



DEEPHEALTH

Deep Learning pipeline on histopathology images: detection of prostatic tumor

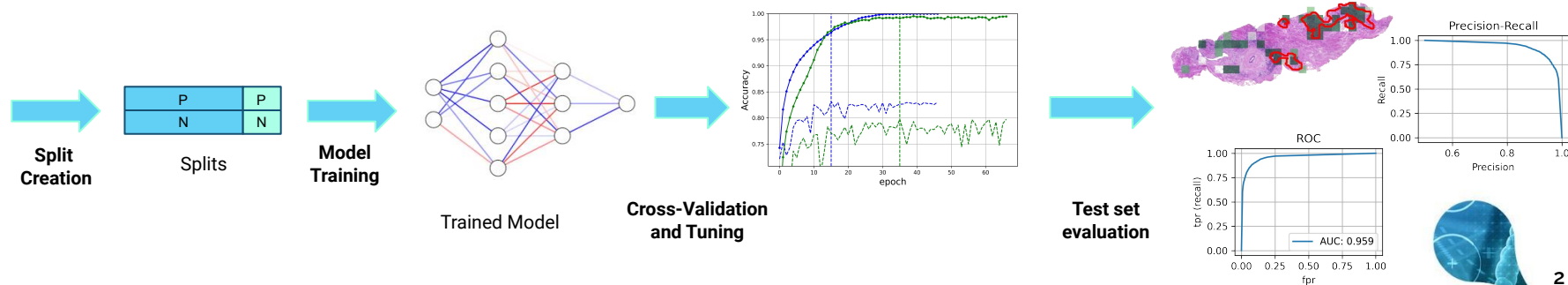
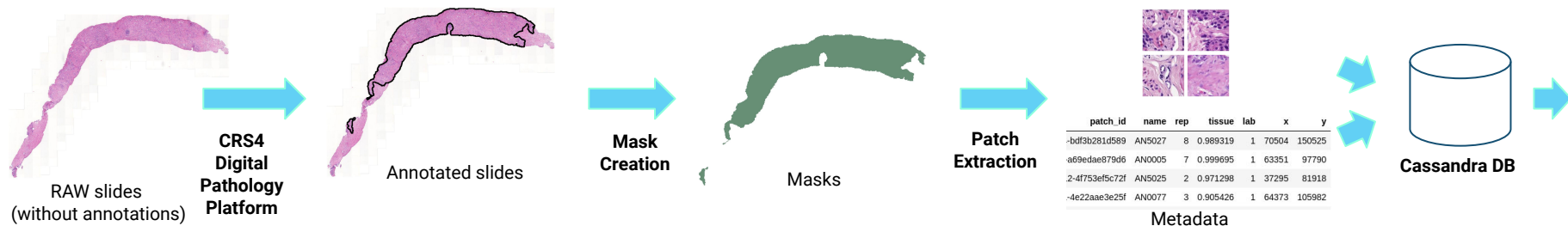
Francesco Versaci, **Giovanni Busonera**



The project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 825111.

DeepHealth Winter School 2022

Pipeline: Predictive model





Creation of the splits

P	P
N	N

training
split

validation
split

All **splits** are **balanced**, meaning that the number of tumor patches is roughly the same of the normal ones.

We used the **80%** of the whole dataset for **training data**. The remaining **20%** is used for the **validation set**. The **split_ratios** parameter is set to **[8, 2]**. (**47782** and **11558** patches respectively)

The **test-set** data is organized as a **single split** obtained from different Cassandra tables (**44416** patches)

```
ap = PlainTextAuthProvider(username='xxxx', password='xxxx')
cd = CassandraDataset(ap, ['cassandra-db'])

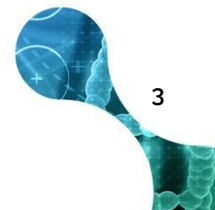
cd.init_listmanager(table= args.ids_table,
                    metatable= args.metatable,
                    id_col='patch_id',
                    num_classes=args.num_classes,
                    label_col=args.label,
                    grouping_cols=['sample_name']
)

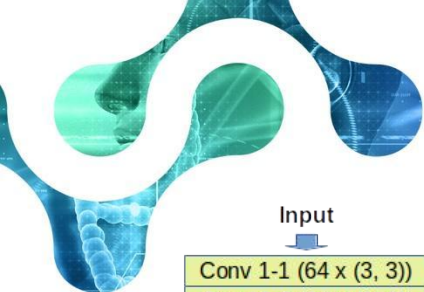
cd.read_rows_from_db()

min_class = np.min(cd._clm._stats.sum(axis=0))
data_size = min_class*cd._clm.num_classes

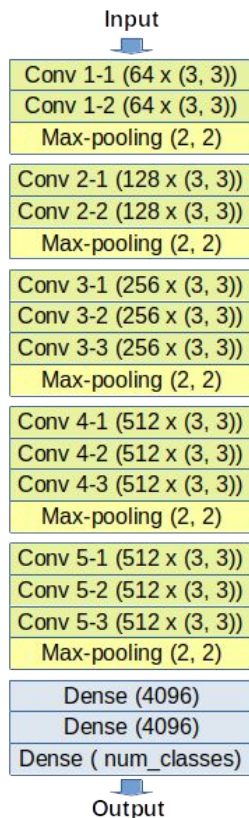
cd.init_datatable(table=args.table)
cd.split_setup(split_ratios=args.split_ratios, max_patches=data_size, augs=[])

cd.save_splits(args.out_fn)
```

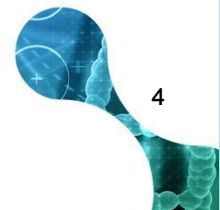


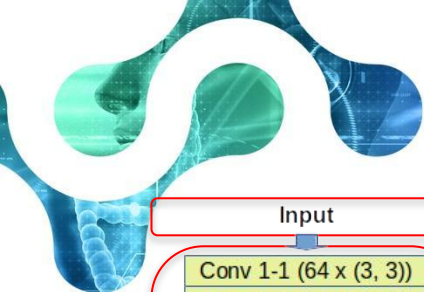


VGG16

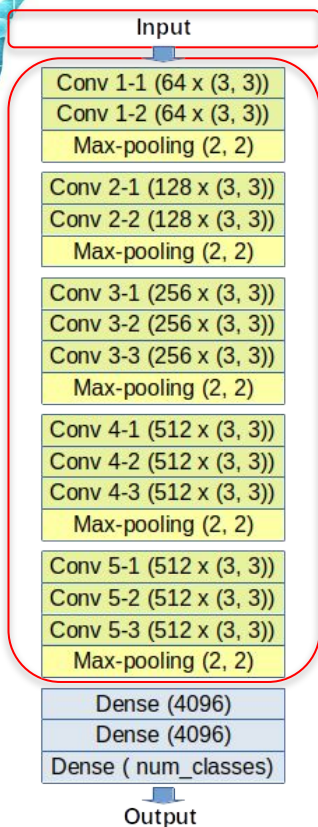


- **VGG16** is a **convolutional neural network** with a top **classifier** made of **dense layers**
- The number of **trainable parameters** is about **166M** (images **256x256** pixel instead of 224x224)
- **Pros:** Classical CNN architecture, **easy to implement** from scratch and although not currently the state of the art, it **was one of the best performing architectures** in ImageNet Large Scale Visual Recognition Challenge (ILSVRC) challenge 2014
- **Cons:** It is **slow to train** compared to more recent ANNs. The **size of VGG16 weights** is about **600MB** which potentially takes a lot of disk space (e.g. if you store the model after each epoch) and hence a **long time to synchronize** in a distributed training setting
- [EDDL provides a VGG16 neural network pretrained on the Imagenet dataset](#)
- We implemented a VGG16 by using PyEDDL to have more flexibility (e.g.: train from scratch with **different initialization methods** or **pretrained net**)



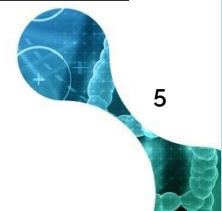


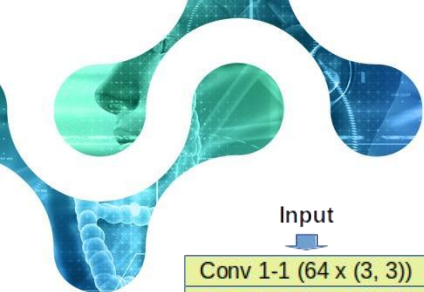
VGG16



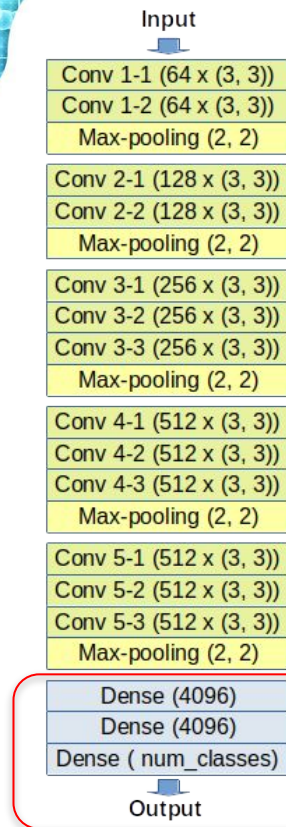
```
in_ = eddl.Input([3, patch_size[0], patch_size[1]])
```

```
def VGG16(in_, num_classes, seed=1234, init=eddl.HeNormal, l2_reg=None, dropout=None):
    x = in_
    x = eddl.ReLu(init(eddl.Conv(x, 64, [3, 3]), seed))
    x = eddl.MaxPool(eddl.ReLu(init(eddl.Conv(x, 64, [3, 3]), seed)), [2, 2], [2, 2])
    x = eddl.ReLu(init(eddl.Conv(x, 128, [3, 3]), seed))
    x = eddl.MaxPool(eddl.ReLu(init(eddl.Conv(x, 128, [3, 3]), seed)), [2, 2], [2, 2])
    x = eddl.ReLu(init(eddl.Conv(x, 256, [3, 3]), seed))
    x = eddl.ReLu(init(eddl.Conv(x, 256, [3, 3]), seed))
    x = eddl.MaxPool(eddl.ReLu(init(eddl.Conv(x, 256, [3, 3]), seed)), [2, 2], [2, 2])
    x = eddl.ReLu(init(eddl.Conv(x, 512, [3, 3]), seed))
    x = eddl.ReLu(init(eddl.Conv(x, 512, [3, 3]), seed))
    x = eddl.MaxPool(eddl.ReLu(init(eddl.Conv(x, 512, [3, 3]), seed)), [2, 2], [2, 2])
    x = eddl.ReLu(init(eddl.Conv(x, 512, [3, 3]), seed))
    x = eddl.ReLu(init(eddl.Conv(x, 512, [3, 3]), seed))
    x = eddl.MaxPool(eddl.ReLu(init(eddl.Conv(x, 512, [3, 3]), seed)), [2, 2], [2, 2])
```





VGG16



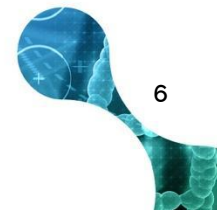
```
x = eddl.Reshape(x, [-1])
x = eddl.Dense(x, 4096)
if dropout:
    x = eddl.Dropout(x, dropout, iw=False)
if l2_reg:
    x = eddl.L2(x, l2_reg)
x = init(x, seed)
x = eddl.ReLU(x)
x = eddl.Dense(x, 4096)
if dropout:
    x = eddl.Dropout(x, dropout, iw=False)
if l2_reg:
    x = eddl.L2(x, l2_reg)
x = init(x, seed)
x = eddl.ReLU(x)
x = eddl.Softmax(eddl.Dense(x,
num_classes))

return x
```

Building the model:

- optimizer and its params
- loss function
- metric
- device

```
net = eddl.Model([in_], [out])
eddl.build(
    net,
    eddl.rmsprop(lr),
    ["soft_cross_entropy"],
    ["categorical_accuracy"],
    eddl.CS_GPU(gpus, mem=mem) if
gpus else eddl.CS_CPU()
)
```





Training and validation

Training

```
### Main loop across epochs
for e in range(args.epochs):
    # ...
    ### Training
    ## Set training mode
    eddl.set_mode(net, 1) # TRMODE = 1, TSMODE = 0
    cd.current_split = 0
    cd.rewind_splits(shuffle=True)
    # ...
    ### Looping across batches of training data
    pbar = tqdm(range(num_batches_tr))

    for b_index, b in enumerate(pbar):
        x, y = cd.load_batch()
        rescale_tensor(x) # to range [-1,1] if tf net

        tx, ty = [x], [y]
        eddl.train_batch(net, tx, ty)
        #...
    pbar.close()
```

For each epoch:

- **Set the network mode to Training.** This makes the **dropout** layers work (if indicated)
- **Set the training split** by selecting the **split** with index **0**
- Randomly reshuffle samples in each split
- Load samples along with their labels by means of the **load_batch** **Cassandra Data loader** method
- **Train the network** by using the current batch and labels



Training and validation

Validation

```
eddl.set_mode(net, 0) # Set test mode.
cd.current_split = 1 ## Set validation split
pbar = tqdm(range(num_batches_val))

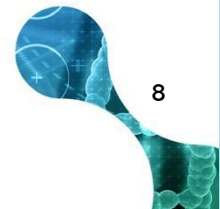
for b_index, b in enumerate(pbar):
    x, y = cd.load_batch()
    rescale_tensor(x) # to range [-1,1] if tf net

    eddl.forward(net, [x])
    output = eddl.getOutput(out)

    result = output.getdata()
    target = y.getdata()

#...
pbar.close()
```

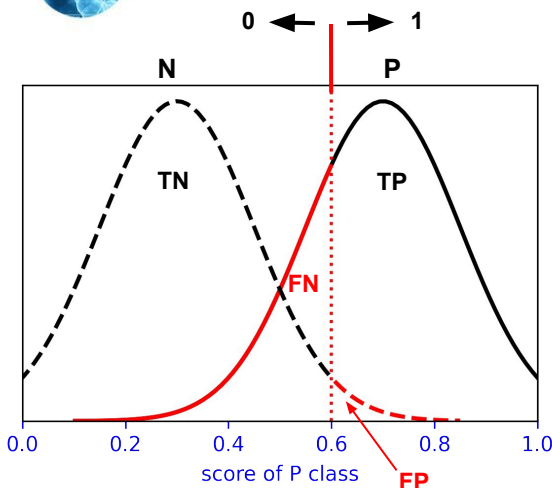
- **Set the network mode to Test.** This deactivates the **dropout** layers (if indicated)
- **Set the validation split** by selecting the **split** with index 1
- Load samples along with their labels by means of the ***load_batch* Cassandra Data loader** method
- Perform a **forward pass** and get **results** from the output of the network





Metrics

Performance Measure: Accuracy, Recall and Precision



Binarization of predictions (th=0.6)

Positive Sample - Pred(P) = 0.75 → [1] → TP
 Positive Sample - Pred(P) = 0.58 → [0] → FN
 Negative Sample - Pred(P) = 0.62 → [1] → FP
 Negative Sample - Pred(P) = 0.22 → [0] → TN

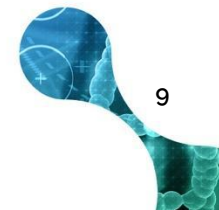
		Actual values	
		N	P
Predicted values	N	TN	FN
	P	FP	TP

Confusion Matrix

$$\text{Accuracy} = \frac{\text{Correct predictions}}{\text{Total}} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Recall} = \frac{\text{Correct pos preds}}{\text{Positive samples}} = \frac{TP}{TP + FN}$$

$$\text{Precision} = \frac{\text{Correct pos preds}}{\text{Positive preds}} = \frac{TP}{TP + FP}$$





Metrics

Performance Measure: Accuracy

We use Accuracy to evaluate the validation set during the training process (Recall and Precision are used for cross-validation and test set evaluation)

- **Pros: Intuitive measure**, a model with an accuracy of 90% means that it correctly classifies 90% of observations

- **Cons:**
 - **No information** on single class performance.



Predictions	Actuals	
	N	P
N	40	0
P	10	50
Rec = 100%		
Prec = 83%		

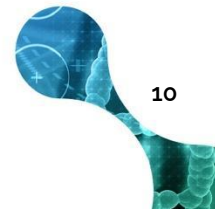
Predictions	Actuals	
	N	P
N	50	10
P	0	40
Rec = 80%		
Prec = 100%		



Acc = 90%

- It **does not work** with **imbalanced datasets**.

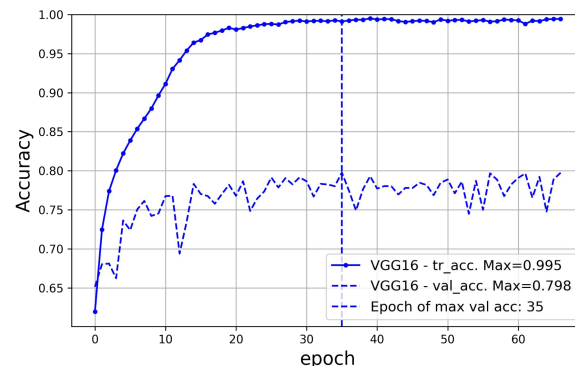
For example, consider we have 98 dogs and 2 cats, and our model provides only dogs as output. The accuracy is $98/100 = 98\%$. **The model has high accuracy but** is obviously not a good model and it **fails to generalize**.



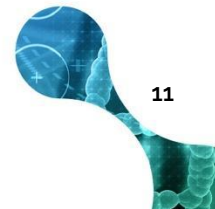


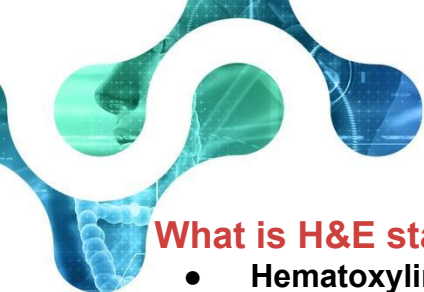
Improving Performance

The Accuracy of the validation set is generally lower than the training set. This is due to overfitting. The model is able to learn the relationship between the training set images and their labels but lacks generalization capabilities.



- **Get More (or better) training data:** Time consuming and demanding task because a pathologist has to annotate new slides
- **Early stopping:** Don't wait for the training loss to converge, just get the model with the best validation accuracy.
- **L1, L2 Regularization:** Add a penalty term to the loss function
- **Dropout:** Some neurons of the selected layers are switched off at training time (remember to reactivate them at evaluation time)
- **Data augmentation:** Random image transformations (e.g.: flip, rotation)





Improving Performance

Stain Normalization

What is H&E stain:

- **Hematoxylin** and **eosin (H&E)** stain is used to emphasize cell and tissue structure details to help pathologists in analyzing biopsies of suspected cancer.
- It is the most widely used stain.
- The hematoxylin stains cell nuclei a purplish blue, and eosin stains the extracellular matrix and cytoplasm pink. The other structures take on different shades, hues, and combinations of these colors.

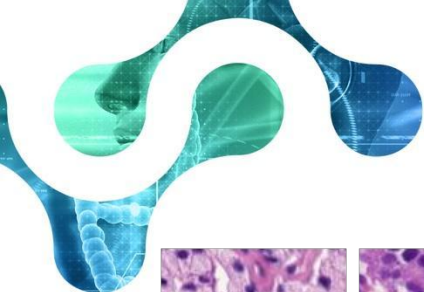
Why Normalize:

- Histological slides from **different labs** (and even across batches within the same lab) **show heterogeneous appearance** as a result of the different preparation and staining procedures.
- This **variability** can heavily **affect the result of automatic image analysis algorithms**.

How we normalize:

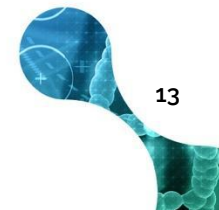
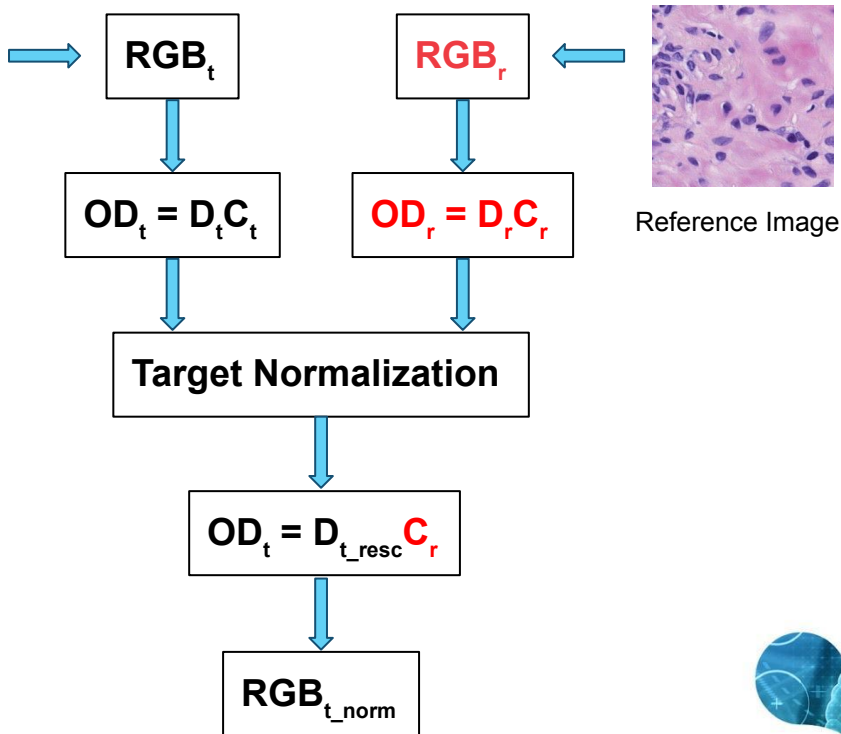
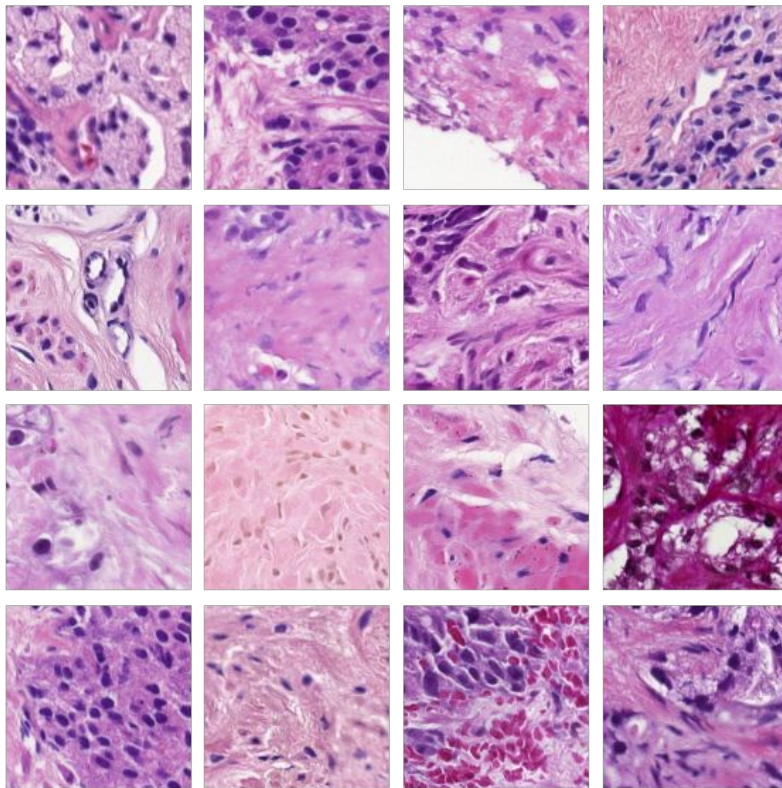
- We used the **Stain separation method** proposed in the paper “**Structure-preserving color normalization and Sparse Stain Separation for Histological Images**”⁽¹⁾.
- **We normalized patches** within a Cassandra data-table **storing the normalized version to a new data-table**. This allowed us to make a direct comparison by using the same splits but taking patches either with or without normalization

(1) A. Vahadane et al., "Structure-Preserving Color Normalization and Sparse Stain Separation for Histological Images," in IEEE Transactions on Medical Imaging, vol. 35, no. 8, pp. 1962-1971, Aug. 2016, doi: 10.1109/TMI.2016.2529665.



Improving Performance

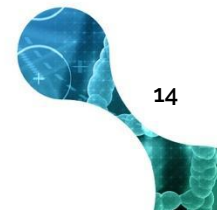
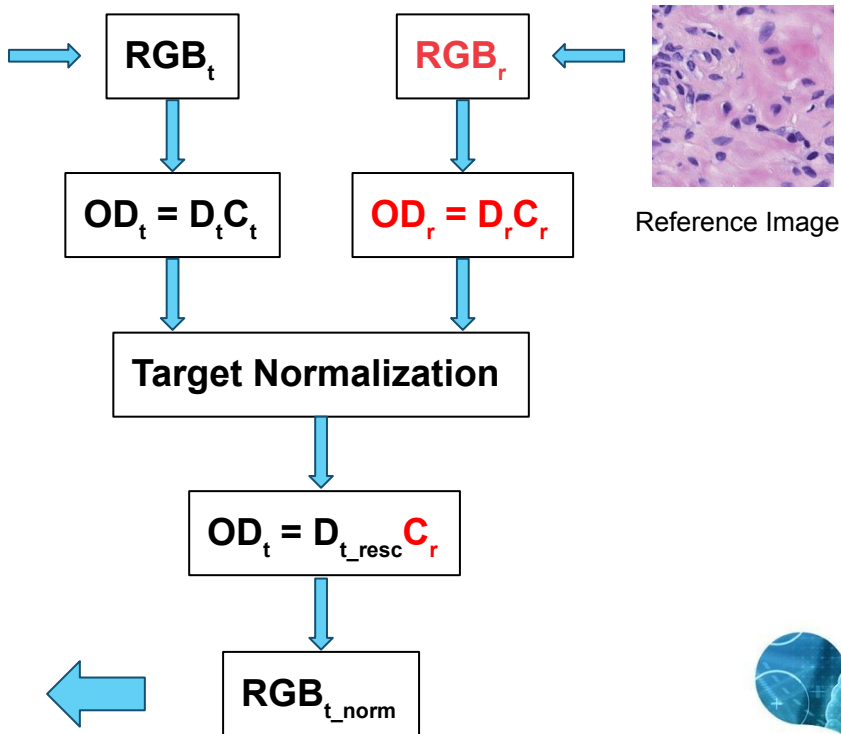
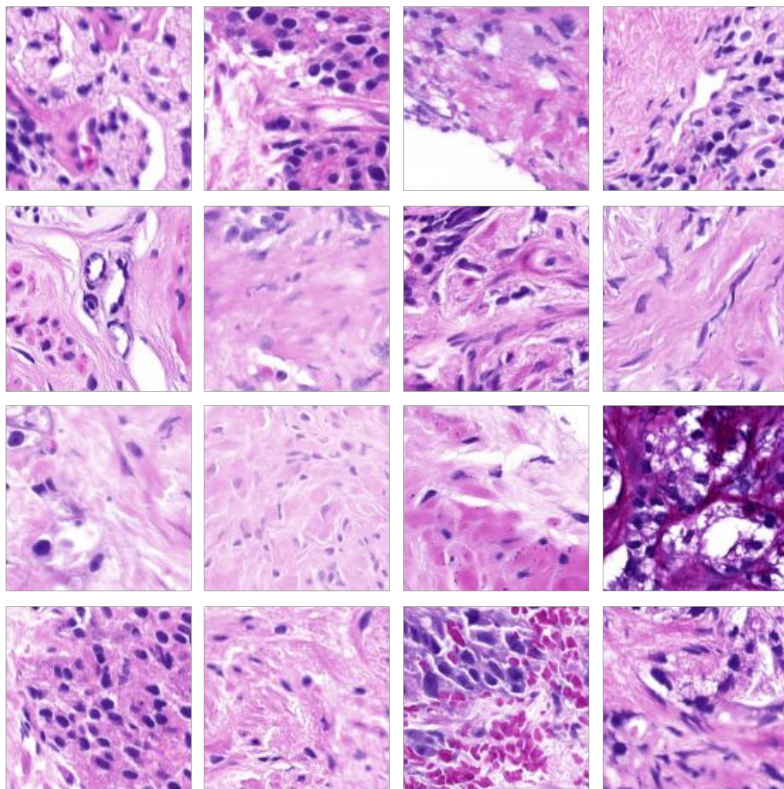
Stain Normalization: Structure-preserving color normalization





Improving Performance

Stain Normalization: Structure-preserving color normalization





Cross-validation data loading

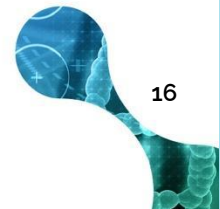
Training

```
for b_index, b in enumerate(pbar):  
    x, y = cd.load_batch_cross(not_splits=val_splits)  
    rescale_tensor(x) ### [0,1] or [-1,1]  
  
    tx, ty = [x], [y]  
    eddl.train_batch(net, tx, ty)  
    #...  
pbar.close()
```

Validation

```
for b_index, b in enumerate(pbar):  
    x, y = cd.load_batch_cross(not_splits=train_splits)  
    rescale_tensor(x)  
  
    eddl.forward(net, [x])  
    output = eddl.getOutput(out)  
    #...  
pbar.close()
```

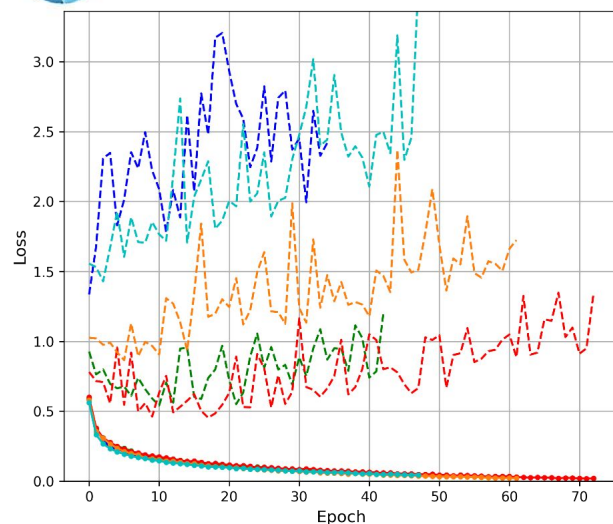
- To make **cross-validation** work and load the right data during the training and validation steps, the **load_batch** function must be replaced with the **load_batch_cross** function.
- **load_batch_cross** takes a list of split indexes as the input argument and load batches from all the splits except those specified in the list
- For example, we have 5 splits and we want to use split 2 for validation. We will provide the list [2] to **load_batch_cross** during the training phase and the list [0,1,3,4] during the validation step



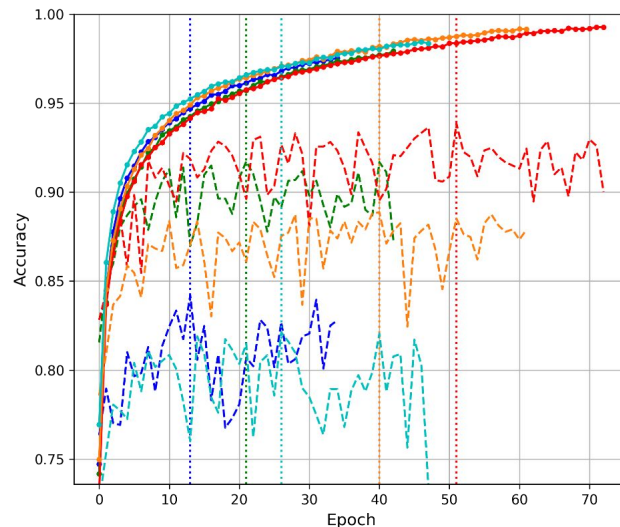


Model evaluation

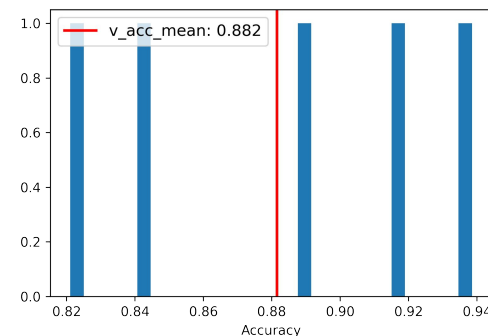
Cross validation results

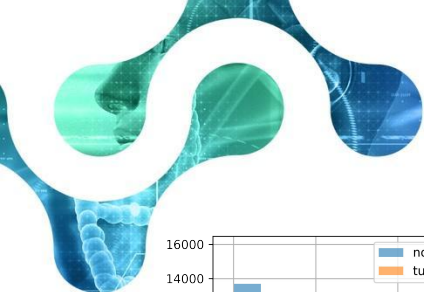


val_0 - Loss min=0.068	val_0 - Val Loss min=1.335
val_1 - Loss min=0.058	val_1 - Val Loss min=0.542
val_2 - Loss min=0.020	val_2 - Val Loss min=0.455
val_3 - Loss min=0.021	val_3 - Val Loss min=0.865
val_4 - Loss min=0.044	val_4 - Val Loss min=1.429



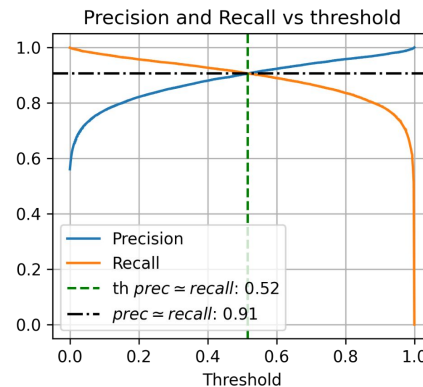
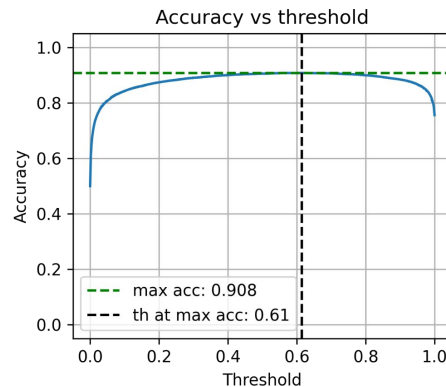
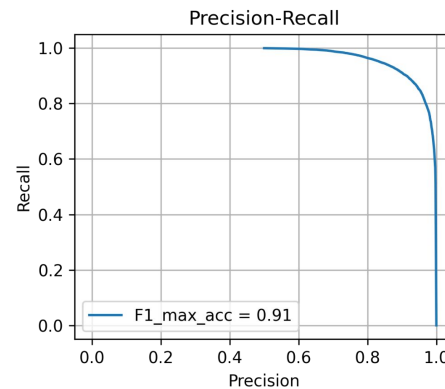
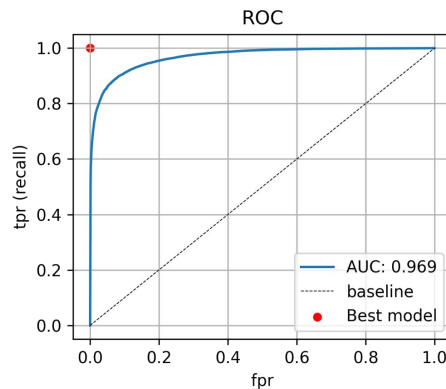
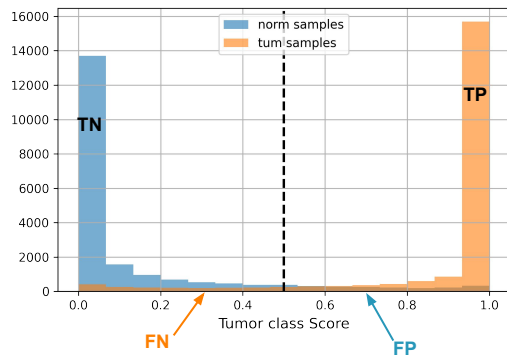
val_0 - tr_acc. Max=0.975	val_0 - val_acc. Max=0.842	Max epoch: 13
val_1 - tr_acc. Max=0.979	val_1 - val_acc. Max=0.917	Max epoch: 21
val_2 - tr_acc. Max=0.993	val_2 - val_acc. Max=0.939	Max epoch: 51
val_3 - tr_acc. Max=0.992	val_3 - val_acc. Max=0.889	Max epoch: 40
val_4 - tr_acc. Max=0.984	val_4 - val_acc. Max=0.821	Max epoch: 26





Model evaluation

Test Set Results

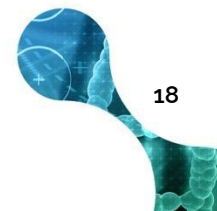


$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

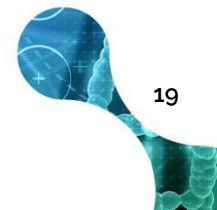
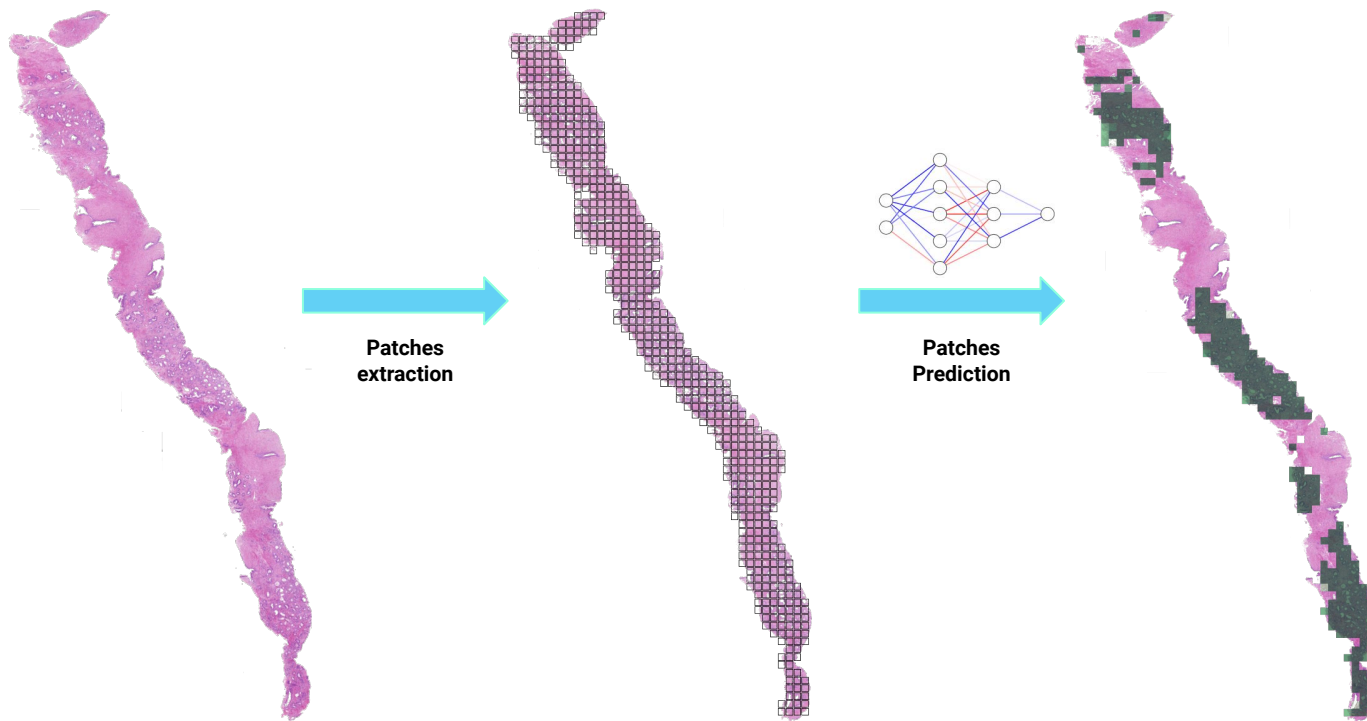
$$\text{Recall(TPR)} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{F1 score} = \frac{2 \times \text{Prec} \times \text{Recall}}{\text{Prec} + \text{Recall}}$$



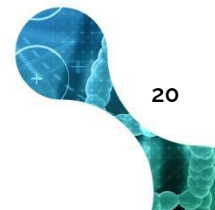
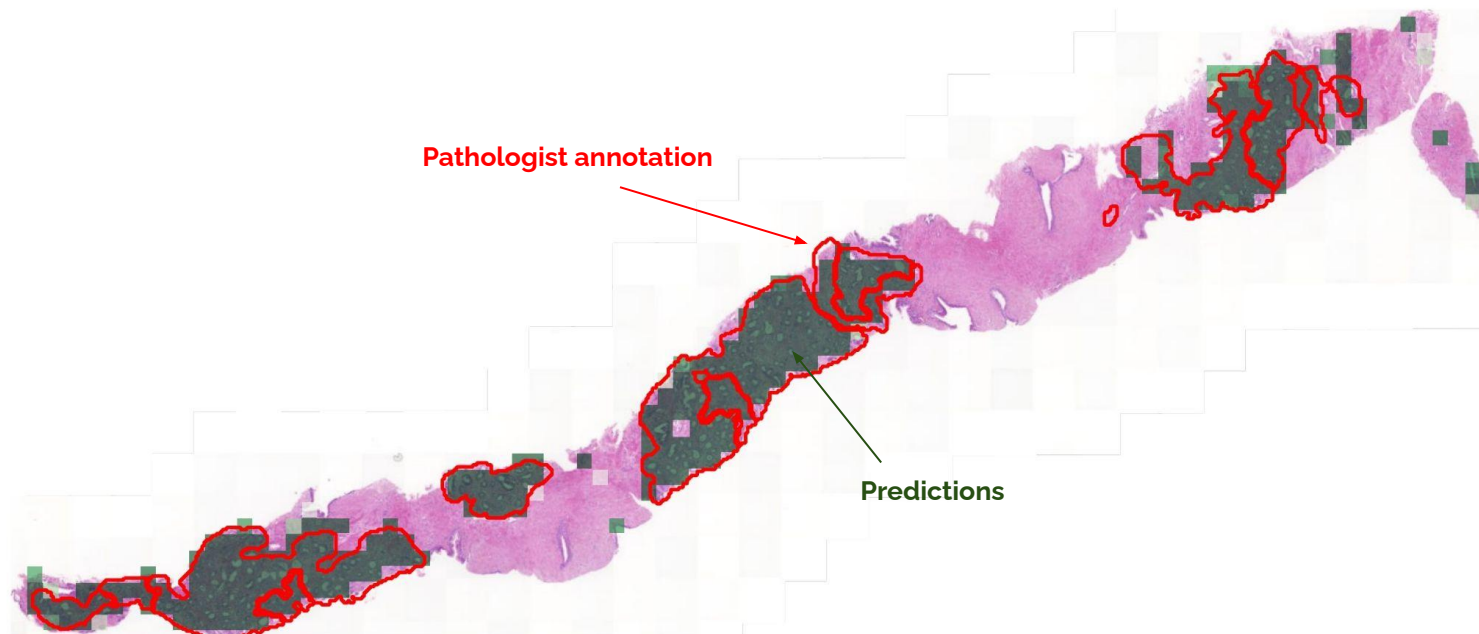


Slide prediction and visualization



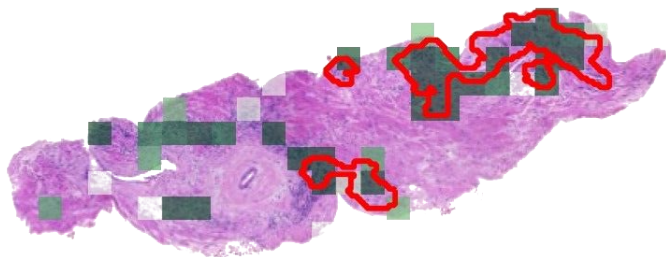


Slide prediction and visualization

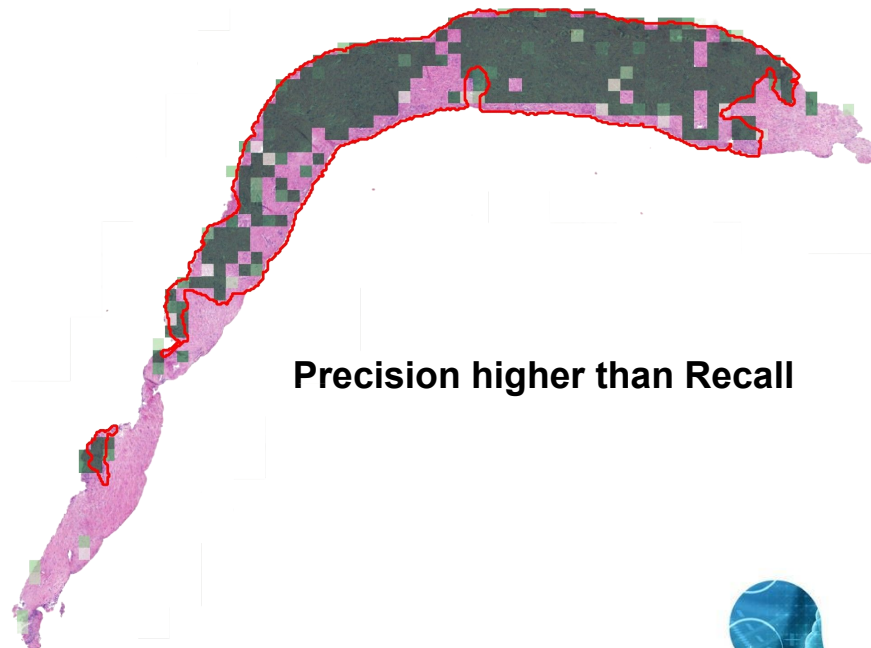




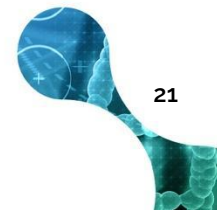
Slide prediction and visualization



Precision lower than Recall




Precision higher than Recall

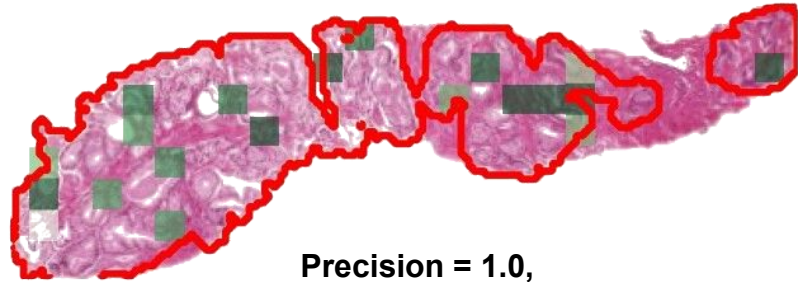




Slide prediction and visualization



**Recall ≈ 1.0 ,
almost no False negatives.
Poor Precision,
plenty of False positive.**

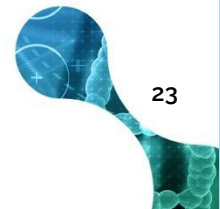


**Precision = 1.0,
no False positive.
Poor recall,
plenty of False negative.**



How to fight mispredictions

- **Data Preprocessing:**
 - Remove noisy patches
 - More sophisticated normalization algorithms (e.g.: GAN based)
- **Try more recent neural network** architectures (e.g.: ResNet, Xception, EfficientNet)
- Use different types of **data augmentations**
- Increase the dataset size by using **more uncorrelated images** (better to have a lot of cases with a small number of slides each rather than having a small number of cases each with many slides)
- **Increase the number of patches** of the class for which the classifier shows worse performance





DEEPHEALTH
Winter School 2022

Thank you!
Q&A



The project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 825111.