
32
33            #Removing batch dimension
34            sample_entropy = torch.cat(sample_entropy, dim=0)
35            return sample_entropy

```

251

We also need to select some criteria for which samples are "uncertain enough" to be flagged for human analysis. We can do this by setting some cut-off value for the entropy or flagging a percentile of the samples with the highest entropy. Below we show the 4 validation samples with the highest uncertainty, along with their entropy, true label and predicted label:

```

1 ollice department suspends use of pepper pellet gun At least one big-city
  ↳ police department has suspended use of pepper-spray pellet guns blamed
  ↳ for the death of a 21-year-old college student
2 Entropy: 1.3659993410110474
3 True Label: 1
4 Predicted Label: 1
5 -----
6 Does a NY Supreme Court judge say that the Central Park lawn is more
  ↳ important than free speech? Here is the coverage in The New York Times
  ↳ and the New York Law Journal about the denial of United for Peace and
  ↳ Justice's Sunday rally in Central Park (United for Peace and Justice v.
  ↳ Bloomberg, 111893/04), and here is what NY Supreme Court Justice
  ↳ Silberman said: ... the evidence established that the
  ↳ department's determination was based on entirely content-neutral
  ↳ factors, to wit: that the Great Lawn2 was not an appropriate venue for
  ↳ a demonstration of this magnitude. ... The Parks Department
  ↳ appropriately applied content-neutral regulations while leaving
  ↳ plaintiff with a reasonable alternate site.
7 Entropy: 1.3181042671203613
8 True Label: 3
9 Predicted Label: 2
10 -----
11 Teen accused of making threats about school, possessing weapons A
  ↳ 17-year-old charged with making terrorist threats against his high
  ↳ school, fellow students and a police officer said at his arraignment
  ↳ that he would do anything quot;to help the community quot; in the
  ↳ wake of his arrest.
12 Entropy: 1.1795392036437988
13 True Label: 3
14 Predicted Label: 3
15 -----
16 Why do Workouts Work? Most machines don't improve with use. Old pickup
  ↳ trucks don't gradually become Ferraris just by driving them fast, and
  ↳ a pocket calculator won't change into a supercomputer by crunching
  ↳ lots of numbers. The human body is different...
17 Entropy: 1.1735734939575195
18 True Label: 3
19 Predicted Label: 1

```

260

The first sample doesn't really fit into any of the classes, so it makes sense that the model is uncertain (although it gets it right). For the second sample, the uncertainty likely comes from the fact that the text is poorly structured with random numbers and symbols placed all around the place. The same applies to the third sample. The fourth sample is likely predicted as "Sports" due to the word "Workout" being used, and the model is uncertain due to the title being a bit non-sensical.

## 2.8 Task 2h)

For the Kaggle competition we try a couple of different options for our model. The resulting accuracies of all these approaches are shown in table 1, where the first entry is simply the pre-trained model from task 2c which gave an accuracy of 94.03%.

Next we tried using a RoBERTa (Robustly optimized BERT pretraining Approach) model [11], which is a slightly larger model than the regular BERT model and it is also pre-trained for longer and in a *smarter* way. RoBERTa models can match or exceed the performance of regular BERT models [11]. There are 4 key differences in the pre-training process of RoBERTa models: (1) dynamic masking instead of static masking, (2) larger mini-batches, (3) larger byte-pair encoding (larger vocabularies) and (4) no Next Sentence Prediction (NSP) loss during training. We used a batch size of 64, learning rate = 1e-5,

weight decay = 0.01, 5 epochs and AdamW as the optimizer.

Lastly, we tried XLNet-large [12] which is a significantly larger pre-built model than BERT and RoBERTa. We used the same hyperparameters as for RoBERTa, but only 3 epochs due to the long training time (~2 hours per epoch). The key difference between XLNet and the previous models we used is the pre-training method. XLNet uses a generalized autoregressive method that maximizes the expected likelihood over all possible orders that the masked tokens are predicted in the autoregressive decomposition [12].

As we can see in table 1, RoBERTa clearly performed the best both on the test and validation set. The models were imported using the following code:

Model	Val. Acc.	Test Acc.
PT-BERT (Task 2c)	94.70%	94.04%
RoBERTa	94.95%	94.96%
XLNet-Large	94.75%	94.64%

Table 1. Validation and test accuracy of the used models.

```

1 from transformers import XLNetTokenizer, XLNetForSequenceClassification
2 tokenizer = XLNetTokenizer.from_pretrained('xlnet-large-cased')
3 model =
  ↳ XLNetForSequenceClassification.from_pretrained('xlnet-large-cased',
  ↳ num_labels=4)

```

```

1 from transformers import RobertaTokenizer,
  ↳ RobertaForSequenceClassification
2 tokenizer = RobertaTokenizer.from_pretrained("roberta-base")
3 model = RobertaForSequenceClassification.from_pretrained("roberta-base",
  ↳ num_labels=4)

```

The training loops and data pre-processing were identical to the ones shown in section 2.2.

## 3 Problem 3

### 3.1 Task 3a)

ResNet50 [13] is a so called *Residual Network* that implements the idea of *skip connections*. The benefit of skip connections is that optimization of deep networks becomes more stable, since one avoids vanishing/exploding gradients. This is especially important when we are training the network from scratch (instead of finetuning).

Our task is to classify CT images (512x512 pixels) of a liver into 3 classes (no liver, liver present, tumor present), which requires the model to capture relatively high-level context from the image, like organ boundaries and tumor regions. ResNet50 is a convolutional neural network (CNN), and is therefore well suited when it comes to extracting local textures, patterns and shapes. A medium-sized model like ResNet50 should be capable enough to capture patterns of such complexity. ResNet50 is also widely

328 used and is a well-tested architecture, which makes  
329 it a natural starting point for our task.

### 330 3.2 Task 3b)

331 We start by setting up the pre-built ResNet50 model  
332 with randomly initialized weights:

```
1 #Randomly initialized ResNet50
2 model = resnet50(weights=None)
3 #Changing number of output classes to 3
4 n_classes = 3
5 model.fc = nn.Linear(model.fc.in_features, n_classes)
```

333  
334 The only change we have to make to the archi-  
335 tecture is the number of output classes in the final  
336 linear classifier. Note that ResNet50 expects an in-  
337 put with dimensions  $(N, C = 3, H, W)$  where  $C$  is  
338 the number of channels, thus we have to add this  
339 extra dimension to our dataset. Another thing to  
340 note is that the validation data has pixelwise labels,  
341 which we have to convert to a single label per image.  
342 We do this with the following function:

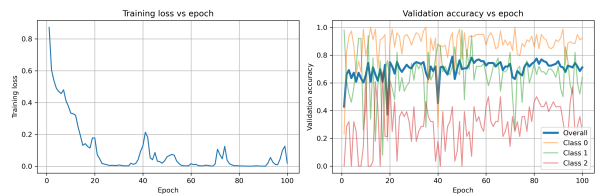
```
1 def pixelwise_to_class(img_dir, label_dir):
2     X = []
3     y = []
4     img_files = sorted(img_dir.glob("*.npy"))
5     label_files = sorted(label_dir.glob("*.npy"))
6
7     for img_path, label_path in zip(img_files, label_files):
8         img = np.load(img_path)
9         pixelwise = np.load(label_path)
10        X.append(img)
11        #If any pixel labelled 2 => class 2
12        if np.any(pixelwise == 2):
13            y.append(2)
14        #If no pixel has 2, but 1 exists => class 1
15        elif np.any(pixelwise == 1):
16            y.append(1)
17        else:
18            y.append(0)
19        #Converting X and y to tensors
20        X = np.stack(X)
21        y = np.array(y)
22        X = torch.from_numpy(X).float()
23        y = torch.from_numpy(y).long()
24        X = X.squeeze(1).repeat(1,3,1,1)
25    return X, y
```

343  
344 For the training we use the following hyperpa-  
345 rameters: **batch\_size=32**, **learning rate=1e-4**  
346 and **100 epochs**. As our optimizer we use AdamW  
347 [6], because the decoupled weight decay makes it  
348 ideal for training deep networks (especially from  
349 scratch). We use cross-entropy loss [14] as our loss  
350 function/criterion, which is a standard choice for  
351 single-label classification problems. The training  
352 loop is given in the following code:

```
1 for epoch in tqdm(range(n_epochs)):
2     model.train()
3     curr_loss = 0
4     samples = 0
5
6     #Training
7     for batch_x, batch_y in train_loader:
8         batch_x = batch_x.to(device)
9         batch_y = batch_y.to(device)
10        optimizer.zero_grad()
11        #forward pass
12        logits = model(batch_x)
13        #Using Cross-entropy loss as loss function
14        loss = loss_func(logits, batch_y)
15        #Backpass + backprop
16        loss.backward()
17        optimizer.step()
18        #Multiplying by size of current batch and dividing by total set
19        curr_loss += loss.item() * batch_x.size(0)
20        samples += batch_x.size(0)
21    training_loss.append(curr_loss/samples)
22
23    #Validation
24    model.eval()
25    total_samples = 0
26    total_correct = 0
27    correct = np.array([0,0,0])
28    n_class = np.array([0,0,0])
29
30    with torch.no_grad():
```

```
32     for batch_x, batch_y in val_loader:
33         batch_x = batch_x.to(device)
34         batch_y = batch_y.to(device)
35         logits = model(batch_x)
36         preds = torch.argmax(logits, dim=1)
37         #Counting number of samples and correcting preds
38         total_samples += batch_x.size(0)
39         total_correct += (preds == batch_y).sum().item()
40         #Counting classwise correct preds
41         for pred, label in zip(preds, batch_y):
42             if label == 0:
43                 n_class[0] += 1
44                 if pred == 0:
45                     correct[0] += 1
46             elif label == 1:
47                 n_class[1] += 1
48                 if pred == 1:
49                     correct[1] += 1
50             elif label == 2:
51                 n_class[2] += 1
52                 if pred == 2:
53                     correct[2] += 1
54     val_acc0.append(correct[0]/n_class[0])
55     val_acc1.append(correct[1]/n_class[1])
56     val_acc2.append(correct[2]/n_class[2])
57     val_acc_tot.append(total_correct/total_samples)
58     current_val_acc = val_acc_tot[-1]
59     if current_val_acc > best_val_acc:
60         best_val_acc = current_val_acc
61         best_epoch = epoch + 1
62         best_checkpoint = {
63             "epoch": n_epochs,
64             "model_state_dict": model.state_dict(),
65             "train_losses": training_loss,
66             "val_acc_tot": val_acc_tot,
67             "val_acc0": val_acc0,
68             "val_acc1": val_acc1,
69             "val_acc2": val_acc2}
70     torch.save(best_checkpoint, "best_problem3_checkpoint.pth")
```

353 Since 100 epochs is potentially too many, we use  
354 *checkpointing* in order to save the model state with  
355 the highest validation accuracy. In fact, as we can  
356 see in figure 6 the total validation accuracy peaks  
357 at epoch 46 with **78.85%**. The class-wise accuracies  
358 for this model state were 96% (class 0), 72% (class 1)  
359 and 43% (class 2). This big spread in performance  
360 across classes likely comes from the fact that there  
361 is a class imbalance in the training (and validation)  
362 set. 50% of the samples are class 0, while 33.8%  
363 and 16.2% are class 1 and 2 respectively. Thus, the  
364 model will do reasonably well if it just guesses that  
365 most samples are class 0.



366 **Figure 6.** Training results from scratch. Left: Training  
367 loss as a function of epochs. Right: Total validation  
368 accuracy (blue), class 0 (yellow), class 1 (green), class 2  
369 (red).

### 370 3.3 Task 3c)

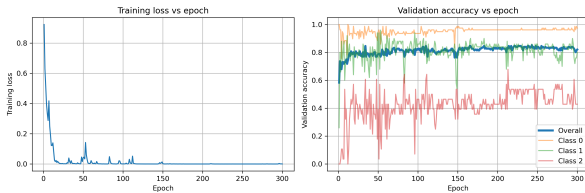
371 We load the pre-trained DINO ResNet50 [15], which  
372 outputs a 2048 dimensional feature vector. Thus,  
373 we need to add an additional linear classification  
374 layer with  $(in\_dim, out\_dim) = (2048, 3)$ :

```
1 dino = torch.hub.load('facebookresearch/dino:main', 'dino_resnet50',
2                       pretrained=True)
3 model = nn.Sequential(dino, nn.Linear(2048, 3))
```

374 Initially we run the training loop with AdamW as  
375 our optimizer. This gave a peak validation accuracy  
376 of 78.18%, which is worse than the model we trained  
377 from scratch. In the documentation [15] they use

378 SGD [16] (with weight decay) as the optimizer, so  
 379 we will try that instead. We also copy their hyper-  
 380 parameters: learning rate = 0.03, weight decay =  
 381 1e-4 and 300 epochs. The results of the training are  
 382 shown in figure 7. We can see that the class 0 and to-  
 383 tal validation accuracy pretty much converge, while  
 384 the class 1 and 2 accuracies never quite stabilize dur-  
 385 ing the training. This is still an improvement from  
 386 the model we trained from scratch, although the  
 387 comparison is not really fair since we have now done  
 388 more epochs and used a different optimizer. The  
 389 total validation accuracy peak was **85.26%** (epoch  
 390 243), which is a significant improvement. The class  
 391 specific accuracies were 96% (+0%), 86% (+14%)  
 392 and 54% (+11%) respectively.

393 Ideally we would run the first model with the same  
 394 optimizer and hyperparameters in order to get a  
 more fair comparison.



**Figure 7.** Training results for pre-trained ResNet50. Left: Training loss as a function of epochs. Right: Total validation accuracy (blue), class 0 (yellow), class 1 (green), class 2 (red).

395

### 396 3.4 Task 3d)

397 The purpose of XAI is to reinforce our trust in a  
 398 given AI model. Most state-of-the-art networks have  
 399 many layers and millions/billions of parameters and  
 400 are essentially "black boxes", meaning we as humans  
 401 cannot comprehend or understand the reasoning of a  
 402 model. This is where explainable AI comes in, which  
 403 is a field of research that aims to create ways for an  
 404 AI model to *explain why* it does something. This  
 405 can be done by for example highlighting regions of  
 406 interest in an image or generating a text-explanation.  
 407 In our case, we are interested in knowing why it clas-  
 408 sifies a certain CT image as having a tumor and want  
 409 it to highlight the ROI. In general, interpretability  
 410 is a very important factor in medical AI settings,  
 411 since mistakes done by the model can be fatal.

412 Gradient-weighted Class Activation Mapping (Grad-  
 413 CAM) [17] is an XAI technique that uses the gradi-  
 414 ents with respect to any target concept, for instance  
 415 "tumor", to highlight influential regions for the clas-  
 416 sification in a given image (can also be used for text,  
 417 etc.). The original paper [17] tests Grad-CAM on  
 418 ResNet-based architectures with good success, mak-  
 419 ing it a natural choice for us.

420 For a chosen class  $c$  we take the gradient of the class

score,  $y^c$  (before softmax is applied), with respect to  
 the feature map activation of the final convolutional  
 layer  $A^k$ . The gradients are then global average  
 pooled:

$$\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}^k} \quad (25)$$

where  $Z$  is a normalization constant and  $(i, j)$  de-  
 notes a pixel. The heatmap is then constructed by  
 taking a linear combination of these "weights" and  
 the feature maps:

$$L_{Grad-CAM}^c = \sum_k \alpha_k^c A^k. \quad (26)$$

This matrix will highlight both positive and neg-  
 ative influences of the given classification, but we  
 are interested in only the positive ones. Thus, the  
 last step is passing  $L_{Grad-CAM}^c$  through a ReLU  
 function.

### 3.5 Task 3e)

We implement GradCAM using the PyTorch Grad-  
 CAM library:

```
1 backbone = model[0] #extracting ResNet50 part of our model
2 layer = backbone.layer4[-1]
3 grad_cam = GradCAM(model=model, target_layers = [layer])
```

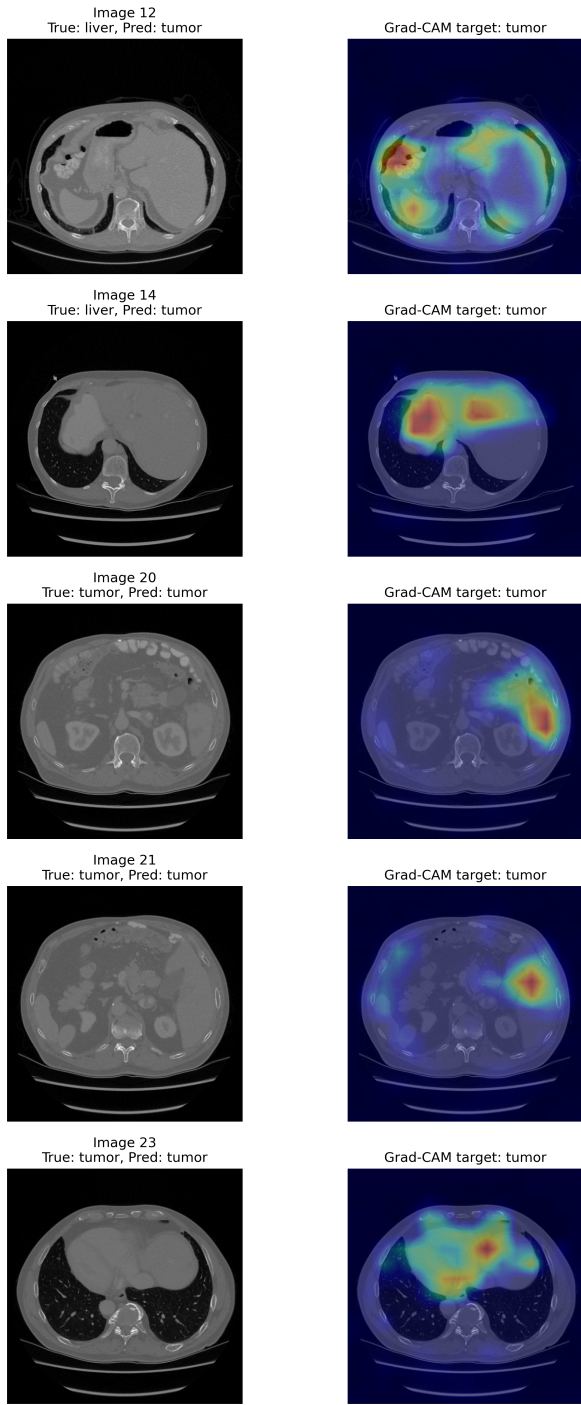
We are using the gradients from the final con-  
 volutional layer to create the heatmap (eq. 26),  
 which is layer 4 in this case. These gradients will be  
 used in equation 25 to determine the contribution  
 (weight) of a particular region in the image for the  
 given classification. We apply GradCAM to 5 val-  
 idation images that were labeled as tumors, using  
 the following code:

```
1 for i in range(0,5):
2     idx = val_idxr[i]
3     img_path = val_file_names[idx]
4     img = np.load(img_path)
5     lbl = val_y[idx]
6     #Selecting target class for GradCAM as the predicted class
7     target = val_preds[idx]
8     targets = [ClassifierOutputTarget(target)]
9     #Applying grayscale GradCAM to tensor representation of image
10    X = torch.tensor(img,
11    ↪ dtype=torch.float32).unsqueeze(0).unsqueeze(0).to(device)
12    gs_cam = grad_cam(input_tensor = X, targets = targets)
13    # ... plotting
```

The resulting images are shown in figure 8. The  
 first two examples were wrong predictions, with  
 "liver" being the true label, while the last three ex-  
 amples are correctly predicted as "tumor". To an  
 untrained eye, the "blobs" that the model activates  
 on seem rather similar in all of the examples. With-  
 out medical expertise it is hard to tell what exactly  
 the model is doing wrong in these examples.

### 3.6 Task 3f)

In order to turn our GradCAM output into a mask  
 for pixelwise classification, we will start by normal-  
 izing the output. We will then choose a threshold  
 in  $[0, 1]$  that will determine if we classify a pixel as  
 tumor/liver or background. In order to choose this



**Figure 8.** GradCAM applied to 5 validation images that were classified as tumors by the DINO-ResNet50 model. Left: Original image. Right: GradCAM applied.

463 threshold, we will try 8 different ones for each class,  
464 and choose the ones that give the best IoU score.  
465 The implementation of the GradCAM pixelmask  
466 and IoU-score computation:

```

1 def gradcam_pixelmask(gradcam, threshold):
2     #Normalizing the gradcam heatmap
3     gradcam = (gradcam - gradcam.min()) / (gradcam.max() - gradcam.min() +
4     ↪ 1e-9)
5     #gradcam value above threshold -> mask=1, else 0
6     mask = np.where(gradcam >= threshold, 1, 0)
7     #Smoothing the mask
8     #Removing noise outside main object, 3x3 kernel

```

467

```

8     mask = binary_opening(mask, structure=np.ones((5,5)))
9     #Filling small holes inside object, 3x3 kernel
10    mask = binary_closing(mask, structure=np.ones((5,5)))
11    return mask
12
13 def iou_score(pred, true):
14     #Turning masks into boolean arrays so we can use np.logical funcs
15     pred_mask = pred.astype(bool)
16     true_mask = true.astype(bool)
17     #Computing intersection and union
18     inter = np.logical_and(pred_mask, true_mask).sum()
19     union = np.logical_or(pred_mask, true_mask).sum()
20     iou = inter / (union + 1e-9)
21     return iou

```

468

The thresholds were determined using the following code:

469

470

```

1 thresholds = np.linspace(0.5, 0.99, 20)
2 best_iou_liver = 0
3 best_iou_tumor = 0
4 best_threshold_tumor = 0
5 best_threshold_liver = 0
6 for threshold in tqdm(thresholds):
7     iou_tumor = []
8     #Looping over validation images with tumor
9     for i in range(len(val_file_names)):
10        if val_y[i] == 2:
11            #Loading image and pixelwise labelling
12            img = np.load(val_file_names[i])
13            true_labels = np.load(val_label_file_names[i])
14            #Converting image to tensor with right dimensions
15            X = torch.tensor(img,
16            ↪ dtype=torch.float32).unsqueeze(0).unsqueeze(0).to(device).repeat(1,3,1,1)
17            #Defining tumor as target
18            target_tumor = ClassifierOutputTarget(2)
19            #GradCAM output for tumor class for sample
20            gradcam_tumor = gradcam(input_tensor=X, targets =
21            ↪ [target_tumor])[0]
22            #Computing true and prediction mask
23            pred_mask_tumor = gradcam_pixelmask(gradcam_tumor, threshold)
24            true_mask_tumor = np.where(true_labels == 2, 1, 0)
25            #Computing IoU
26            iou_tumor.append(iou_score(pred_mask_tumor, true_mask_tumor))
27        iou_tumor = np.array(iou_tumor)
28        mean_iou_tumor = np.mean(iou_tumor)
29        if mean_iou_tumor >= best_iou_tumor:
30            best_iou_tumor = mean_iou_tumor
31            best_threshold_tumor = threshold
32 #...Same for liver class

```

471

As we can see in table 2, the results are not very good with an IoU score of 13.96% and 11.69% for the tumor and liver classes respectively. The reason for this is likely that the representations of the image have been down-sampled too much in the final layer in order to get accurate segmentations. Tumors in CT images are generally small and fuzzy, while Grad-CAM produces low-resolution blobs that are meant to be used as a visual aid. As discussed in [18], GradCAM is meant for visualization and not for image segmentation tasks. For image segmentation there are specialized models that will perform way better, like for example U-Net (specialized for medical images) and SAM [19].

472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

Across all the validation samples the highest IoU achieved was 60.29%, while the worst was 0%. Here we used a threshold of 0.76, which gave the highest mean IoU. These examples are shown in figure 9. In the good example we see that the true position of the tumor aligns very well with the most dense region on the GradCAM heatmap. What is dragging the IoU score down for this sample are the random weak activations in other regions of the image. The example of bad overlap is a bit brutal, because the tumor is extremely small and the highlighted region is just below the true position.

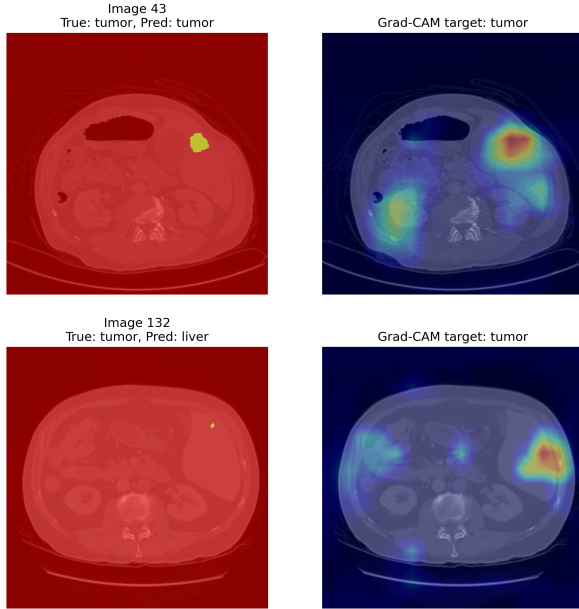
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497

Our conclusion is that GradCAM should only be used as a decision-support and visualization tool in clinical practice, and NOT be used as an automatic replacement of a radiologist.

498  
499  
500  
501

Method (Class)	Mean IoU Score
GradCAM (Tumor)	0.1396
GradCAM (Liver)	0.1169

**Table 2.** Mean IoU scores.



**Figure 9.** Left: True label, tumor highlighted in yellow. Right: Grad-CAM heatmap of same image. Top example shows the best overlap, while the bottom example shows 0 overlap.

### 3.7 Task 3g)

Initially we pass the test data through our model from task 3c. This resulted in an accuracy of 83.116%.

In order to improve the performance, we will use a pre-built model from Huggingface called ConvNeXT [20] which is based on the ConvNet architecture [21]. In [21] (table 10) they show the steps they add to their baseline model, which is ResNet50, in order to end up at the ConvNeXT architecture. They also show how much each of these extra steps improve accuracy on the ImageNET-1k dataset. Some examples of these improvements are: Using GELU [22] instead of ReLU, using fewer activations and norms, depthwise convolutions and a larger kernel size on the filters. Thus, we can assume that using ConvNeXT will be an improvement for our image classification task.

Due to the class imbalance in the training set we will also implement class-weighting into our Cross-Entropy loss function. To try to prevent overfitting and model overconfidence we will also add label smoothing (=0.1). We use the following hyperparameters: learning rate = 1e-4, 50 epochs (with checkpointing), weight decay = 0.01 and batch size

8. We used AdamW [6] as the optimizer because it was used in [21]. The training results are shown in figure 10. The peak validation accuracy was 88.46% (epoch 3). We imported the model using the following code, and changed the final linear classification layer to fit our problem:

```

1 weights = ConvNeXt_Base_Weights.DEFAULT
2 model = convnext_base(weights=weights)
3 #Tuning the linear classifier so it has correct dimensions
4 in_features = model.classifier[2].in_features
5 model.classifier[2] = nn.Linear(in_features, 3)

```

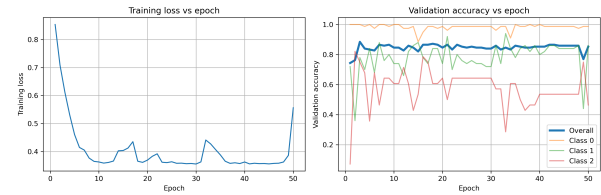
On the Kaggle test data ConvNeXT performed significantly better than our previous model, achieving an accuracy of **90.909%**. The classwise accuracies are shown in table 3.

The key step in our implementation of ConvNeXT is the data pre-processing. Everything else is pretty much the same as for the other model. ConvNeXT is pre-trained on ImageNET, so we have to pre-process the data in the same way as ImageNET. Here is how we do that on the training data:

```

1 #ConvNeXt takes in input with dims (N, C, H, W). Adding C dimension (3
  ↳ channels):
2 train_x = train_x.unsqueeze(1).repeat(1,3,1,1)
3 #Normalizing training set
4 train_x = train_x / 255
5 #ImageNET mean and std RGB values (using same normalization as ImageNET)
6 mean = torch.tensor([0.485, 0.456, 0.406]).view(1, 3, 1, 1)
7 std = torch.tensor([0.229, 0.224, 0.225]).view(1, 3, 1, 1)
8 train_x.sub_(mean)
9 train_x.div_(std)

```



**Figure 10.** Training curves for ConvNeXT model.

Class	Validation Accuracy (%)
Total	88.46%
Background (0)	100%
Liver (1)	78%
Tumor (2)	75%

**Table 3.** Validation accuracy of ConvNeXT model.

## 4 AI Declaration

For this submission Perplexity AI [23] was used as aid. The AI was used to generate code for plotting, like for example figures 3 and 4 (only the tedious matplotlib blocks of code). The AI was also used to generate BibTeX entries for various sources and converting handwritten formulas to LaTeX code.

## References

- [1] S. Kullback and R. A. Leibler. “On Information and Sufficiency”. In: *The Annals of Mathematical Statistics* 22.1 (1951). DOI: [10.1214/aoms/1177729694](https://doi.org/10.1214/aoms/1177729694).
- [2] Wikipedia contributors. *Normal distribution*. [https://en.wikipedia.org/wiki/Normal\\_distribution](https://en.wikipedia.org/wiki/Normal_distribution). Accessed: 2026-04-11. 2025.
- [3] Wikipedia contributors. *Evidence lower bound*. [https://en.wikipedia.org/wiki/Evidence\\_lower\\_bound](https://en.wikipedia.org/wiki/Evidence_lower_bound). Accessed: 2026-04-11. 2026.
- [4] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Ed. by J. Burstein, C. Doran, and T. Solorio. Minneapolis, Minnesota: Association for Computational Linguistics, June 2019, pp. 4171–4186. DOI: [10.18653/v1/N19-1423](https://doi.org/10.18653/v1/N19-1423). URL: <https://aclanthology.org/N19-1423/>.
- [5] Hugging Face. *BERT — Hugging Face Transformers Documentation*. Accessed: 2026-04-14. 2024. URL: [https://huggingface.co/docs/transformers/model\\_doc/bert](https://huggingface.co/docs/transformers/model_doc/bert).
- [6] Y. Hao, Z. Tang, Y. Tian, Y. Zhang, and Z. Zhou. *AdamW*. Cornell University Computational Optimization Open Textbook - Optimization Wiki. ChemE 6800 Fall 2024. 2024. URL: <https://optimization.cbe.cornell.edu/index.php?title=AdamW> (visited on 04/14/2026).
- [7] L. McInnes, J. Healy, and J. Melville. “UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction”. In: *arXiv preprint arXiv:1802.03426* (2018). DOI: [10.48550/arXiv.1802.03426](https://doi.org/10.48550/arXiv.1802.03426). arXiv: [1802.03426](https://arxiv.org/abs/1802.03426) [stat.ML].
- [8] A. Maćkiewicz and W. Ratajczak. “Principal components analysis (PCA)”. In: *Computers & Geosciences* 19.3 (1993), pp. 303–342. ISSN: 0098-3004. DOI: [10.1016/0098-3004\(93\)90090-R](https://doi.org/10.1016/0098-3004(93)90090-R). URL: <https://www.sciencedirect.com/science/article/pii/S009830049390090R>.
- [9] H. Sajjad, F. Dalvi, N. Durrani, and P. Nakov. “On the effect of dropping layers of pre-trained transformer models”. In: *Computer Speech & Language* 77 (2023), p. 101429. ISSN: 0885-2308. DOI: [10.1016/j.csl.2022.101429](https://doi.org/10.1016/j.csl.2022.101429). URL: <https://www.sciencedirect.com/science/article/pii/S0885230822000596>.
- [10] Y. Gal and Z. Ghahramani. “Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning”. In: *Proceedings of The 33rd International Conference on Machine Learning*. Ed. by M. F. Balcan and K. Q. Weinberger. Vol. 48. Proceedings of Machine Learning Research. New York, New York, USA: PMLR, 20–22 Jun 2016, pp. 1050–1059. URL: <https://proceedings.mlr.press/v48/gal16.html>.
- [11] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov. “RoBERTa: A Robustly Optimized BERT Pretraining Approach”. In: *arXiv preprint arXiv:1907.11692* (2019).
- [12] Z. Yang, Z. Dai, Y. Yang, J. G. Carbonell, R. Salakhutdinov, and Q. V. Le. “XLNet: Generalized Autoregressive Pretraining for Language Understanding”. In: *CoRR* abs/1906.08237 (2019). arXiv: [1906.08237](https://arxiv.org/abs/1906.08237). URL: <http://arxiv.org/abs/1906.08237>.
- [13] K. He, X. Zhang, S. Ren, and J. Sun. “Deep residual learning for image recognition”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016, pp. 770–778.
- [14] PyTorch Contributors. *CrossEntropyLoss — PyTorch 2.11 Documentation*. Accessed: 2026-04-23. PyTorch. 2024. URL: <https://docs.pytorch.org/docs/stable/generated/torch.nn.CrossEntropyLoss.html>.
- [15] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. “Emerging Properties in Self-Supervised Vision Transformers”. In: *Proceedings of the International Conference on Computer Vision (ICCV)*. 2021.
- [16] PyTorch Team. *SGD — PyTorch 2.11 Documentation*. Accessed: 2026-04-23. PyTorch. 2024. URL: <https://docs.pytorch.org/docs/stable/generated/torch.optim.SGD.html>.
- [17] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. “Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization”. In: *arXiv preprint arXiv:1610.02391* (2016).
- [18] T. Rheude, A. Wirtz, A. Kuijper, and S. Weisarg. “Leveraging CAM Algorithms for Explaining Medical Semantic Segmentation”. In: *arXiv preprint arXiv:2409.20287* (2024). URL: <https://arxiv.org/abs/2409.20287>.

- 660 [19] W. Yao, J. Bai, W. Liao, Y. Chen, M. Liu,  
661 and Y. Xie. “From CNN to Transformer: A  
662 Review of Medical Image Segmentation Mod-  
663 els”. In: *Journal of Imaging Informatics in*  
664 *Medicine* 37.4 (2024), pp. 1529–1547. DOI: [10.](https://doi.org/10.1007/s10278-024-00981-7)  
665 [1007/s10278-024-00981-7](https://doi.org/10.1007/s10278-024-00981-7). URL: [https :](https://pmc.ncbi.nlm.nih.gov/articles/PMC11300773/)  
666 [// pmc . ncbi . nlm . nih . gov / articles /](https://pmc.ncbi.nlm.nih.gov/articles/PMC11300773/)  
667 [PMC11300773/](https://pmc.ncbi.nlm.nih.gov/articles/PMC11300773/).
- 668 [20] Hugging Face. *ConvNeXT · Hugging Face*  
669 *Transformers Documentation*. Model released  
670 2022-01-10, added to Transformers 2022-02-07.  
671 2022. URL: [https://huggingface.co/docs/](https://huggingface.co/docs/transformers/model_doc/convnext)  
672 [transformers/model\\_doc/convnext](https://huggingface.co/docs/transformers/model_doc/convnext).
- 673 [21] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T.  
674 Darrell, and S. Xie. “A ConvNet for the 2020s”.  
675 In: *Proceedings of the IEEE/CVF Conference*  
676 *on Computer Vision and Pattern Recognition*.  
677 2022.
- 678 [22] D. Hendrycks and K. Gimpel. “Gaussian Er-  
679 ror Linear Units (GELUs)”. In: *arXiv preprint*  
680 *arXiv:1606.08415* (2016). URL: [https : / /](https://arxiv.org/abs/1606.08415)  
681 [arxiv.org/abs/1606.08415](https://arxiv.org/abs/1606.08415).
- 682 [23] Perplexity. *Perplexity AI Assistant (powered*  
683 *by GPT-5.1)*. <https://www.perplexity.ai/>.  
684 Accessed: 2026-04-21. 2026.