

# Presentation 3: Running Multi-Stain Model

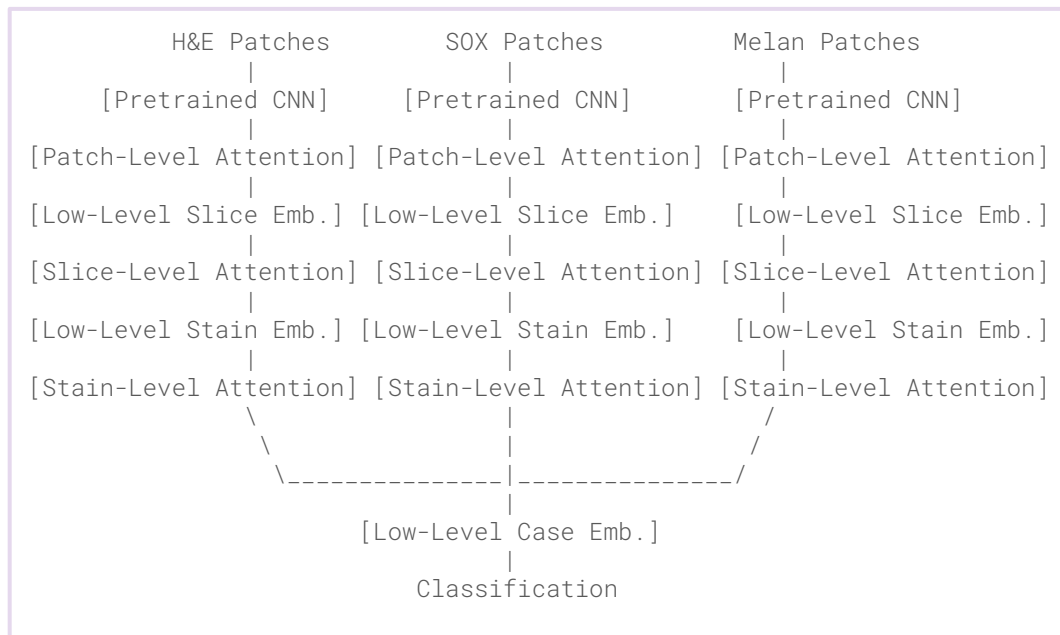
---

STAT 390 | Project 1 | Fall 2025

NOT  
CROSS  
STITCH  
STITCHES

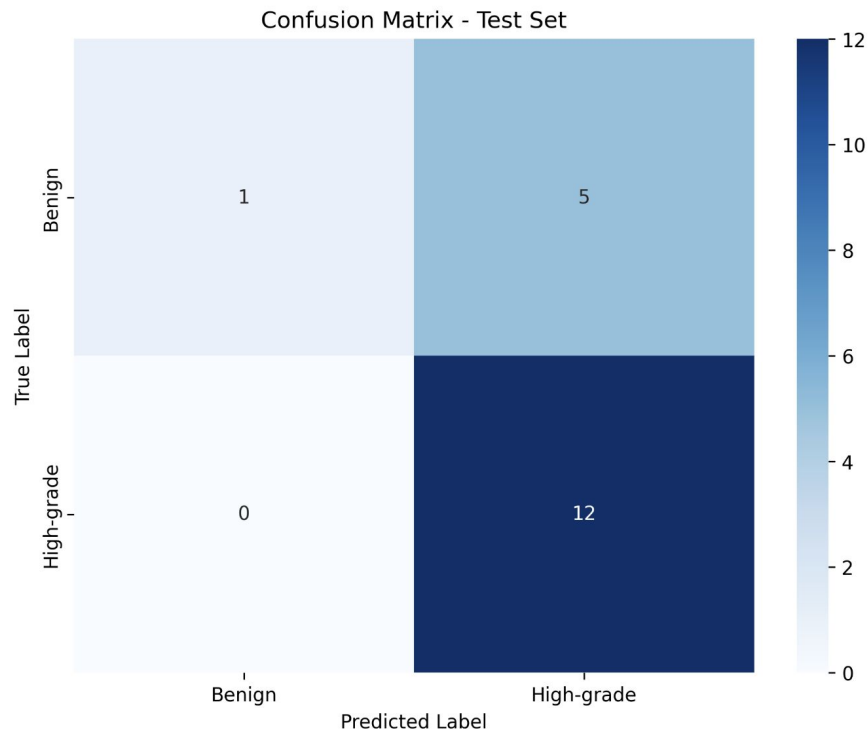


# Multi-Stain Model Overview



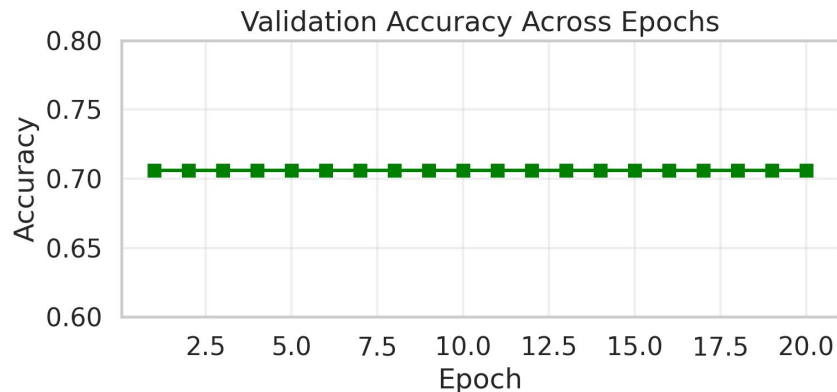
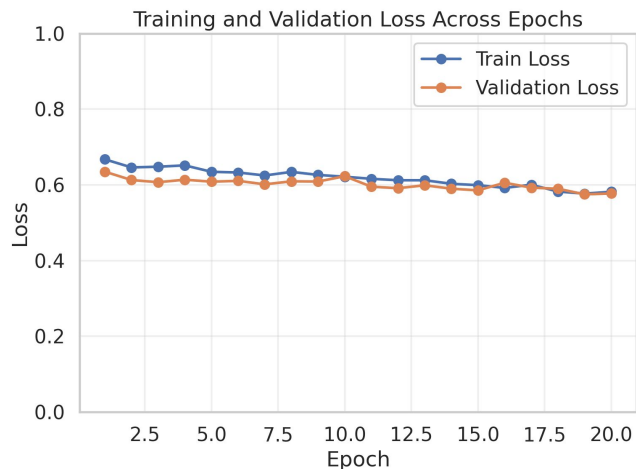
# 1. Initial Results

- Confusion matrix →
- Recall (high grade) =  $12/12 = 1.0$
- Recall (benign) =  $1/6 \approx 0.167$
- Accuracy =  $13/18 \approx 0.722$
- Label distribution per split:  
label 0 1  
split  
test 6 12  
train 17 34  
val 5 12



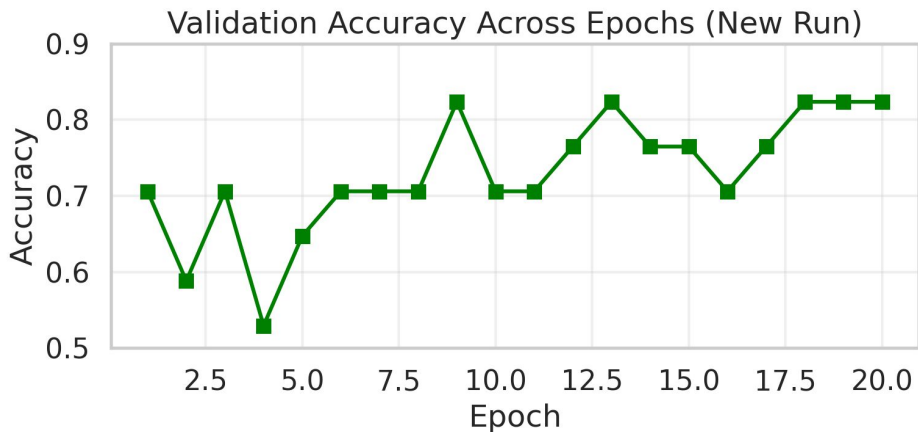
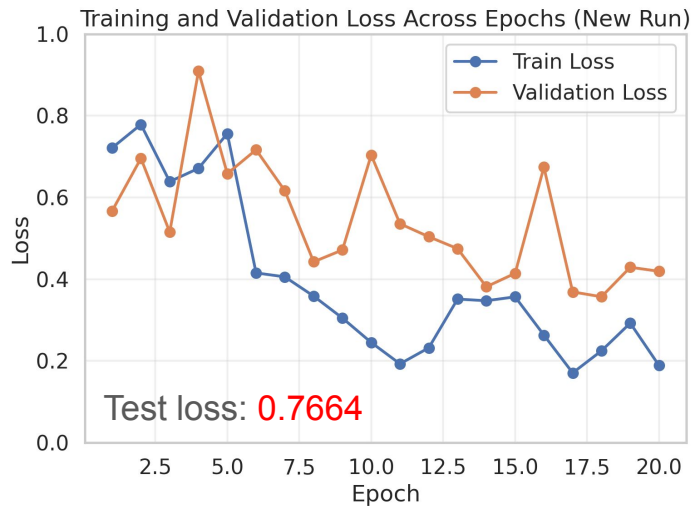
# 1. Initial Results

- Initial model showed strong bias toward predicting label = 1 (High-grade)
  - Misclassified almost all benign samples as high-grade
  - Caused by improper gradient flow
  - Learning rate also too low?



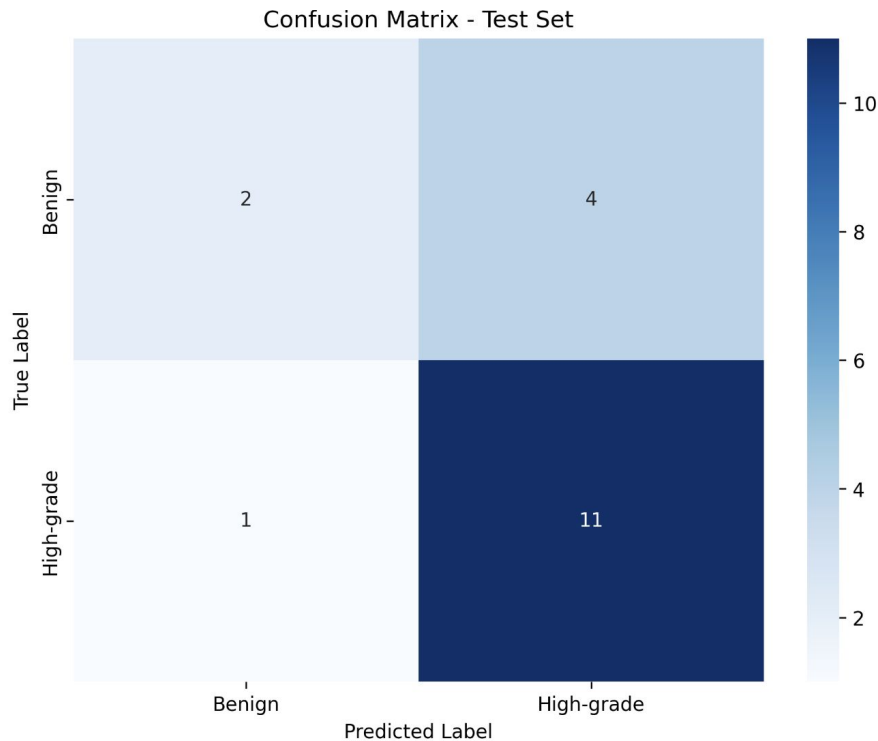
## 2. Updated Results

- The following fixes are made
  - Fixed gradient flow across all three layers (bug)
  - Increased learning rate from  $1e-4$  to  $2e-4$
  - Also 20 epochs

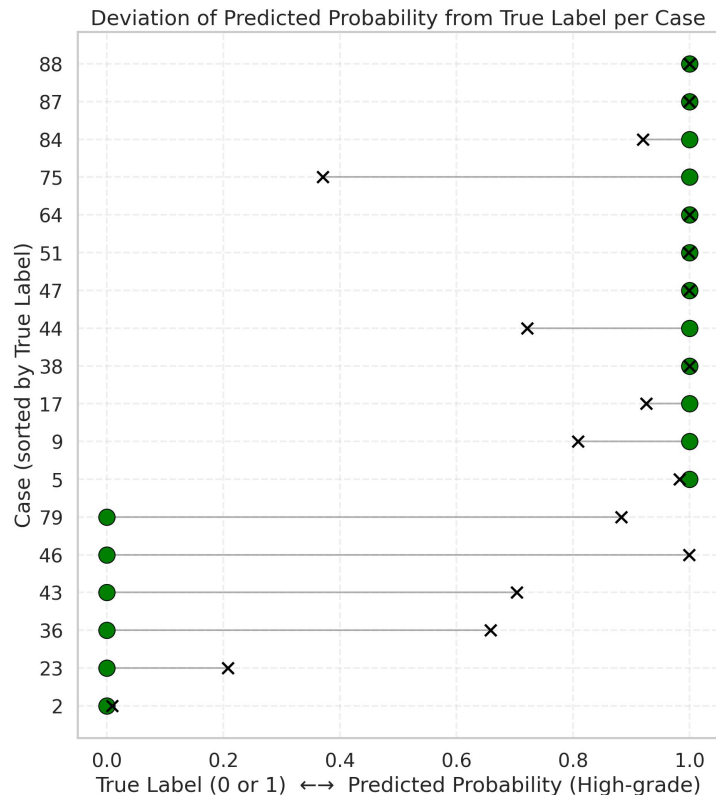


## 2. Updated Results

- Confusion matrix →
- Recall (high grade) =  $11/12 \approx 0.917$
- Recall (benign) =  $2/6 \approx 0.333$
- Accuracy =  $13/18 \approx 0.722$
- Label distribution per split:  
label 0 1  
split  
test 6 12  
train 17 34  
val 5 12



## 2. Updated Results



**Green Dots:** true label

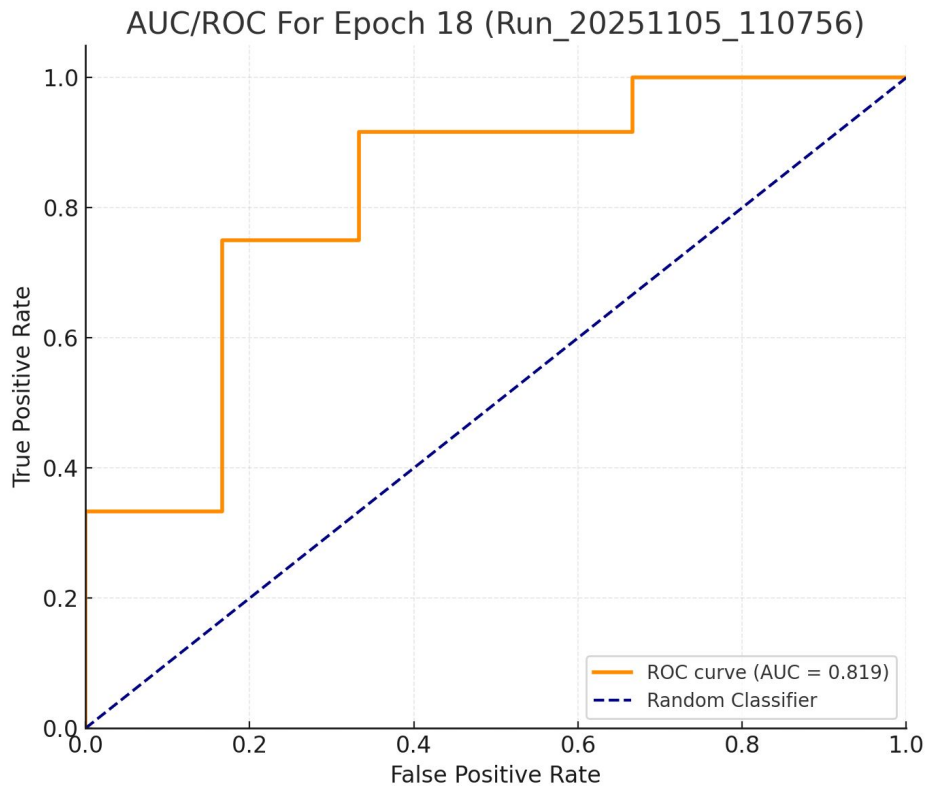
**X:** predicted probability to be positive

**Takeaway:** most predictions tend to skew to the positive end.

*Threshold tuning possible?*



## 2. Updated Results



**AUC = 0.819**

- Solid performance!

## 2. Hyperparameter Tuning: Learning Rate

- Learning rates tested: 1e-5, 2e-5, 5e-5, 1e-4, 2e-4
- Methodology:
  - (1) GPT analyzed baseline val loss/acc
  - (2) Begin with wide grid search → finer grid tuning is running tonight

Learning Rate	Recall	Precision	Accuracy
1e-5	0.917	0.647	0.611
2e-5	0.917	0.647	0.611
5e-5	0.917	0.647	0.611
1e-4	0.917	0.647	0.611
2e-4	0.917	0.733	0.722

Takeaway: Keep trying larger learning rates (necessary with 20 epochs instead of 5)

## 2. Hyperparameter Tuning: Learning Rate

---

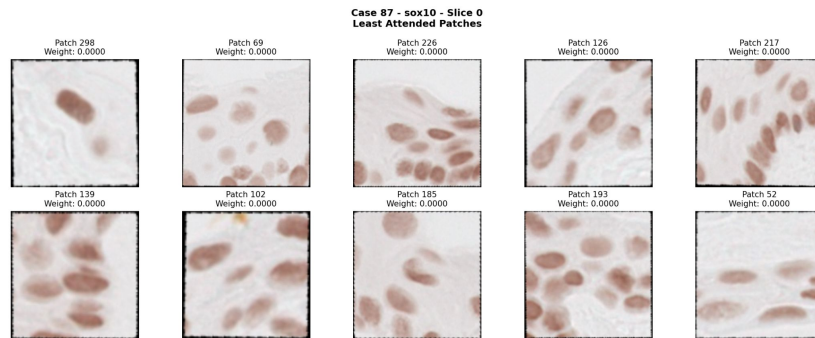
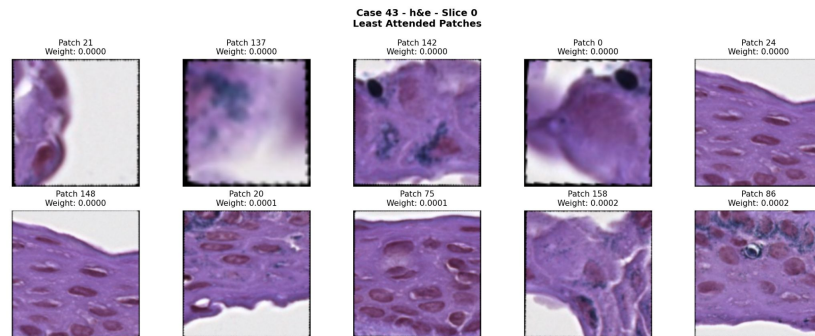
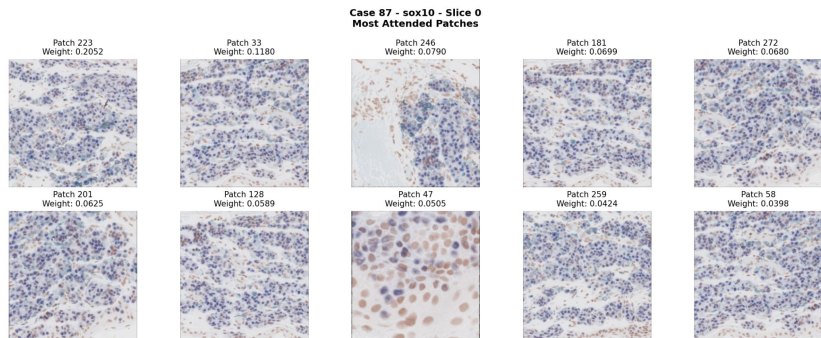
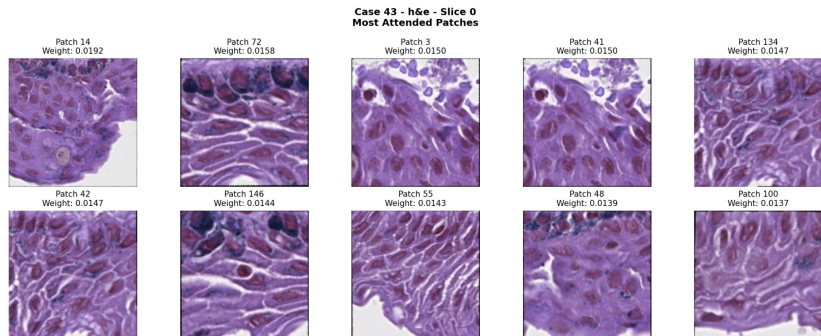
	Fixed LR	Scheduler LR
<b>Pros</b>	<ul style="list-style-type: none"><li>• Stable when properly tuned</li></ul>	<ul style="list-style-type: none"><li>• Reduces overfitting</li><li>• Adapts throughout training process</li></ul>
<b>Cons</b>	<ul style="list-style-type: none"><li>• Not as suitable for multiple epochs &amp; large datasets</li></ul>	<ul style="list-style-type: none"><li>• Longer training times</li></ul>
<b>Next Steps</b>	LRs 1.8e-4, 2.2e-4, 2.5e-4, 3e-4, 5e-4  (code running right now)	<a href="#">Link</a> to code chunk

# 3. Attention Analysis: Patch-Level Results

- Attention mechanism is definitely working
- Patch level (when working really well):

Case 43 (Up): Predicted = 1, Actual = 0

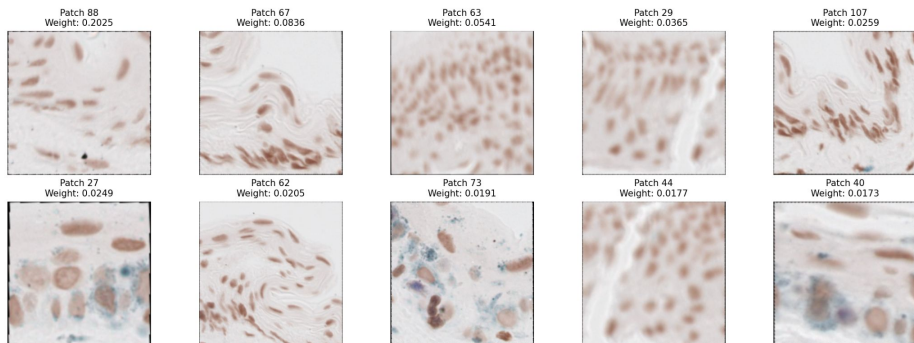
Case 87 (Down): Predicted = 1, Actual = 1



### 3. Attention Analysis: Slice-Level Results

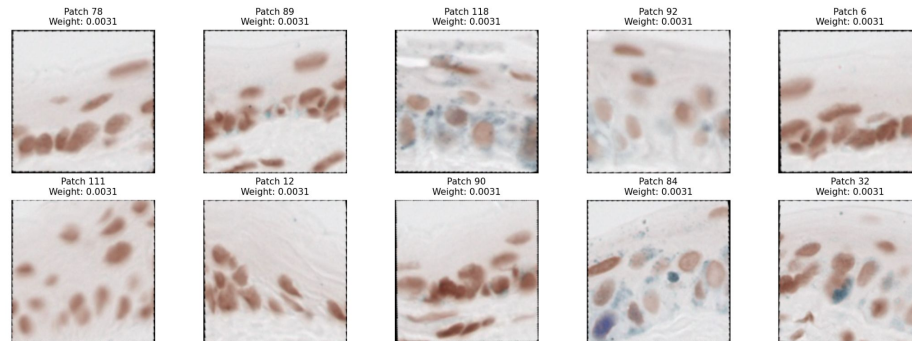
- Patch level (when not that contrastive?):

**Case 75 - sox10 - Slice 0**  
**Most Attended Patches**



Predicted = 0, Actual = 1

**Case 75 - sox10 - Slice 0**  
**Least Attended Patches**



### 3. Attention Analysis: Slice-Level Results

---

- **h&e stain** had highest attention for **15 cases**
- **melan stain** had highest attention for **3 cases**
- **sox10 stain** had highest attention for **0 cases**

From attention\_summary.txt, run\_20251105\_110756

### 3. Attention Analysis: Visualizing Combinations

---

- **Goal:** Visualize patches of each combination–high grade, high attention; low grade, high attention; high grade, low attention; low grade, low attention
- **Approach** [[link to code chunk](#)]:

For each case (high grade, benign):

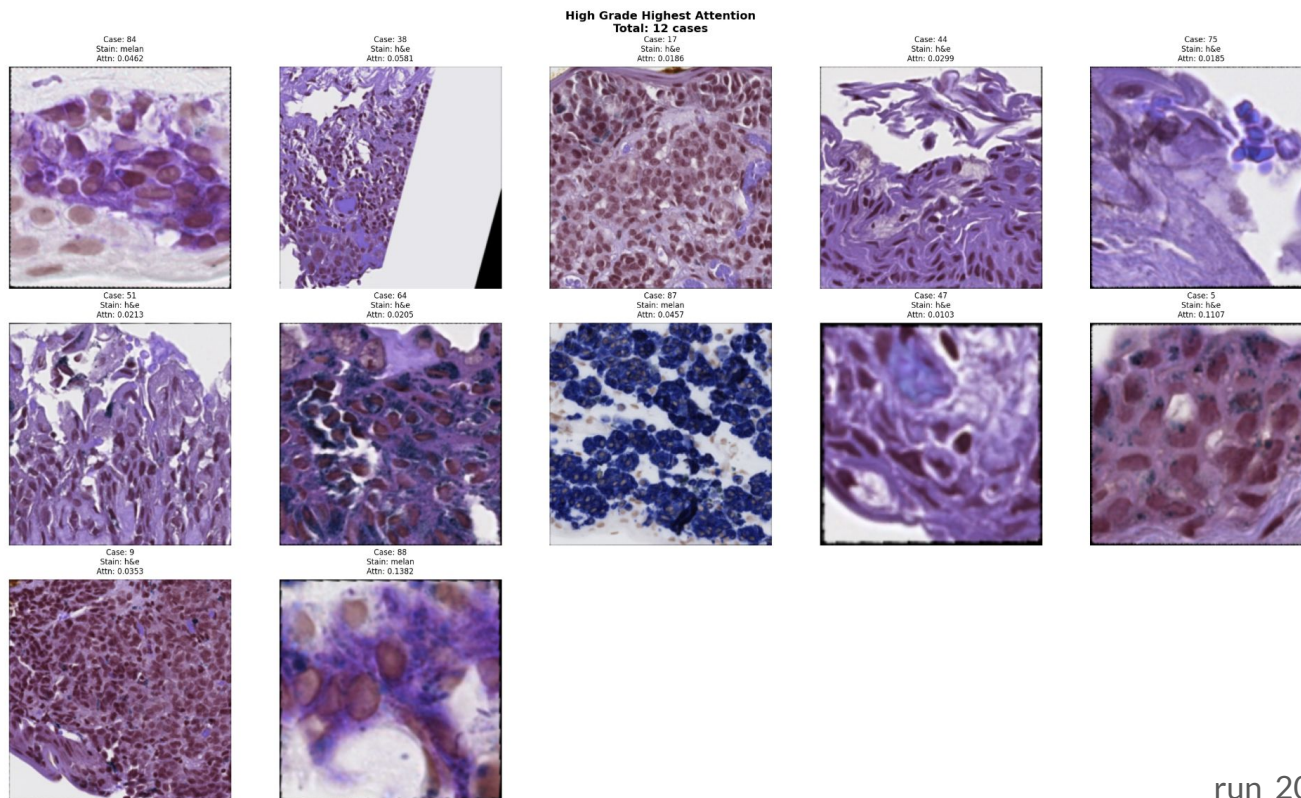
**Highest attention patch:**

- Find highest attention stain (from case\_weights)
- Find highest attention slice within that stain
- Visualize highest attention patch within that slice

**Lowest attention patch:**

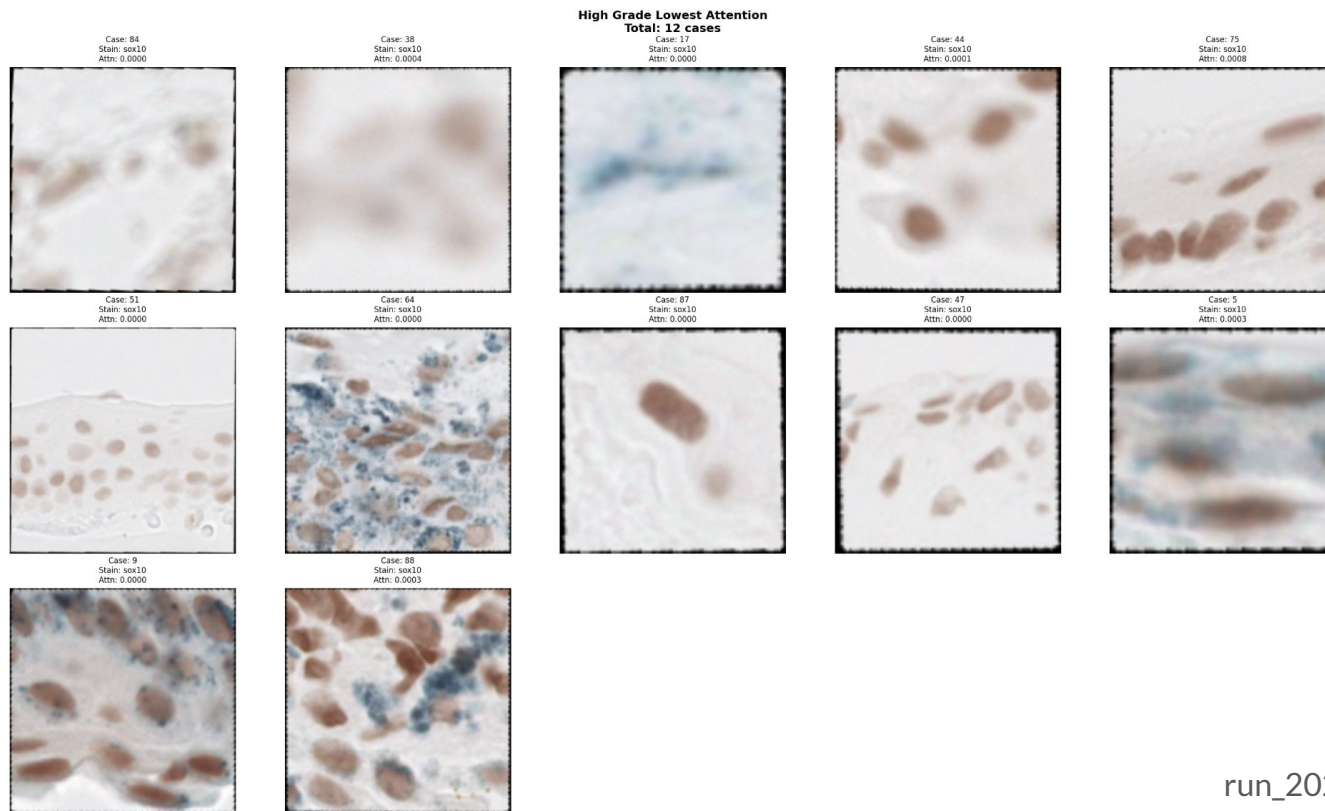
- Find lowest attention stain (from case\_weights)
- Find lowest attention slice within that stain
- Visualize lowest attention patch within that slice

### 3. Attention Analysis: High Grade, High Attention





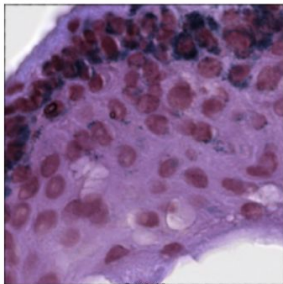
### 3. Attention Analysis: High Grade, Low Attention



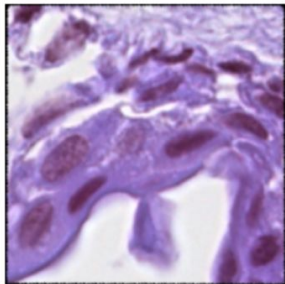
### 3. Attention Analysis: Benign, High Attention

**Benign Grade Highest Attention  
Total: 6 cases**

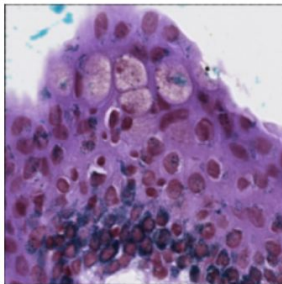
Case: 23  
Stain: h&e  
Attn: 0.0141



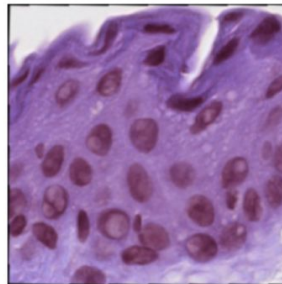
Case: 79  
Stain: h&e  
Attn: 0.0254



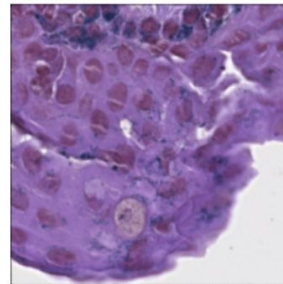
Case: 36  
Stain: h&e  
Attn: 0.0333



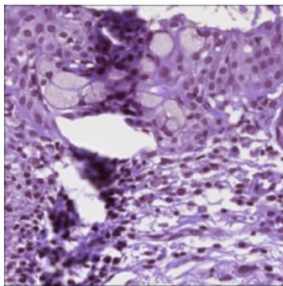
Case: 2  
Stain: h&e  
Attn: 0.0276



Case: 43  
Stain: h&e  
Attn: 0.0116



Case: 46  
Stain: h&e  
Attn: 0.0294

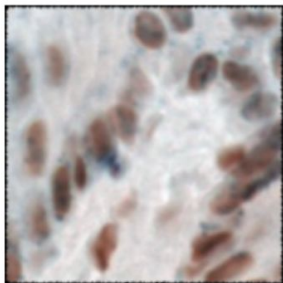


### 3. Attention Analysis: Benign, Low Attention

---

**Benign Grade Lowest Attention  
Total: 6 cases**

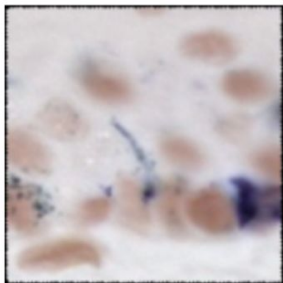
Case: 23  
Stain: sox10  
Attn: 0.0000



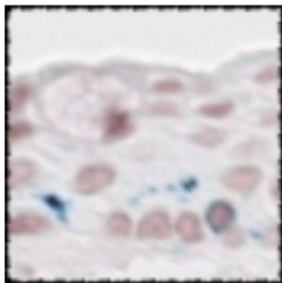
Case: 46  
Stain: sox10  
Attn: 0.0000



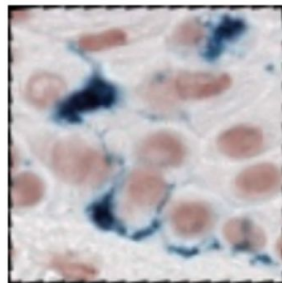
Case: 79  
Stain: melan  
Attn: 0.0009



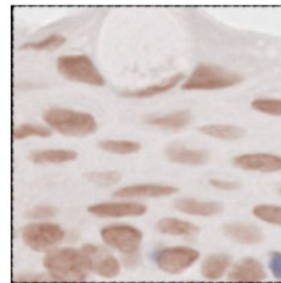
Case: 36  
Stain: melan  
Attn: 0.0092



Case: 2  
Stain: melan  
Attn: 0.0013



Case: 43  
Stain: sox10  
Attn: 0.0000



## 4. Color Experimentation - HED

- HED color space has worse performance than RGB, but also without color normalization
- very close performance for high-grade on slice level

### Patch Level

### Slice Level

### Case Level

#### RGB, 5 epoch, normalization

	precision	recall	f1-score	support	precision	recall	f1-score	support	precision	recall	f1-score	support
Benign	0.31	0.42	0.35	1706	0.60	0.69	0.64	13	0.67	0.67	0.67	6
High-grade CMIL	0.82	0.75	0.78	6287	0.94	0.91	0.92	66	0.83	0.83	0.83	12
accuracy			0.67	7993			0.87	79			0.78	18
macro avg	0.57	0.58	0.57	7993	0.77	0.80	0.78	79	0.75	0.75	0.75	18
weighted avg	0.71	0.67	0.69	7993	0.88	0.87	0.88	79	0.78	0.78	0.78	18

#### Hematoxylin-Eosin-DAB, 5 epoch, no normalization

	precision	recall	f1-score	support	precision	recall	f1-score	support	precision	recall	f1-score	support
Benign	0.26	0.71	0.38	1706	0.53	0.62	0.57	13	0.43	0.50	0.46	6
High-grade CMIL	0.85	0.46	0.59	6287	0.92	0.89	0.91	66	0.73	0.67	0.70	12
accuracy			0.51	7993			0.85	79			0.61	18
macro avg	0.56	0.58	0.49	7993	0.73	0.75	0.74	79	0.58	0.58	0.58	18
weighted avg	0.73	0.51	0.55	7993	0.86	0.85	0.85	79	0.63	0.61	0.62	18

## 4. Color Experimentation - Code Implementation

- **Goal:** try out different color space with new implementation. HSV, HED, LAB
- [link to code chunk](#)
- **Next Steps:**
  - get normalization parameters for each color space using a sample of training patches from each stain (need to migrate code into Quest)
  - Visualize transformed images in attention analysis

```
def compute_hed_mean_std(dataset, sample_bags=300, sample_patches=15):  
    """  
    Compute mean and std of HED channels from a subset of the training dataset.  
    dataset: SliceMILDataset with RGB2HED transform applied  
    sample_bags: how many bags to sample  
    sample_patches: how many patches to take from each bag  
    """  
    all_pixels = []  
  
    # Randomly choose some bags from dataset  
    sampled_indices = random.sample(range(len(dataset)), min(sample_bags, len(dataset)))  
  
    for idx in tqdm(sampled_indices, desc="Sampling HED patches"):  
        bag_imgs, _ = dataset[idx] # (num_patches, 3, H, W)  
        num_patches = bag_imgs.size(0)  
  
        # Randomly pick up to sample_patches from each bag  
        chosen = random.sample(range(num_patches), min(sample_patches, num_patches))  
        sampled_imgs = bag_imgs[chosen] # (M, 3, H, W)  
  
        # Flatten and append (reshape to Nx3)  
        all_pixels.append(sampled_imgs.permute(0, 2, 3, 1).reshape(-1, 3))  
  
    all_pixels = torch.cat(all_pixels, dim=0) # shape: (total_pixels, 3)  
  
    # Compute channel-wise mean and std  
    mean = all_pixels.mean(dim=0)  
    std = all_pixels.std(dim=0)  
  
    return mean, std
```

```
hed_mean, hed_std = compute_hed_mean_std(train_ds, sample_bags=100, sample_patches=10)  
print("HED Mean:", hed_mean)  
print("HED Std:", hed_std)
```

## 5. Next Steps

---

- 1) Investigate ways to make the gradient descent smoother and more stable
- 2) Investigate how slice-level attention varies, especially for slices split into sub-slices  
→ Can we remove slice level for some / all cases?
- 3) Investigate transformations / large grey areas
- 4) More complex attention modules; add more fully-connected layers at the end classification

# References

GenAI tools