# Supplementary Material for *Atlas of Digital Pathology: A Generalized Hierarchical Histological Tissue Type-Annotated Database for Deep Learning*

## A    Supplementary Materials

This document compiles the supplementary materials for the CVPR submission paper ID-6981 under the title of *Atlas of Digital Pathology: A Generalized Hierarchical Histological Tissue Type-Annotated Database for Deep Learning*. In the main paper, we studied the quality of the patch annotations by: (1) training a predictive convolutional neural network, and (2) collecting feedback on annotations from an expert pathologist. While the convolutional neural network learns to associate labels with observed visual patterns and makes consistent predictions while lacking high-level knowledge of label correctness, the pathologist is affected by human inconsistency but can draw on high-level histological knowledge. A label which was erroneously omitted by the ground-truth labeler would be detected by both the neural network and the pathologist, but a label consistently mislabeled as another type would only be detected by the pathologist. In the following three sections, we cover (a) the training error of the neural network, (b) the statistical analysis of the neural network and pathologist validation, and (c) association rule learning of individual WSIs.

### A.1    Training Error Analysis

In Figure 1, we display the plots of the training and validation set accuracy and loss over the 80-epoch training process. All three network architectures (i.e. VGG16, ResNet18, and Inception-V3) are still improving in training accuracy and loss at 80 epochs but validation accuracy and loss have already converged. Also note that VGG16 starts the training slower than both ResNet18 and Inception-V3 but converges to a superior validation accuracy around the $50^{th}$ epoch.

### A.2    Statistical Analysis of Neural Network and Pathologist Validation

#### A.2.1    Overview

In this section, we consider the discordances of both the neural network (VGG16-layer-3+HBR with optimal thresholds on the test set of 1767 patches) and the pathologist (with their suggested label additions and subtractions on 1000 random patches) with the ground-truth labels, which enables us to find possible ground truth labeling errors. Ideally, perfect ground-truth labels would result in perfect confusion matrix metrics.

We also analyze possible reasons for the discordances by examining the residuals in the other classes using a novel metric we call the *mean prediction residuals (MPR)*. In cases where the model predicts a False Positive or a False Negative, a large residual error in another class going in the same direction (i.e. positive for FP, negative for FN) could indicate strong mutual inter-class support and possible label omission. Likewise, a large residual error in the opposite direction (i.e. negative for FP, positive for FN) could indicate strong mutual inter-class opposition and possible label swapping.

#### A.2.2    Confusion Matrix Metrics

In Tables 1, 2, 3, 4, 5, 6, we analyze the confusion matrix metrics of the neural network and the pathologist for those classes with at least one ground-truth exemplar. The neural network does not give predictions for "Undifferentiated" tissue types (i.e. with codes ending in "X"). Overall, the confusion matrix metrical performance for both models is very good and the worst discordances exist for classes with either known consistent mislabeling errors (according to the pathologist) or few training examples (which disadvantages the neural network but not the pathologist).

#### A.2.3    Mean Prediction Residual

**Overview**

The Mean Prediction Residual (MPR) is a metric that we devised to measure the prediction residuals in the other (consequent) classes whenever a discordance exists for a given (antecedent) class. For each discordance (either a false positive or a false negative), the prediction residual is the difference between a target label and its predicted score, where their mean is the MPR. There are two types of MPR:

1. FP-MPR (for false positives)

2. FN-MPR (for false negatives)

(a) Training Accuracy, Level 1    (b) Training Accuracy, Level 2    (c) Training Accuracy, Level 2+HBP    (d) Training Accuracy, Level 3    (e) Training Accuracy, Level 3+HBP

(f) Validation Accuracy, Level 1    (g) Validation Accuracy, Level 2    (h) Validation Accuracy, Level 2+HBP    (i) Validation Accuracy, Level 3    (j) Validation Accuracy, Level 3+HBP

(k) Training Loss, Level 1    (l) Training Loss, Level 2    (m) Training Loss, Level 2+HBP    (n) Training Loss, Level 3    (o) Training Loss, Level 3+HBP

(p) Validation Loss, Level 1    (q) Validation Loss, Level 2    (r) Validation Loss, Level 2+HBP    (s) Validation Loss, Level 3    (t) Validation Loss, Level 3+HBP

Figure 1. *Training progress plots for all three network architectures across all five training configurations and four metrics: the rows of the figure correspond to different metrics and the columns correspond to different training configurations. Original values are shown as solid lines, smoothed values are shown as translucent lines.*

Table 1. *True Positive Rate (TPR)*

| Confusion Matrix Metric | |
| --- | --- |
| Neural Network | Agrees well with ground-truth positive labels except for N.G.M (which has few training examples) |
| Pathologist | Agrees with ground-truth positive labels except for E.M.U, E.T.U, E.T.O, H.X, S.M.S, S.C.X, G.N, and G.X (which are known to have systematic mislabeling errors) |

Table 2. *False Positive Rate (FPR)*

| Confusion Matrix Metric | |
|---|---|
| |  FPR bar chart: Neural Network vs Pathologist across Existing Ground-Truth Labels (E.M.S, E.M.U, E.M.O, E.T.S, E.T.U, E.T.O, E.T.X, E.P, C.D.I, C.D.R, C.L, C.X, H.E, H.K, H.Y, H.X, S.M.C, S.M.S, S.E, S.C.H, S.C.X, S.R, A.W, A.B, A.M, M.M, M.K, N.P, N.R.B, N.R.A, N.G.M, N.G.W, N.G.X, G.O, G.N, G.X, T) |
| Neural Network | Predicts more false positives in general (perhaps it is overly sensitive to small regions of tissue than humans) |
| Pathologist | Agrees well with most ground-truth positive labels except for E.M.O (which are known to be consistently mislabeled as E.M.U in the ground truth) |

Table 3. *True Negative Rate (TNR)*

| Confusion Matrix Metric | |
|---|---|
| |  TNR bar chart: Neural Network vs Pathologist across Existing Ground-Truth Labels (E.M.S, E.M.U, E.M.O, E.T.S, E.T.U, E.T.O, E.T.X, E.P, C.D.I, C.D.R, C.L, C.X, H.E, H.K, H.Y, H.X, S.M.C, S.M.S, S.E, S.C.H, S.C.X, S.R, A.W, A.B, A.M, M.M, M.K, N.P, N.R.B, N.R.A, N.G.M, N.G.W, N.G.X, G.O, G.N, G.X, T) |
| Neural Network | Highly agrees with ground-truth negative labels |
| Pathologist | Highly agrees with ground-truth negative labels, except for E.M.O (related to consistent mislabeling with E.M.U mentioned above) |

The FP-MPR between an antecedent class $A$ and a consequent class $B$ is defined as follows:

$$\text{FP-MPR}(A, B) = \mathbb{E}[\text{Pred}_B(i) - \text{Targ}_B(i)]\forall i \\ \text{s.t. } \text{Targ}_A(i) = 0, \text{Pred}_A(i) = 1. \quad (1)$$

The FN-MPR between an antecedent class $A$ and a consequent class $B$ is defined as follows:

$$\text{FN-MPR}(A, B) = \mathbb{E}[\text{Pred}_B(i) - \text{Targ}_B(i)]\forall i \\ \text{s.t. } \text{Targ}_A(i) = 1, \text{Pred}_A(i) = 0. \quad (2)$$

These pairwise relationships between antecedent and consequent classes can be displayed in a matrix, with each row corresponding with an antecedent for which a discordance exists and the columns corresponding to the consequent classes for which the prediction residual is calculated - this is known as the MPR matrix. Antecedent classes without any discordances can be shown as black rows but still re-appear as consequent classes in the columns. Another way to understand the MPR matrix is to show the consequent classes with the highest and lowest MPR value for each row of the MPR matrix. This is called the Max/Min MPR Table.

**False Positive-Mean Prediction Residual (FP-MPR)**

In Table 7, we show the FP-MPR matrices for both the neural network and the pathologist in the second column, and the max/min FP-MPR tables in the third column. For the neural network, we observe that the most mutually-supporting classes already co-occur frequently in the data while while those in strong opposition are similar in appearance. A similar pattern is observed for the pathologist False Positives. However, since the neural network picks up inconsistencies in the labeling and the pathologist picks up both inconsistencies and systematic mislabeling (as mentioned earlier), then classes with much more negative FP-MPR values for the pathologist than the neural network are likely to have systematic mislabeling, such as the FP-MPR from G.O to G.N, which drops from just $-0.16129$ for the neural network to $-0.75758$ for the pathologist.

**False Negative-Mean Prediction Residual (FN-MPR)**

Table 4. *False Negative Rate (FNR)*



**FNR**

Existing Ground-Truth Labels

Legend: Neural Network, Pathologist

| Confusion Matrix Metric | |
|---|---|
| Neural Network | Agrees well with ground-truth negative labels except for N.G.M (which has few training examples) |
| Pathologist | Generally agrees with the ground-truth positive labels, except for E.M.U, E.T.U, E.T.O, H.X, S.M.S, S.C.X, G.N, and G.X (which are known to have systematic mislabeling errors) |

Table 5. *Accuracy (ACC)*



**ACC**

Existing Ground-Truth Labels

Legend: Neural Network, Pathologist

| Confusion Matrix Metric | |
|---|---|
| Neural Network | Highly accurate across all classes |
| Pathologist | Highly accurate across all classes |

In Table 8, we show the FN-MPR matrices for both the neural network and the pathologist in the second column, and the max/min FN-MPR tables in the third column. For the neural network and the pathologist, we observe (similarly to FP) that classes with strong mutual support tend to co-occur frequently in the data while those in strong opposition are similar in appearance. Again, as for FP-MPR, those classes with much more positive FP-MPR values for the pathologist than the neural network are likely to have systematic mislabeling. For example, the FN-MPR from E.M.U to E.M.O rises from just $0.21277$ for the neural network to $0.79487$ for the pathologist.

## B  Association Rule Learning of WSI Scoring

In this section we study (1) co-occurence network and (2) associate rule learning (ARL) (introduced in the submitted paper draft) using six different WSIs shown in in Table 9 (first row). Note that slide numbers 1, 2, and 3 here are the same slides used to demonstrates the heatmap representation in Figure 5 of the submitted paper draft. The circular co-occurrence of each slide is demonstrated in Table 9 using different number of image patches per slide, where the labels of each patch (extracted from individual WSI) are predicted by different levels of VGG16 trained network. The nodes of co-occurrence network here share the similar connections of the ADP co-occurrence shown in Figure 2 of the submitted paper draft.

The results of applying Apriori ARL algorithm to the predicted labels (driven from VGG16-level-3+HBP) of each WSI are also shown in Tables 10 to Tables 15. The selected consequent labels here are mainly similar to the ARL results of ADP Atlas demonstrated in Table 2 of the submitted paper draft. However, we notice dissimilarities in the antecedent itemsets between the selected WSIs here and the ones used to populate the ADP database. This is mainly because (a) they are related to different tissue cases; (b) different levels of confidence are selected as the best candidates. In fact, the majority of WSIs here are selected from GI tracts. In conclusion, we observe the following

1. there are no G.N for E.M.U (because no endocrine glands in GI)

Table 6. *F1 Score (F1)*



| Confusion Matrix Metric | |
|---|---|
| Neural Network | Has low F1 score for E.M.S, E.M.O, H.K, and N.G.M (which were accidentally omitted in the ground truth) |
| Pathologist | Has low F1 score for E.M.U, E.M.O, E.T.U, S.C.X, N.R.A, G.N, and G.X (which are known to have systematic mislabeling errors) |

2. there are no M.K, T for C.D.I and instead has M.M, A.W, and T (because no skeletal muscle in GI)

3. there are no N.P and N.R.B (because generally no nervous tissue in GI)

4. there are no H.Y for G.O and T (because they have less lymphocytes near exocrine gland and transport vessels)

In particular, between the ARL of Slide #4 shown in Table 13 and the other slides, we observe

1. Slide #4 ARL has E.M.S, C.L, H.K instead of E.M.U, E.T.U, H.E, H.K for H.Y (because Slide #4 lacks the exocrine glands with those epithelia - H.Y occurs with C.L instead)

2. Slide #4 ARL has H.K, H.Y, M.M instead of H.K, H.E, E.M.S, E.M.O for C.L (because Slide #4 again lacks the exocrine glands with those epithelia - C.L occurs with M.M instead)

Table 7. *FP-MPR Matrices and Max/Min Tables for Neural Network and Pathologist*

| Model | FP-MPR Matrix | Max/Min FP-MPR Table |
|---|---|---|
| Neural Network |  | |
| Pathologist |  | |

Neural Network Max/Min FP-MPR Table:

| Antecedent Class | Max FP-MPR Class | Max FP-MPR | Min FP-MPR Class | Min FP-MPR |
|---|---|---|---|---|
| E.M.S | T | 0.25974 | C.D.I | $-0.02597$ |
| E.M.U | E.T.U | 0.053571 | E.M.O | $-0.13393$ |
| E.M.O | H.Y | 0.118421 | E.T.O | $-0.10526$ |
| E.T.S | C.D.I | 0.25 | E.M.U | $-0.25$ |
| E.T.U | E.M.U | 0.09322 | E.T.O | $-0.08475$ |
| E.T.O | H.Y | 0.166667 | E.M.O | $-0.16667$ |
| C.D.I | E.T.U | 0.1 | C.L | $-0.33333$ |
| C.L | H.Y | 0.072727 | C.D.I | $-0.29091$ |
| H.E | E.T.U | 0.112676 | E.T.O | $-0.04225$ |
| H.K | E.M.O | 0.105263 | H.E | $-0.10526$ |
| H.Y | E.M.O | 0.098765 | C.D.I | $-0.03704$ |
| S.M.S | E.M.S | 0 | E.M.S | 0 |
| S.R | M.K | 1 | M.M | $-1$ |
| A.W | M.M | 0.285714 | C.L | $-0.28571$ |
| M.M | E.M.S | 0.098361 | C.L | $-0.03279$ |
| M.K | S.R | 0.5 | M.M | $-0.5$ |
| N.P | T | 0.285714 | E.M.U | $-0.42857$ |
| N.R.B | E.M.U | 0 | N.G.M | $-0.58333$ |
| N.G.M | N.R.B | 0.111111 | E.M.S | $-0.11111$ |
| G.O | E.M.U | 0.096774 | G.N | $-0.16129$ |
| G.N | E.M.S | 0 | E.M.S | 0 |
| T | E.M.S | 0.148438 | C.L | $-0.04688$ |

Pathologist Max/Min FP-MPR Table:

| Antecedent Class | Max FP-MPR Class | Max FP-MPR | Min FP-MPR Class | Min FP-MPR |
|---|---|---|---|---|
| E.M.S | T | 0.33333 | E.M.U | $-0.11111$ |
| E.M.U | G.O | 0.22222 | E.T.U | $-0.55556$ |
| E.M.O | G.O | 0.054054 | E.M.U | $-0.71429$ |
| E.T.S | E.M.S | 0 | E.T.U | $-1$ |
| E.T.O | E.M.S | 0 | E.T.U | $-1$ |
| C.D.I | E.M.S | 0 | C.L | $-0.4$ |
| C.D.R | E.M.S | 0 | E.M.S | 0 |
| C.L | E.M.S | 0.038462 | C.D.I | $-0.57692$ |
| H.E | E.M.O | 0.11111 | E.T.U | $-0.22222$ |
| H.K | E.M.O | 0.16667 | E.M.U | $-0.5$ |
| H.Y | E.M.O | 0.5 | E.M.U | $-0.5$ |
| S.M.C | A.W | 0.33333 | S.M.S | $-1$ |
| A.W | C.L | 0.16667 | C.D.I | $-0.33333$ |
| M.M | E.M.O | 0.4 | M.K | $-0.6$ |
| M.K | E.M.S | 0 | E.M.S | 0 |
| N.R.B | N.G.X | 1 | E.M.S | 0 |
| N.R.A | C.L | 0.16667 | C.D.I | $-0.16667$ |
| N.G.X | N.R.B | 0.083333 | C.D.I | $-0.083333$ |
| G.O | E.M.O | 0.42424 | G.N | $-0.75758$ |
| T | E.M.S | 0.15789 | E.T.U | $-0.10526$ |

Table 8. *FN-MPR Matrices and Max/Min Tables for Neural Network and Pathologist*

| Model | FN-MPR Matrix | Max/Min FN-MPR Table |
|---|---|---|
| Neural Network |  | |
| Pathologist |  | |

**Neural Network — Max/Min FN-MPR Table**

| Antecedent Class | Max FN-MPR Class | Max FN-MPR | Min FN-MPR Class | Min FN-MPR |
|---|---|---|---|---|
| E.M.S | E.M.U | 0.12632 | T | −0.34737 |
| E.M.U | E.M.O | 0.21277 | E.T.U | −0.042553 |
| E.M.O | E.M.U | 0.3 | H.K | −0.1 |
| E.T.S | E.M.O | 0.2 | E.M.S | −0.2 |
| E.T.U | E.M.O | 0.11905 | H.K | −0.071429 |
| E.T.O | E.T.U | 0.45833 | H.K | −0.041667 |
| E.P | E.M.O | 0.33333 | H.E | −0.33333 |
| C.D.I | C.L | 0.25397 | H.E | −0.031746 |
| C.D.R | C.D.I | 0.75 | C.L | −0.5 |
| C.L | C.D.I | 0.10638 | H.Y | −0.042553 |
| H.E | E.T.U | 0.05814 | T | −0.12791 |
| H.K | E.M.U | 0.125 | E.M.O | −0.053571 |
| H.Y | E.M.U | 0.072289 | E.M.O | −0.024096 |
| S.M.S | E.M.S | 0 | E.M.S | 0 |
| S.E | E.M.S | 0 | E.M.S | 0 |
| S.R | E.M.S | 0 | E.M.S | 0 |
| A.W | T | 0.11111 | E.M.U | −0.11111 |
| A.B | T | 1 | M.K | −1 |
| A.M | E.M.S | 0 | E.M.S | 0 |
| M.M | C.D.I | 0.1087 | C.L | −0.086957 |
| M.K | M.M | 0.28571 | H.Y | −0.14286 |
| N.P | E.M.S | 0 | C.D.I | −1 |
| N.R.B | H.E | 0.11111 | E.M.S | −0.11111 |
| N.R.A | H.E | 0.2 | E.M.S | −0.2 |
| N.G.M | N.R.B | 0.17021 | T | −0.17021 |
| G.O | T | 0.1 | E.M.O | −0.1 |
| G.N | G.O | 0.29412 | E.M.S | −0.11765 |
| T | E.M.U | 0.088235 | E.M.S | −0.25 |

**Pathologist — Max/Min FN-MPR Table**

| Antecedent Class | Max FN-MPR Class | Max FN-MPR | Min FN-MPR Class | Min FN-MPR |
|---|---|---|---|---|
| E.M.S | E.M.S | 0 | E.M.S | 0 |
| E.M.U | E.M.O | 0.79487 | E.T.U | −0.36325 |
| E.T.S | E.M.S | 0 | E.M.S | 0 |
| E.T.U | E.M.O | 0.52055 | E.M.U | −0.3653 |
| E.T.O | E.M.O | 0.45455 | E.T.U | −0.54545 |
| E.T.X | E.M.O | 1 | E.M.S | 0 |
| C.D.I | C.L | 0.34091 | E.M.U | −0.022727 |
| C.D.R | E.M.S | 0 | E.M.S | 0 |
| C.L | C.D.I | 0.30769 | E.T.U | −0.076923 |
| C.X | E.M.S | 0 | E.M.S | 0 |
| H.E | E.M.O | 0.2 | T | −0.4 |
| H.Y | E.M.S | 0 | C.L | −1 |
| H.X | H.E | 0.33333 | S.M.S | −0.33333 |
| S.M.S | S.M.C | 0.6 | C.D.I | −0.2 |
| S.C.X | E.M.S | 0 | S.M.S | −0.5 |
| S.R | E.M.S | 0 | S.M.S | −0.5 |
| A.W | E.M.S | 0 | E.M.S | 0 |
| A.M | E.M.S | 0 | S.M.S | −1 |
| M.M | E.M.O | 0.18182 | E.M.U | −0.18182 |
| M.K | M.M | 0.6 | E.T.U | −0.2 |
| N.G.M | N.R.A | 1 | E.M.S | 0 |
| N.G.X | E.M.S | 0 | E.M.S | 0 |
| G.N | G.O | 0.64103 | E.T.U | −0.051282 |
| G.X | G.O | 0.75 | E.M.U | −0.25 |
| T | E.M.S | 0 | C.D.I | −0.13333 |

Table 9. *Circular co-occurrence networks of six different WSIs predicted by three different levels of network training (including augmented HBP layers). Within each WSI, different number of image patches are extracted for label prediction.*

| # patch | Slide-1 | Slide-2 | Slide-3 | Slide-4 | Slide-5 | Slide-6 |
|---|---|---|---|---|---|---|
| | 4347 | 1734 | 4919 | 3236 | 2565 | 4587 |

Table 10. *Results of applying to the Apriori Association Rule Learning algorithm to the predicted labels of all WSI #1, displaying only the most significant rule for each unique consequent label where such a rule exists.*

| Antecedent Itemsets | ⇒ | Consequent Labels | Confidence |
|---|---|---|---|
| {A.W, M.M, E.M.S} | ⇒ | C.D.I | 1 |
| {H.E, E.M.O} | ⇒ | C.L | 0.5540 |
| {C.L, G.O} | ⇒ | E.M.O | 0.5382 |
| {C.D.I, H.E, A.W, M.M, T} | ⇒ | E.M.S | 0.8811 |
| {C.D.I, C.L, G.O, E.M.S} | ⇒ | E.M.U | 0.8795 |
| {C.D.I, E.M.O} | ⇒ | E.T.U | 0.6776 |
| {E.M.O} | ⇒ | G.O | 1 |
| {H.Y, E.M.S, C.L} | ⇒ | H.E | 1 |
| {E.M.S, C.D.I, H.E} | ⇒ | M.M | 0.9421 |
| {E.M.S} | ⇒ | T | 1 |

Table 11. *Results of applying to the Apriori Association Rule Learning algorithm to the predicted labels of all WSI #2, displaying only the most significant rule for each unique consequent label where such a rule exists.*

| Antecedent Itemsets | ⇒ | Consequent Labels | Confidence |
|---|---|---|---|
| {A.W, M.M} | ⇒ | C.D.I | 1 |
| {H.K, E.M.S, E.M.O} | ⇒ | C.L | 0.9487 |
| {E.T.U, C.L, H.K} | ⇒ | E.M.O | 0.9429 |
| {E.M.U, H.K,T} | ⇒ | E.M.S | 1 |
| {C.D.I, C.L, H.Y, M.M, G.O, E.M.S} | ⇒ | E.M.U | 0.6061 |
| {C.D.I, G.O} | ⇒ | E.T.U | 0.5994 |
| {E.M.O} | ⇒ | G.O | 1 |
| {H.K} | ⇒ | H.E | 1 |
| {E.M.U, H.K} | ⇒ | H.Y | 1 |
| {T, E.M.U, E.T.U, C.D.I, C.L} | ⇒ | M.M | 0.8478 |
| {E.M.S} | ⇒ | T | 1 |

Table 12. *Results of applying to the Apriori Association Rule Learning algorithm to the predicted labels of all WSI #3, displaying only the most significant rule for each unique consequent label where such a rule exists.*

| Antecedent Itemsets | ⇒ | Consequent Labels | Confidence |
|---|---|---|---|
| {C.D.I} | ⇒ | A.W | 0.5402 |
| {A.W, T} | ⇒ | C.D.I | 0.9861 |
| {H.K, G.O} | ⇒ | C.L | 0.9485 |
| {H.K, G.O, E.M.S} | ⇒ | E.M.O | 0.8933 |
| {H.E, H.Y, A.W} | ⇒ | E.M.S | 0.9811 |
| {C.D.I, M.M, G.O} | ⇒ | E.M.U | 0.5000 |
| {C.D.I, H.E, E.M.O} | ⇒ | E.T.U | 0.6000 |
| {E.M.U, E.M.O} | ⇒ | G.O | 1 |
| {H.K} | ⇒ | H.E | 1 |
| {H.Y, M.M, T, E.M.O} | ⇒ | H.K | 0.5926 |
| {E.M.S, C.D.I, H.K} | ⇒ | H.Y | 0.9844 |
| {E.M.S, C.D.I, H.K} | ⇒ | M.M | 0.9219 |
| {E.M.S} | ⇒ | T | 1 |

Table 13. *Results of applying to the Apriori Association Rule Learning algorithm to the predicted labels of all WSI # 4, displaying only the most significant rule for each unique consequent label where such a rule exists.*

| Antecedent Itemsets | ⇒ | Consequent Labels | Confidence |
|---|---|---|---|
| {H.K, H.Y, M.M} | ⇒ | C.L | 0.9303 |
| {C.D.I, H.E, H.K, M.M, T} | ⇒ | E.M.S | 0.8596 |
| {H.K, E.M.S} | ⇒ | H.E | 0.9932 |
| {H.Y, M.M, C.L, H.E} | ⇒ | H.K | 0.7462 |
| {E.M.S, C.L, H.K} | ⇒ | H.Y | 0.9790 |
| {E.M.S, C.D.I} | ⇒ | M.M | 0.8780 |
| {E.M.S} | ⇒ | T | 1 |

Table 14. *Results of applying to the Apriori Association Rule Learning algorithm to the predicted labels of all WSI #5, displaying only the most significant rule for each unique consequent label where such a rule exists.*

| Antecedent Itemsets | ⇒ | Consequent Labels | Confidence |
|---|---|---|---|
| {H.Y, A.W, T} | ⇒ | C.D.I | 0.9796 |
| {E.M.S, E.M.O} | ⇒ | C.L | 1 |
| {A.W, M.M} | ⇒ | E.M.S | 0.9643 |
| {C.D.I, C.L, G.O} | ⇒ | E.M.U | 0.9231 |
| {H.E, H.K, M.M, E.M.U} | ⇒ | E.T.U | 0.6923 |
| {E.M.O} | ⇒ | G.O | 1 |
| {E.M.S, E.M.U} | ⇒ | H.E | 1 |
| {M.M, E.M.U, E.T.U} | ⇒ | H.K | 0.9231 |
| {E.M.S, H.K} | ⇒ | H.Y | 1 |
| {E.M.S, C.D.I} | ⇒ | M.M | 0.8223 |
| {E.M.S} | ⇒ | T | 1 |

Table 15. *Results of applying to the Apriori Association Rule Learning algorithm to the predicted labels of all WSI #6, displaying only the most significant rule for each unique consequent label where such a rule exists.*

| Antecedent Itemsets | ⇒ | Consequent Labels | Confidence |
|---|---|---|---|
| {C.D.I } | ⇒ | A.W | 0.5163 |
| {H.Y,A.W,E.M.S } | ⇒ | C.D.I | 1 |
| {H.Y,E.M.O,E.T.U } | ⇒ | C.L | 0.9655 |
| {C.L,H.Y,G.O } | ⇒ | E.M.O | 0.5437 |
| {H.E,A.W,M.M,T } | ⇒ | E.M.S | 0.9617 |
| {H.Y,M.M,G.O,E.M.S } | ⇒ | E.M.U | 0.8571 |
| {C.D.I,G.O } | ⇒ | E.T.U | 0.5288 |
| {E.M.U,E.M.O } | ⇒ | G.O | 1 |
| {H.Y,E.M.S,E.M.U } | ⇒ | H.E | 1 |
| {E.M.O,E.T.U,C.L,H.E} | ⇒ | H.Y | 0.5625 |
| {T,E.M.U,E.T.U,C.D.I} | ⇒ | M.M | 0.8727 |
| {E.M.S,E.M.U } | ⇒ | T | 1 |