

1. Design your own architecture and justify its components (width, depth, convolutions, activations and any other operations)

```

MyNetwork.py x
hw_cnn > MyNetwork.py > % MyNetwork > __init__
40 # begin: [28 x 28 x 3]
41 if self.architecture == 'default':
42     self.conv = torch.nn.Sequential(
43         torch.nn.Conv2d(in_channels=self.nInputChannels, kernel_size=3, out_channels=32, padding='valid'),
44         # conv out: [26 x 26 x 32]
45         torch.nn.ReLU(),
46         torch.nn.BatchNorm2d(32),
47         torch.nn.Dropout(0.2),
48         torch.nn.Conv2d(in_channels=32, kernel_size=3, out_channels=64, padding='valid'),
49         # conv out: [24 x 24 x 64]
50         torch.nn.ReLU(),
51         torch.nn.BatchNorm2d(64),
52         torch.nn.MaxPool2d(kernel_size=3)
53     )
54     # max pool out: [8 x 8 x 64]
55     # 8*8*64 = 4096 (i.e., the flattened vector length)
56
57     ## Add your fully connected architecture
58     self.fc = torch.nn.Sequential(
59         torch.nn.Dropout(0.5),
60         torch.nn.Linear(4096, 100),
61         torch.nn.ReLU(),
62         torch.nn.Dropout(0.5),
63         torch.nn.Linear(100, self.nOutputClasses)
64     )
65
66 if self.architecture == "arch2":
67     # begin [28 x 28 x 3]
68     self.conv = torch.nn.Sequential(
69         torch.nn.Conv2d(self.nInputChannels, out_channels=32, kernel_size=(3,3)),
70         # conv out: [26 x 26 x 32]
71         torch.nn.ReLU(),
72         torch.nn.BatchNorm2d(32),
73         torch.nn.Dropout(0.2),
74         torch.nn.Conv2d(32, out_channels=32, kernel_size=(3,3)),
75         torch.nn.ReLU(),
76         torch.nn.BatchNorm2d(32),
77         # conv out: [24 x 26 x 32]
78         torch.nn.MaxPool2d(kernel_size=2),
79         torch.nn.Dropout(0.2),
80     )
81     self.fc = torch.nn.Sequential(
82         torch.nn.Linear(4608, 200),
83         torch.nn.ReLU(),
84         torch.nn.BatchNorm1d(200),
85         torch.nn.Dropout(0.5),
86         torch.nn.Linear(200, 100),
87         torch.nn.ReLU(),
88         torch.nn.BatchNorm1d(100),
89         torch.nn.Dropout(0.5),
90         torch.nn.Linear(100, self.nOutputClasses)
91     )
92

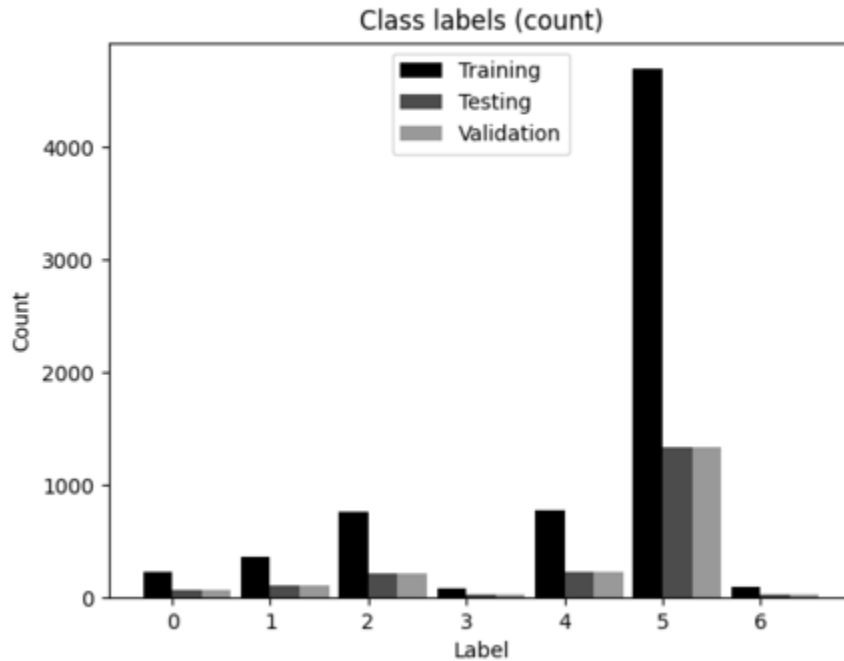
```

**Fig 1: Model architectures.** Architecture 1 (**Arch1**) is the top-half of code and architecture 2 (**Arch2**) is the bottom-half of code.

The **two architectures** share common themes of using **batch normalization** to control the scale of input data to a neuron across batches, and using **dropouts** to prevent overfitting. The **main differences** between the architectures is the **dimensions of the feature maps** and the **pooling kernel shape** used for dimensionality reduction at the end of the convolutional layer, and **no batch normalization in the fully-connected layer of Arch1**.

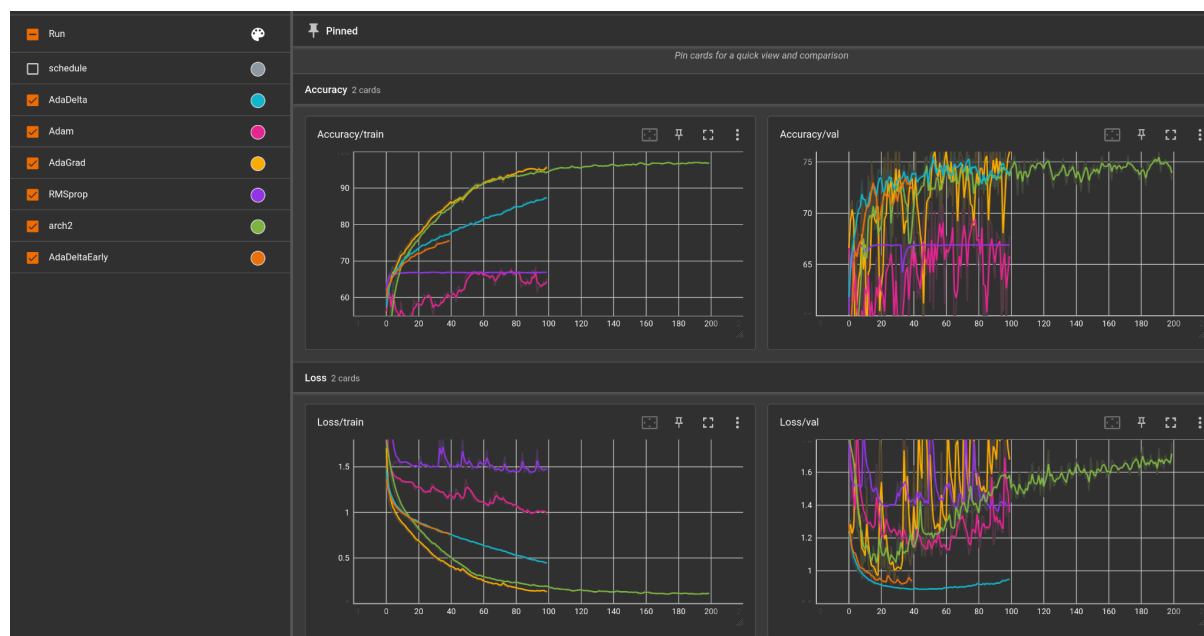
The choice of depth, channels, convolutions, and activations was chosen by following similar architectures found from similar examples in literature (e.g., [https://link.springer.com/chapter/10.1007/978-3-030-76773-0\\_15#Tab4](https://link.springer.com/chapter/10.1007/978-3-030-76773-0_15#Tab4)). However, model architectures can be tweaked extensively making it hard to conclude what the best approach is without any additional knowledge. Given more time, it would be worth training on randomized values of output features and tweaking the model architecture to further minimize loss on the validation set during training.

2. Justify your choices of loss function and optimizer
3. Provide the curves (preferably from Tensorboard) with the evolution of the training and validation losses and accuracies during training



**Fig 3: Class balance for MedMNIST dataset**

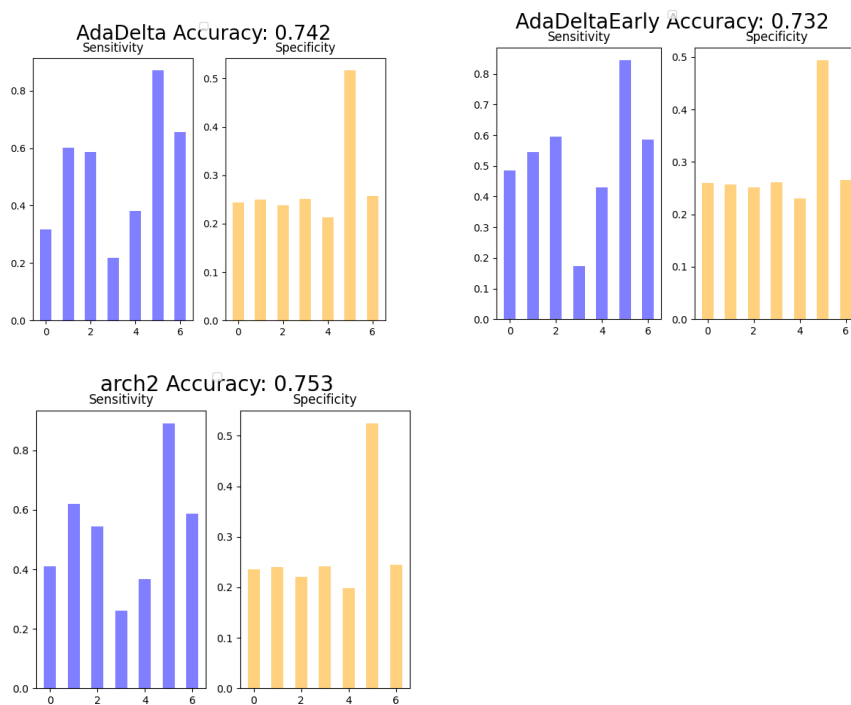
A **categorical cross entropy loss function** was used since this function can penalize multiple classes for a single label. **Class weights** were assigned as  $w_k = (1 - N_k/N)$  where  $N_k$  is the number of samples for a given  $k$ th class. Ultimately, this will weigh classes by the complement of their frequency in the dataset such that a highly represented class (e.g., class 5 [Fig 3]) will have a lesser effect on the loss function to prevent optimization of the highly represented class.



**Fig 4: Loss/accuracy for various optimization functions and architectures.** AdaDelta, Adam, AdaGrad, and RMSprop optimizers used during training with **Arch1**. The **Arch2** model used the AdaDelta optimizer.

The **AdaDelta** optimizer was chosen since it had the lowest loss in the validation set compared to other optimization functions.

4. Evaluate the performance on the test dataset by providing the overall accuracy of the model, and the individual sensitivity and specificity (and any other metrics that you may consider relevant) on every class. Interpret your results.



**Fig 5 model performance on test.** The **accuracy, sensitivity, and specificity** is reported for 3 models AdaDelta, AdaDelta with an early stop at 40 epochs, and Arch2 (from Fig1).

The most performant model in terms of accuracy was architecture 2, and the possible reasoning for the slight performance boost of Arch2 could be one of the following aspects that make it different from the other models: **batch normalization on the fully connected layer, different feature map size, and/or the increased number of epochs** (Fig 4).