# Progress Report

Daniel Biro

January 11, 2021

## Contents

# 1 Experiments - workflow version: 1.1.0; experiments revision: YYY

## 1.1 Methods

The main script for running experiments (run_experiment.py) has been updated and simplified. It now also includes shuffling the data before each model training (the initial shuffling was removed from building_dataset.py), as well as model testing. The plotting and result aggregation functionality for showcasing means and errors were moved to two separate scripts, produce_plots.py in visualization and aggregate_result.py in utils.

Plotting now happens using Pandas dataframes and Seaborn lineplots for all plots except ROC curves which still use matplotlib. Now each trained model has an ROC curve on validation and test sets, while for the other metrics there is a single plot for their progression over time on the training and validation sets. AUC values for each class are now also recorded and a separate plot shows their progression over time on the validation set.

Model testing is called similarly to model training from the model_testing.py script, this creates the new metrics files which are then aggregated to get a single mean and error on each, as well as the ROC curves for each model.

New Snakemake rules and bash scripts were added for processing the genomic regions of the 20th chromosome as the test set, which this new set containing 5264 germline variants and 2269 somatic variants along with 7533 randomly sampled healthy positions (normal class). This test set has around 600 datapoints more than the training and validations sets combined.

All Sphinx documentation has been updated and the new code uploaded to GitHub.

## 1.2 Results

Runtimes have been measured for each model categories and they are longer than the estimate I have given earlier. Everything had to be re-run for shuffling, new plots, extra metrics and not all of it has finished by the time of writing this report. The previously most promising models have been prioritized.

The new approximate runtimes per model trained are as follows: basic model types: 2hrs, bi-directional models: 4hrs, Transformer network (2 encoder layers): 24.5hrs. Adding dropout did not significantly increase runtime.

Below the new validation and test set results of the previously and newly best peforming models are presented. There are two caveats with the Transformer model: the previous results were of a model with 1 encoder layer, this new model has 2 layers, with the new model performing worse. This is accidental as the scripts were usually started up for overnight training, but it was planned to compare how performance changes with an additional layer for later on. The 1 layer model type will be trained and evaluated as well. The new model's results are also only 4 runs due to the long training time. Observing some of the individual model results there is room for improvement with more models being trained (and for the 1 layer model type as well).

| Model | Accuracy | Precision | Recall | F1 Score |
|---|---|---|---|---|
| GRU | $94.51\% \pm 0.79$ | $90.95\% \pm 3.79$ | $89.14\% \pm 2.88$ | $90.01\% \pm 3.35$ |
| Bi-GRU + 5% Dropout | $95.62\% \pm 0.06$ | $95.26\% \pm 0.07$ | $92.56\% \pm 0.1$ | $93.89\% \pm 0.07$ |
| Bi-GRU + 10% Dropout | $95.46\% \pm 0.09$ | $95.03\% \pm 0.14$ | $92.32\% \pm 0.2$ | $93.65\% \pm 0.16$ |
| Transformer + 10% Dropout | $90.74\% \pm 0.07$ | $89.41\% \pm 0.05$ | $85.09\% \pm 0.14$ | $87.2\% \pm 0.09$ |
| **Perceptron** | $\mathbf{96.24}\% \pm \mathbf{0.01}$ | $\mathbf{96.3}\% \pm \mathbf{0.01}$ | $\mathbf{93.72}\% \pm \mathbf{0.03}$ | $\mathbf{94.99}\% \pm \mathbf{0.02}$ |

Table 1: Performance of best models on the validation set

| Model | Normal | Germline Variant | Somatic Variant |
|---|---|---|---|
| GRU | $0.929\% \pm 0.02$ | $0.9308\% \pm 0.01$ | $0.904\% \pm 0.05$ |
| Bi-GRU + 5% Dropout | $0.945\% \pm 0.0$ | $0.9331\% \pm 0.0$ | $0.9591\% \pm 0.0$ |
| Bi-GRU + 10% Dropout | $0.9495\% \pm 0.0$ | $0.9339\% \pm 0.0$ | $0.96\% \pm 0.0$ |
| Transformer + 10% Dropout | $0.8288\% \pm 0.0$ | $0.8208\% \pm 0.01$ | $0.8674\% \pm 0.0$ |
| **Perceptron** | $\mathbf{0.9575}\% \pm \mathbf{0.0}$ | $\mathbf{0.9706}\% \pm \mathbf{0.0}$ | $\mathbf{0.9591}\% \pm \mathbf{0.0}$ |

Table 2: Class AUCs of best models on the validation set

| Model | Accuracy | Precision | Recall | F1 Score |
|---|---|---|---|---|
| GRU | $93.16\% \pm 0.78$ | $89.88\% \pm 3.33$ | $85.7\% \pm 2.61$ | $87.73\% \pm 2.96$ |
| **Bi-GRU + 5% Dropout** | $\mathbf{94.18}\% \pm \mathbf{0.03}$ | $\mathbf{93.91}\% \pm \mathbf{0.05}$ | $\mathbf{88.67}\% \pm \mathbf{0.1}$ | $\mathbf{91.21}\% \pm \mathbf{0.07}$ |
| Bi-GRU + 10% Dropout | $94.11\% \pm 0.04$ | $93.67\% \pm 0.07$ | $88.56\% \pm 0.12$ | $91.04\% \pm 0.09$ |
| Transformer + 10% Dropout | $79.93\% \pm 10.02$ | $70.44\% \pm 12.51$ | $67.33\% \pm 11.96$ | $68.78\% \pm 12.22$ |
| Perceptron | $94.15\% \pm 0.01$ | $94.15\% \pm 0.01$ | $88.41\% \pm 0.01$ | $91.19\% \pm 0.01$ |

Table 3: Performance of best models on the test set

| Model | Normal | Germline Variant | Somatic Variant |
|---|---|---|---|
| GRU | $0.9054\% \pm 0.02$ | $0.9197\% \pm 0.0$ | $0.8808\% \pm 0.04$ |
| Bi-GRU + 5% Dropout | $0.9172\% \pm 0.0$ | $0.9235\% \pm 0.0$ | $0.9266\% \pm 0.0$ |
| Bi-GRU + 10% Dropout | $0.9246\% \pm 0.0$ | $0.9237\% \pm 0.0$ | $0.9273\% \pm 0.0$ |
| Transformer + 10% Dropout | $0.7371\% \pm 0.13$ | $0.7176\% \pm 0.14$ | $0.661\% \pm 0.16$ |
| **Perceptron** | $\mathbf{0.928}\% \pm \mathbf{0.0}$ | $\mathbf{0.9375}\% \pm \mathbf{0.0}$ | $\mathbf{0.9224}\% \pm \mathbf{0.0}$ |

Table 4: Class AUCs of best models on the test set

(a)

(b)

(c)

(d)

(e)

(f)

Figure 1: Figures for GRU model

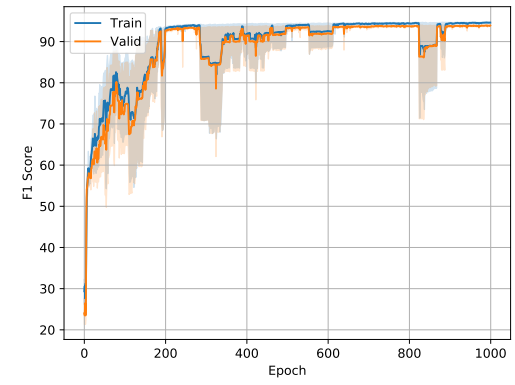Figure 2: Figures for Bi-GRU + 5% Dropout model

(a)

(b)
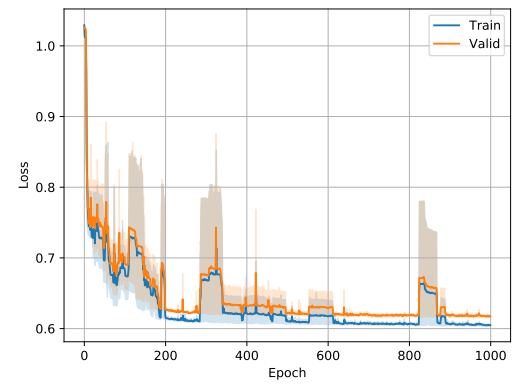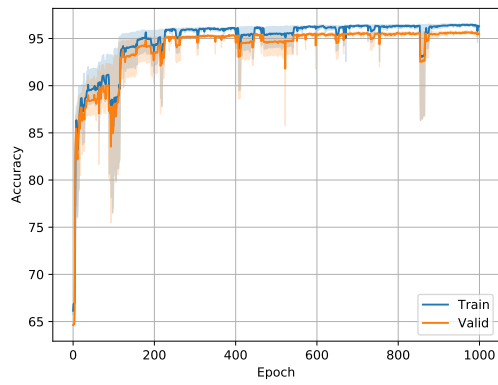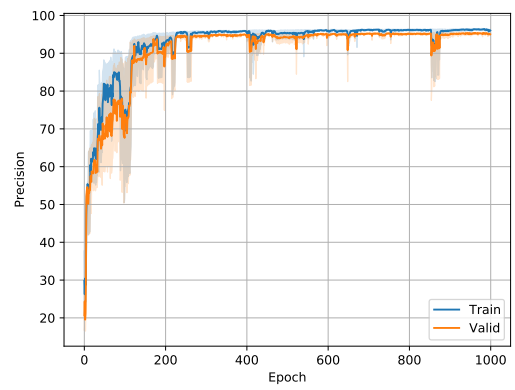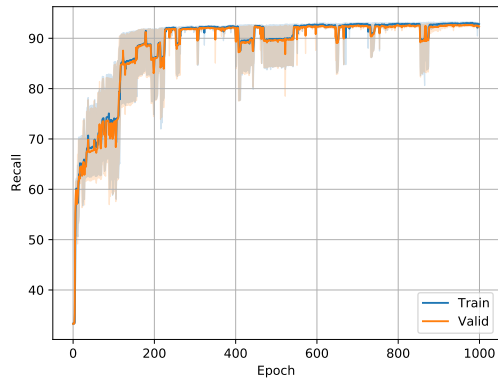
(c)

(d)

(e)

(f)

Figure 3: Figures for Bi-GRU + 10% Dropout model

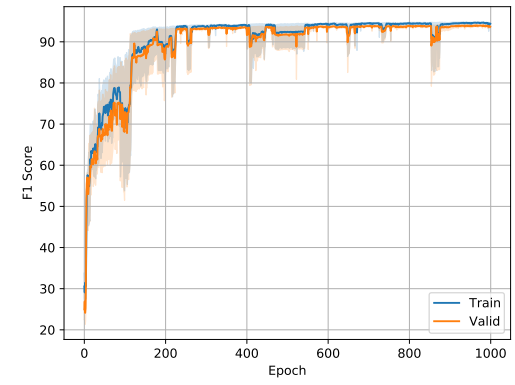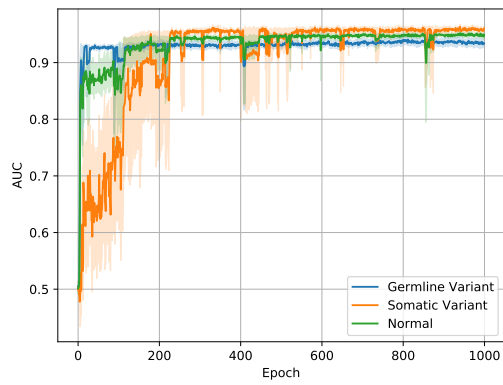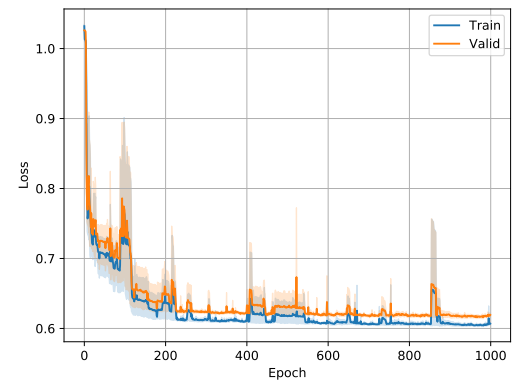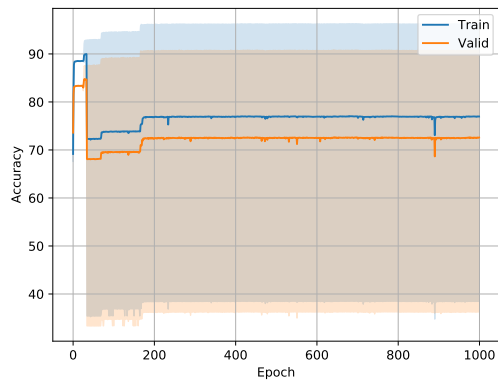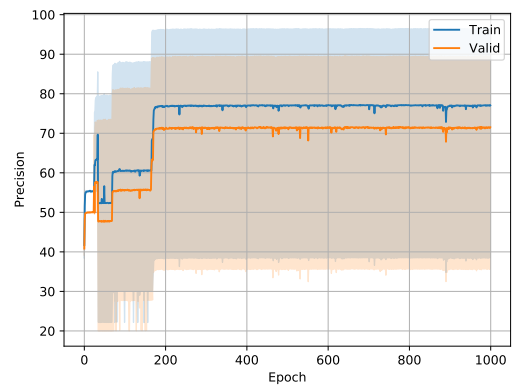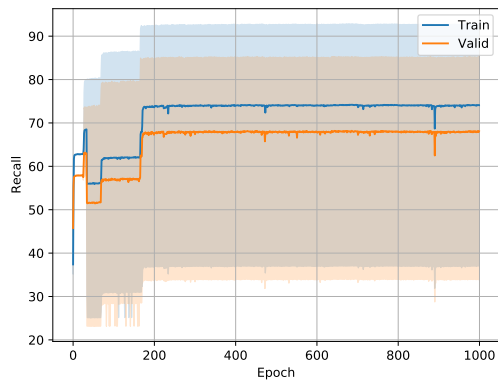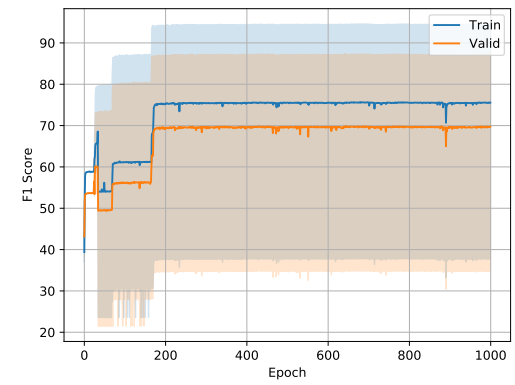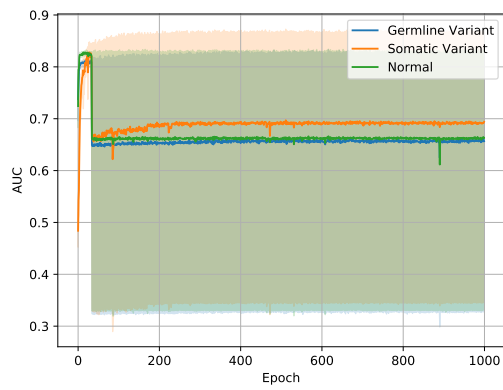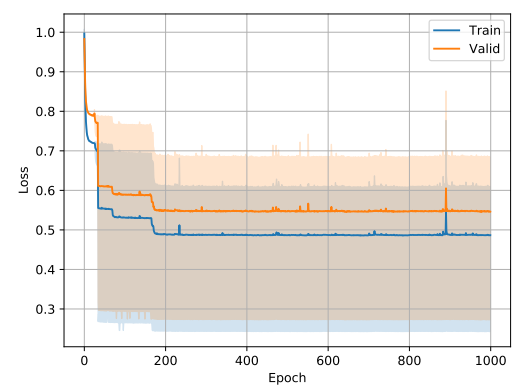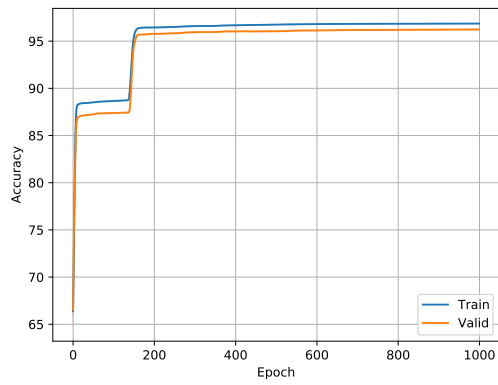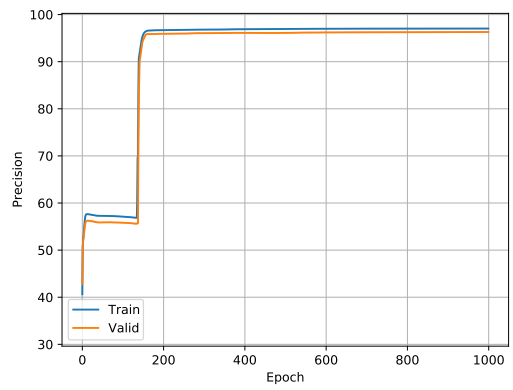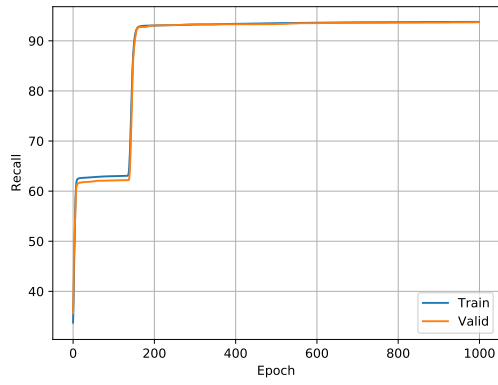(a)

(b)

(c)

(d)
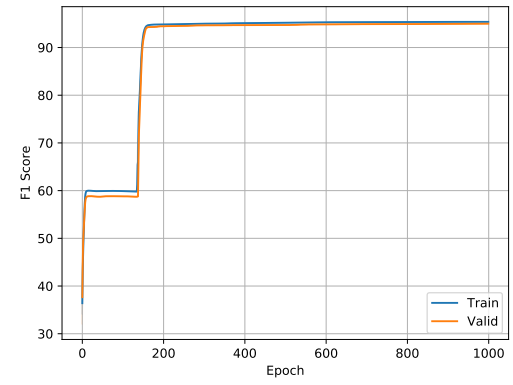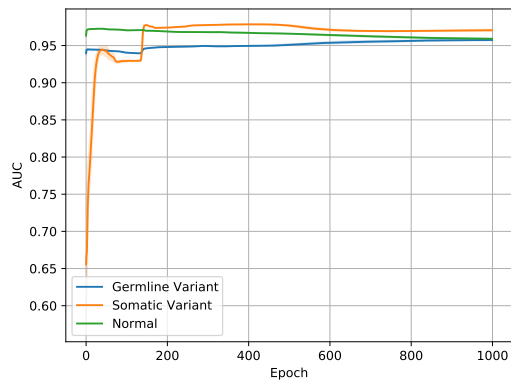
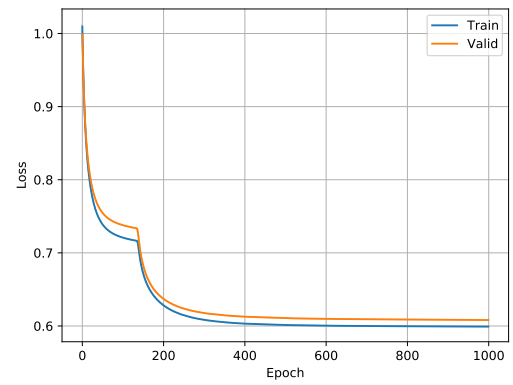(e)

(f)

Figure 4: Figures for Transformer model

(a)

(b)

(c)

(d)

(e)

(f)

Figure 5: Figures for Perceptron Dropout model