# SUPPLEMENTAL MATERIAL

## Dataset Statistics
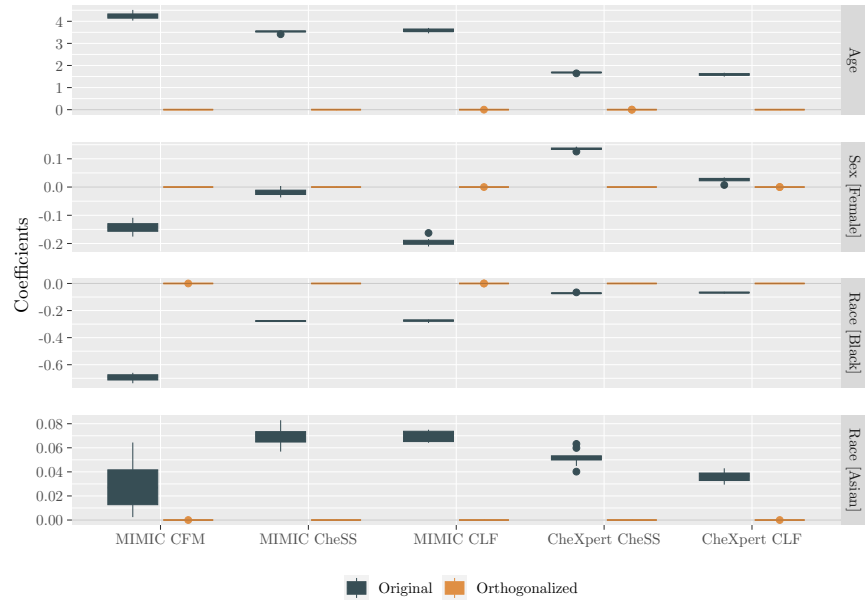
**MIMIC**

Training set

| | All | White | Black | Asian | Male | Female |
|---|---|---|---|---|---|---|
| Patients | 42,148 | 31,936 (75.77%) | 8,398 (19.93%) | 1,814 (4.30%) | 20,123 (47.74%) | 22,025 (52.26%) |
| Scans | 181,342 | 140,445 (77.45%) | 33,906 (18,70%) | 6,991 (3.86%) | 97,361 (53.69%) | 83961 (46.21%) |
| Age | $62.6 \pm 16.6$ | $63.9 \pm 16.3$ | $57.7 \pm 16.7$ | $62.1 \pm 17.8$ | $62.32 \pm 15.8$ | $63.0 \pm 17.5$ |

Test set

| | All | White | Black | Asian | Male | Female |
|---|---|---|---|---|---|---|
| Patients | 257 | 205 (79.77%) | 45 (17.51%) | 7 (2.72%) | 141 (54.86%) | 116 (45.14%) |
| Scans | 3,041 | 2,235 (73.50%) | 676 (22.22%) | 130 (4.27%) | 1,658 (54.52%) | 1,383 (45.48%) |
| Age | $65.8 \pm 12.1$ | $66.2 \pm 12.3$ | $64.1 \pm 11.9$ | $67.4 \pm 9.5$ | $66.0 \pm 11.6$ | $65.4 \pm 12.8$ |

**CheXpert**

Training set

| | All | White | Black | Asian | Male | Female |
|---|---|---|---|---|---|---|
| Patients | 25,730 | 20,034 (77.86%) | 1,751 (6.81%)) | 3,945 (15.33%) | 14,165 (55.05%) | 11,565 (44.95%) |
| Scans | 76,205 | 59,238 (77.73%) | 5,596 (7.34%) | 11,371 (14.92%) | 44,774 (58.75%) | 31,431 (41.25%) |
| Age | $63.1 \pm 17.4$ | $64.3 \pm 17.2$ | $55.7 \pm 17.4$ | $61.6 \pm 17.4$ | $62.5 \pm 17.0$ | $63.8 \pm 17.9$ |

Test set

| | All | White | Black | Asian | Male | Female |
|---|---|---|---|---|---|---|
| Patients | 12,866 | 9,956 (77.38%) | 879 (6.83%) | 2,031 (15.79%) | 7,091 (55.11%) | 5,775 (44.89%) |
| Scans | 38,240 | 29,844 (78.04%) | 2746 (7.18%) | 5,650 (14.278%) | 22,265 (58.22%) | 15,975 (41.78%) |
| Age | $63.3 \pm 17.2$ | $64.2 \pm 17.1$ | $57.4 \pm 16.3$ | $61.1 \pm 17.6$ | $62.8 \pm 16.4$ | $63.9 \pm 18.3$ |

Table T.1: Statistics of the utilized MIMIC and CheXpert subsets per split and subgroups.
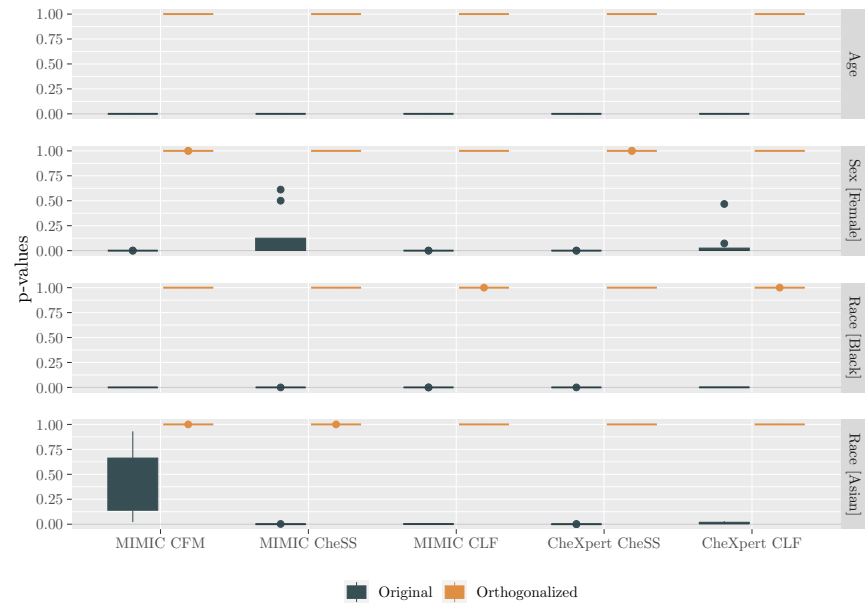
# Influence of Protected Features on Model Prediction

The following figures are an extension of the results presented in Section  and visualize Table 1. This includes the distribution of coefficients and p-values over ten randomly initialized runs obtained from the *evaluation model*. The results are shown for the three exemplary labels *Pleural Effusion*, *Cardiomegaly*, and *No Finding*.
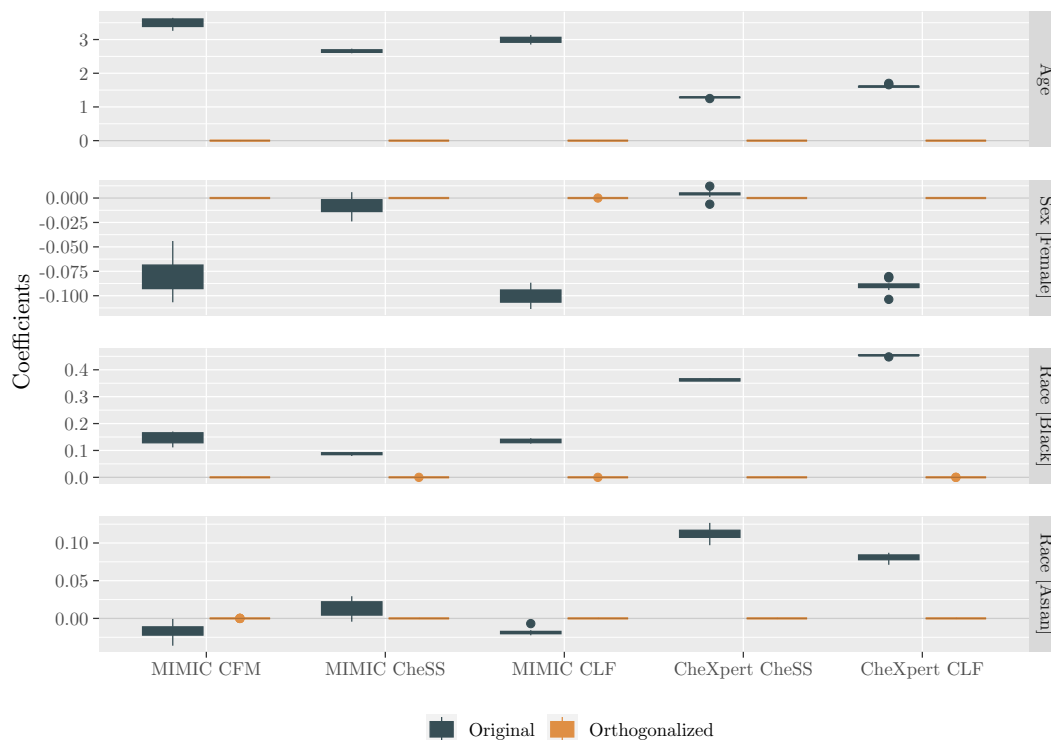
## Pathology: Pleural Effusion



(a) Coefficients for *Pleural Effusion*.



(b) p-values associated with the respective coefficients.

Figure F.1: Distribution of derived coefficients and p-values for 10 downstream models per embedding and protected feature category on the label *Pleural Effusion*.

# Pathology: Cardiomegaly



(a) Coefficients for *Cardiomegaly*.



(b) p-values associated with the respective coefficients.

Figure F.2: Distribution of derived coefficients and p-values for 10 downstream models per embedding and protected feature category on the label *Cardiomegaly*.

**Pathology: No Finding**
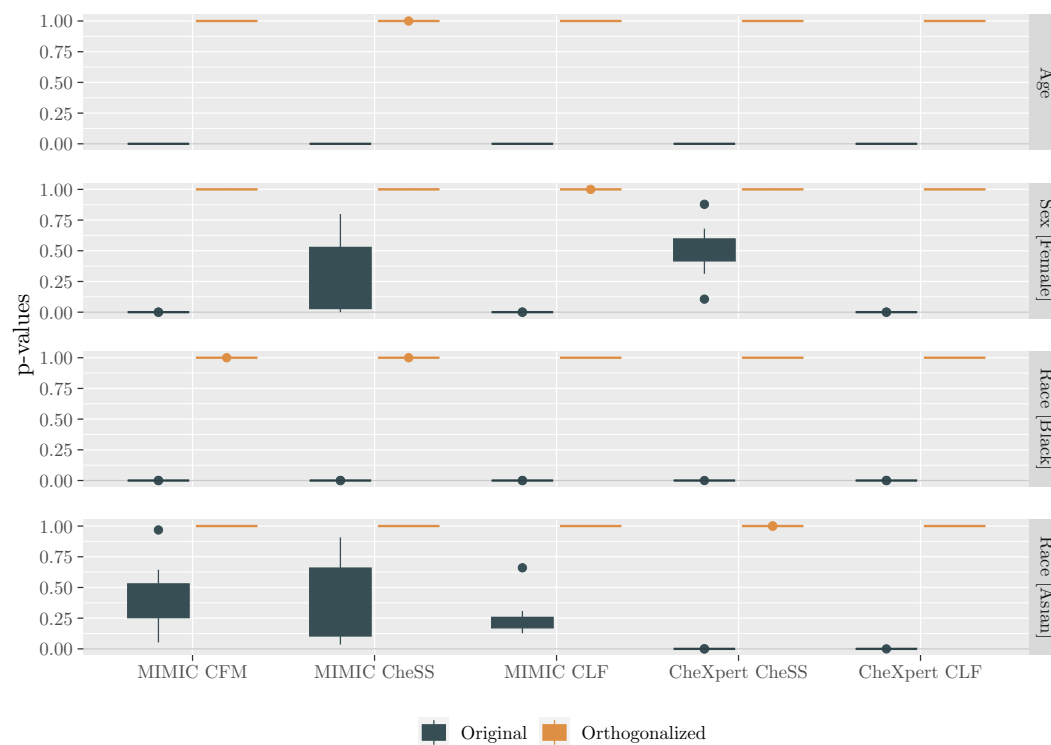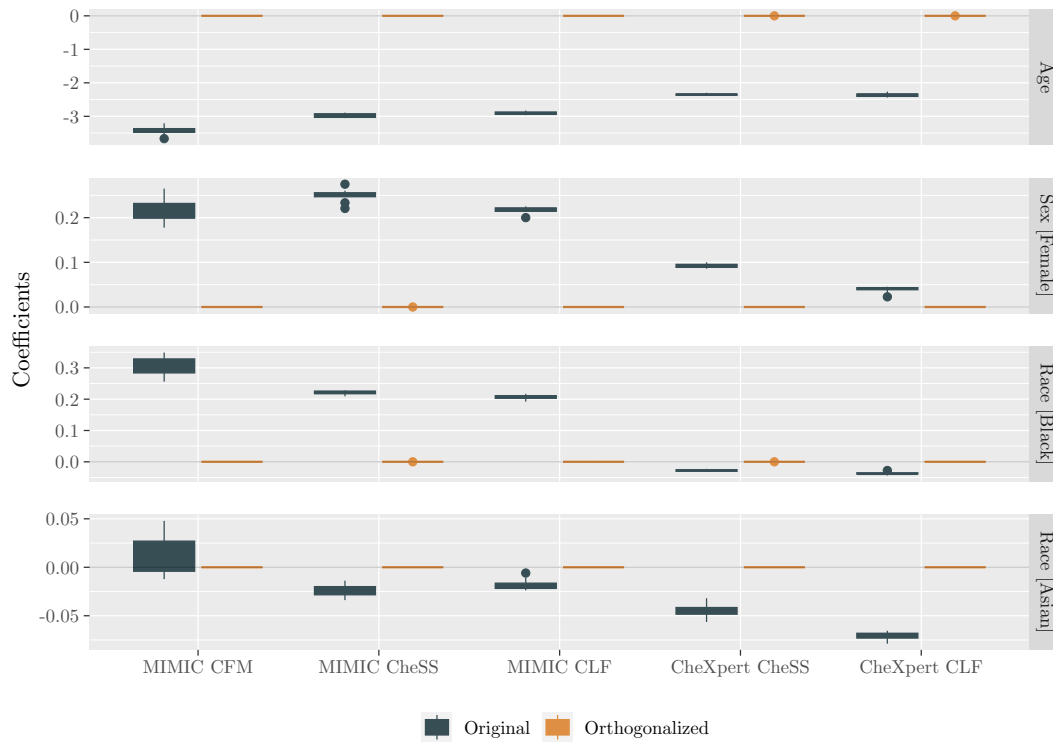


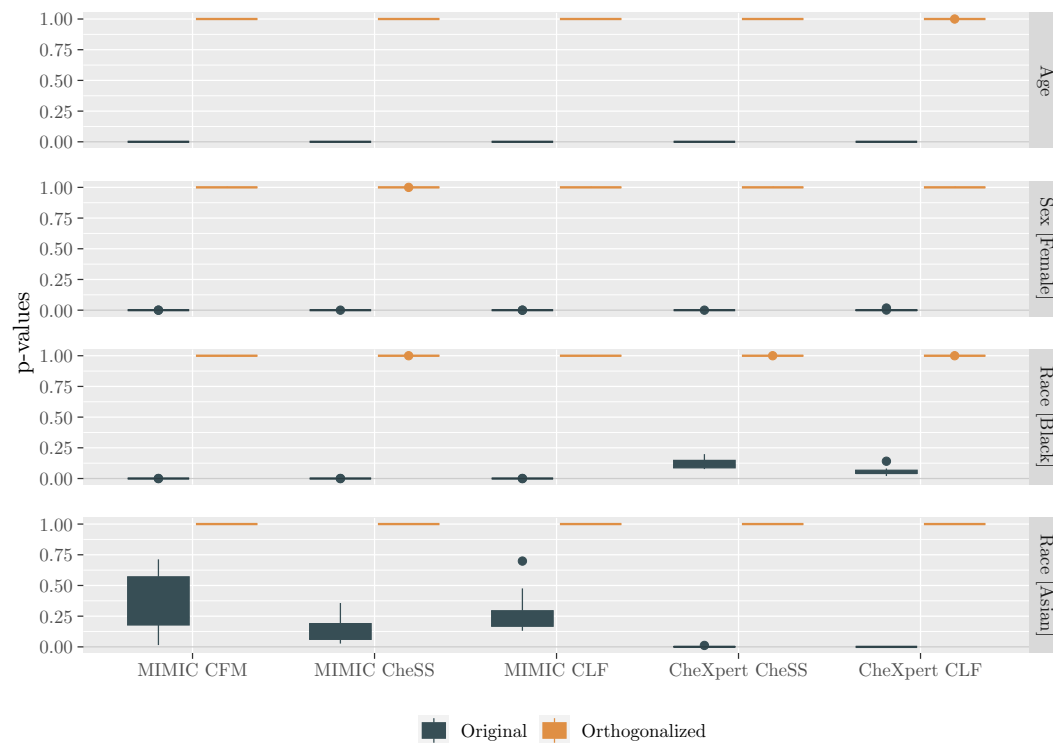(a) Coefficients for *No Finding*.



(b) p-values associated with the respective coefficients.

Figure F.3: Distribution of derived coefficients and p-values for 10 downstream models per embedding and protected feature category on the label *No Finding*.

# Predicting Protected Information

| | Emb. | Orthogonalized? | Age MAE | Age $R^2$ | Sex AUC | Sex Sens. | Sex Spec. | Race [White] AUC | Race [White] Sens. | Race [White] Spec. | Race [Black] AUC | Race [Black] Sens. | Race [Black] Spec. | Race [Asian] AUC | Race [Asian] Sens. | Race [Asian] Spec. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MIMIC | CFM | ✗ | 8.040 ±1.456 | 0.314 ±0.241 | 0.979 ±0.000 | 0.933 ±0.025 | 0.902 ±0.039 | 0.848 ±0.004 | 0.940 ±0.016 | 0.497 ±0.049 | 0.871 ±0.004 | 0.521 ±0.055 | 0.944 ±0.015 | 0.870 ±0.005 | 0.179 ±0.028 | 0.991 ±0.005 |
| | | ✓ | 10.969 ±0.433 | -0.294 ±0.087 | 0.507 ±0.031 | 0.775 ±0.077 | 0.224 ±0.100 | 0.501 ±0.045 | 1.000 ±0.000 | 0.000 ±0.000 | 0.509 ±0.042 | 0.000 ±0.000 | 1.000 ±0.000 | 0.464 ±0.076 | 0.000 ±0.000 | 1.000 ±0.000 |
| | CheSS | ✗ | 7.893 ±0.066 | 0.331 ±0.010 | 0.945 ±0.001 | 0.907 ±0.011 | 0.820 ±0.016 | 0.761 ±0.004 | 0.975 ±0.008 | 0.158 ±0.039 | 0.768 ±0.003 | 0.143 ±0.039 | 0.968 ±0.010 | 0.733 ±0.004 | 0.013 ±0.008 | 0.997 ±0.002 |
| | | ✓ | 9.908 ±0.082 | -0.083 ±0.016 | 0.482 ±0.037 | 0.982 ±0.024 | 0.013 ±0.025 | 0.489 ±0.060 | 1.000 ±0.000 | 0.000 ±0.000 | 0.478 ±0.068 | 0.000 ±0.000 | 1.000 ±0.000 | 0.514 ±0.048 | 0.000 ±0.000 | 1.000 ±0.000 |
| | CLF | ✗ | 8.816 ±0.063 | 0.161 ±0.011 | 0.832 ±0.000 | 0.804 ±0.015 | 0.702 ±0.016 | 0.631 ±0.005 | 0.988 ±0.004 | 0.048 ±0.012 | 0.652 ±0.005 | 0.050 ±0.012 | 0.987 ±0.005 | 0.663 ±0.007 | 0.000 ±0.000 | 1.000 ±0.000 |
| | | ✓ | 9.941 ±0.095 | -0.093 ±0.025 | 0.456 ±0.039 | 0.994 ±0.007 | 0.007 ±0.006 | 0.500 ±0.048 | 1.000 ±0.000 | 0.000 ±0.000 | 0.498 ±0.056 | 0.000 ±0.000 | 1.000 ±0.000 | 0.493 ±0.084 | 0.000 ±0.000 | 1.000 ±0.000 |
| CheXpert | CheSS | ✗ | 9.333 ±0.103 | 0.529 ±0.009 | 0.946 ±0.000 | 0.912 ±0.008 | 0.821 ±0.018 | 0.780 ±0.001 | 0.984 ±0.007 | 0.136 ±0.038 | 0.762 ±0.001 | 0.018 ±0.006 | 0.999 ±0.001 | 0.816 ±0.001 | 0.173 ±0.047 | 0.983 ±0.007 |
| | | ✓ | 13.875 ±0.035 | -0.009 ±0.005 | 0.499 ±0.010 | 1.000 ±0.000 | 0.000 ±0.000 | 0.501 ±0.007 | 1.000 ±0.000 | 0.000 ±0.000 | 0.498 ±0.015 | 0.000 ±0.000 | 1.000 ±0.000 | 0.500 ±0.007 | 0.000 ±0.000 | 1.000 ±0.000 |
| | CLF | ✗ | 10.693 ±0.047 | 0.385 ±0.005 | 0.868 ±0.000 | 0.854 ±0.008 | 0.696 ±0.014 | 0.721 ±0.001 | 0.993 ±0.003 | 0.040 ±0.012 | 0.688 ±0.002 | 0.003 ±0.002 | 1.000 ±0.000 | 0.764 ±0.000 | 0.049 ±0.016 | 0.993 ±0.003 |
| | | ✓ | 13.866 ±0.035 | -0.009 ±0.005 | 0.496 ±0.010 | 1.000 ±0.000 | 0.000 ±0.000 | 0.501 ±0.008 | 1.000 ±0.000 | 0.000 ±0.000 | 0.497 ±0.021 | 0.000 ±0.000 | 1.000 ±0.000 | 0.503 ±0.009 | 0.000 ±0.000 | 1.000 ±0.000 |

Table T.2: Regression/Classification performance for deriving protected features from an embedding vector with mean and standard deviation over 10 randomly initialized runs. The displayed metrics include mean absolute error (MAE), $R^2$ for age regression as well as AUC, sensitivity (sens.), and specificity (spec.) for classification.

| Dataset: | MIMIC CFM ✗ | CFM ✓ | CFM Δ | CheSS ✗ | CheSS ✓ | CheSS Δ | CLF ✗ | CLF ✓ | CLF Δ | CheXpert CheSS ✗ | CheSS ✓ | CheSS Δ | CLF ✗ | CLF ✓ | CLF Δ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Enl. Cardiomed. | 0.728 ± 0.009 | 0.721 ± 0.018 | -0.97 % | 0.636 ± 0.003 | 0.643 ± 0.007 | +1.09 % | 0.601 ± 0.002 | 0.593 ± 0.004 | -1.35 % | 0.636 ± 0.003 | 0.621 ± 0.001 | +2.36 % | 0.634 ± 0.000 | 0.639 ± 0.001 | -0.78 % |
| Cardiomegaly | 0.780 ± 0.002 | 0.775 ± 0.003 | -0.65 % | 0.750 ± 0.001 | 0.751 ± 0.001 | +0.13 % | 0.737 ± 0.001 | 0.736 ± 0.001 | -0.14 % | 0.789 ± 0.000 | 0.791 ± 0.001 | +0.25 % | 0.799 ± 0.000 | 0.793 ± 0.000 | -0.76 % |
| Lung Opacity | 0.696 ± 0.003 | 0.684 ± 0.005 | -1.75 % | 0.626 ± 0.002 | 0.627 ± 0.004 | +0.16 % | 0.623 ± 0.001 | 0.612 ± 0.002 | -1.80 % | 0.684 ± 0.000 | 0.685 ± 0.000 | -0.15 % | 0.695 ± 0.000 | 0.690 ± 0.000 | -0.72 % |
| Lung Lesion | 0.731 ± 0.006 | 0.718 ± 0.007 | -1.81 % | 0.623 ± 0.003 | 0.630 ± 0.004 | +1.11 % | 0.576 ± 0.009 | 0.591 ± 0.014 | +2.54 % | 0.700 ± 0.002 | 0.672 ± 0.002 | +4.00 % | 0.701 ± 0.001 | 0.707 ± 0.002 | -0.85 % |
| Edema | 0.843 ± 0.001 | 0.837 ± 0.002 | -0.72 % | 0.804 ± 0.000 | 0.798 ± 0.000 | -0.75 % | 0.791 ± 0.000 | 0.783 ± 0.001 | -1.02 % | 0.789 ± 0.000 | 0.791 ± 0.000 | -0.25 % | 0.788 ± 0.000 | 0.783 ± 0.000 | -0.64 % |
| Consolidation | 0.748 ± 0.008 | 0.742 ± 0.009 | -0.81 % | 0.648 ± 0.001 | 0.650 ± 0.005 | +0.31 % | 0.638 ± 0.003 | 0.640 ± 0.002 | +0.31 % | 0.683 ± 0.001 | 0.669 ± 0.001 | +2.05 % | 0.689 ± 0.001 | 0.692 ± 0.002 | +0.43 % |
| Pneumonia | 0.703 ± 0.005 | 0.704 ± 0.004 | +0.14 % | 0.586 ± 0.005 | 0.597 ± 0.009 | +1.84 % | 0.589 ± 0.003 | 0.607 ± 0.002 | +2.97 % | 0.652 ± 0.003 | 0.610 ± 0.002 | +6.44 % | 0.652 ± 0.001 | 0.659 ± 0.004 | +1.06 % |
| Atelectasis | 0.746 ± 0.002 | 0.734 ± 0.004 | -1.63 % | 0.702 ± 0.000 | 0.696 ± 0.001 | -0.86 % | 0.685 ± 0.001 | 0.671 ± 0.001 | -2.09 % | 0.636 ± 0.001 | 0.631 ± 0.001 | -0.79 % | 0.632 ± 0.000 | 0.633 ± 0.001 | +0.16 % |
| Pneumothorax | 0.843 ± 0.005 | 0.830 ± 0.007 | -1.57 % | 0.649 ± 0.002 | 0.645 ± 0.005 | -0.62 % | 0.634 ± 0.003 | 0.638 ± 0.004 | +0.63 % | 0.749 ± 0.001 | 0.732 ± 0.001 | +2.27 % | 0.730 ± 0.001 | 0.740 ± 0.001 | +1.35 % |
| Pleural Effusion | 0.870 ± 0.001 | 0.859 ± 0.002 | -1.28 % | 0.802 ± 0.000 | 0.792 ± 0.001 | -1.26 % | 0.797 ± 0.000 | 0.781 ± 0.000 | -2.05 % | 0.798 ± 0.000 | 0.792 ± 0.000 | +0.75 % | 0.804 ± 0.000 | 0.801 ± 0.000 | -0.37 % |
| Pleural Other | 0.894 ± 0.009 | 0.874 ± 0.021 | -2.29 % | 0.711 ± 0.005 | 0.746 ± 0.011 | +4.69 % | 0.684 ± 0.006 | 0.693 ± 0.009 | +1.30 % | 0.756 ± 0.004 | 0.723 ± 0.001 | +4.37 % | 0.718 ± 0.004 | 0.718 ± 0.003 | +0.00 % |
| Fracture | 0.752 ± 0.007 | 0.739 ± 0.013 | -1.76 % | 0.643 ± 0.005 | 0.648 ± 0.009 | +0.77 % | 0.642 ± 0.004 | 0.652 ± 0.012 | +1.53 % | 0.682 ± 0.003 | 0.668 ± 0.003 | +2.05 % | 0.667 ± 0.001 | 0.673 ± 0.003 | +0.89 % |
| Support Devices | 0.909 ± 0.001 | 0.905 ± 0.001 | -0.44 % | 0.801 ± 0.000 | 0.801 ± 0.001 | +0.00 % | 0.767 ± 0.001 | 0.763 ± 0.001 | -0.52 % | 0.748 ± 0.000 | 0.731 ± 0.000 | +2.27 % | 0.711 ± 0.000 | 0.721 ± 0.000 | +1.39 % |
| No Finding | 0.801 ± 0.003 | 0.770 ± 0.005 | -4.03 % | 0.746 ± 0.001 | 0.725 ± 0.002 | -2.90 % | 0.747 ± 0.000 | 0.728 ± 0.001 | -2.61 % | 0.824 ± 0.001 | 0.833 ± 0.001 | -1.09 % | 0.854 ± 0.000 | 0.844 ± 0.000 | -1.18 % |
| Total | 0.789 ± 0.005 | 0.778 ± 0.009 | -1.39 % | 0.695 ± 0.003 | 0.696 ± 0.005 | -0.14 % | 0.679 ± 0.003 | 0.678 ± 0.006 | -0.14 % | 0.723 ± 0.002 | 0.710 ± 0.001 | +1.83 % | 0.720 ± 0.001 | 0.721 ± 0.002 | +0.13 % |

Table T.3: Prediction performance original versus orthogonalized data on the MIMIC and CheXpert datasets. The table shows the mean and standard deviation of the AUC over 10 randomly initialized runs. Additionally, $\Delta$ depicts the percentual change from the original to the corrected embedding AUC.

# Downstream Prediction Performance

The following table provides additional metrics for the the labels *Pleural Effusion*, *Cardiomegaly* and *No Finding* and supplements and Table T.3.

| Pathology: | | Pleural Effusion | | Cardiomegaly | | No Finding | |
|---|---|---|---|---|---|---|---|
| Ortho.: | | ✗ | ✓ | ✗ | ✓ | ✗ | ✓ |
| **MIMIC CFM** | AUC | $0.870 \pm 0.001$ | $0.856 \pm 0.002$ | $0.780 \pm 0.001$ | $0.767 \pm 0.003$ | $0.801 \pm 0.003$ | $0.786 \pm 0.005$ |
| | Acc. | $0.804 \pm 0.003$ | $0.784 \pm 0.003$ | $0.753 \pm 0.004$ | $0.751 \pm 0.002$ | $0.839 \pm 0.005$ | $0.821 \pm 0.004$ |
| | Sens. | $0.612 \pm 0.037$ | $0.485 \pm 0.011$ | $0.351 \pm 0.083$ | $0.169 \pm 0.015$ | $0.318 \pm 0.063$ | $0.474 \pm 0.018$ |
| | Spec. | $0.897 \pm 0.016$ | $0.929 \pm 0.003$ | $0.897 \pm 0.034$ | $0.961 \pm 0.005$ | $0.953 \pm 0.018$ | $0.896 \pm 0.008$ |
| | Prec. | $0.744 \pm 0.020$ | $0.768 \pm 0.007$ | $0.559 \pm 0.029$ | $0.608 \pm 0.017$ | $0.613 \pm 0.057$ | $0.500 \pm 0.011$ |
| | F1 | $0.670 \pm 0.016$ | $0.594 \pm 0.008$ | $0.423 \pm 0.059$ | $0.264 \pm 0.019$ | $0.411 \pm 0.047$ | $0.486 \pm 0.008$ |
| **MIMIC CheSS** | AUC | $0.802 \pm 0.000$ | $0.792 \pm 0.001$ | $0.750 \pm 0.000$ | $0.742 \pm 0.001$ | $0.747 \pm 0.001$ | $0.737 \pm 0.001$ |
| | Acc. | $0.755 \pm 0.002$ | $0.737 \pm 0.001$ | $0.742 \pm 0.002$ | $0.740 \pm 0.001$ | $0.816 \pm 0.004$ | $0.797 \pm 0.001$ |
| | Sens. | $0.506 \pm 0.046$ | $0.320 \pm 0.007$ | $0.139 \pm 0.048$ | $0.054 \pm 0.009$ | $0.245 \pm 0.023$ | $0.372 \pm 0.007$ |
| | Spec. | $0.876 \pm 0.022$ | $0.939 \pm 0.002$ | $0.960 \pm 0.015$ | $0.987 \pm 0.002$ | $0.941 \pm 0.010$ | $0.889 \pm 0.002$ |
| | Prec. | $0.666 \pm 0.019$ | $0.716 \pm 0.005$ | $0.556 \pm 0.011$ | $0.597 \pm 0.012$ | $0.476 \pm 0.019$ | $0.423 \pm 0.003$ |
| | F1 | $0.573 \pm 0.023$ | $0.442 \pm 0.006$ | $0.218 \pm 0.061$ | $0.099 \pm 0.015$ | $0.322 \pm 0.016$ | $0.396 \pm 0.005$ |
| **MIMIC CLF** | AUC | $0.797 \pm 0.001$ | $0.780 \pm 0.000$ | $0.737 \pm 0.000$ | $0.727 \pm 0.001$ | $0.747 \pm 0.000$ | $0.736 \pm 0.001$ |
| | Acc. | $0.748 \pm 0.001$ | $0.724 \pm 0.002$ | $0.739 \pm 0.001$ | $0.737 \pm 0.000$ | $0.819 \pm 0.004$ | $0.801 \pm 0.003$ |
| | Sens. | $0.508 \pm 0.026$ | $0.299 \pm 0.012$ | $0.065 \pm 0.019$ | $0.011 \pm 0.002$ | $0.322 \pm 0.014$ | $0.421 \pm 0.013$ |
| | Spec. | $0.864 \pm 0.014$ | $0.930 \pm 0.005$ | $0.982 \pm 0.006$ | $0.998 \pm 0.000$ | $0.928 \pm 0.007$ | $0.884 \pm 0.007$ |
| | Prec. | $0.644 \pm 0.012$ | $0.674 \pm 0.008$ | $0.573 \pm 0.022$ | $0.647 \pm 0.063$ | $0.496 \pm 0.016$ | $0.443 \pm 0.008$ |
| | F1 | $0.567 \pm 0.012$ | $0.414 \pm 0.011$ | $0.115 \pm 0.031$ | $0.021 \pm 0.003$ | $0.390 \pm 0.006$ | $0.431 \pm 0.006$ |
| **CheXpert CheSS** | AUC | $0.793 \pm 0.000$ | $0.798 \pm 0.000$ | $0.789 \pm 0.000$ | $0.794 \pm 0.001$ | $0.833 \pm 0.000$ | $0.825 \pm 0.001$ |
| | Acc. | $0.726 \pm 0.001$ | $0.730 \pm 0.000$ | $0.875 \pm 0.000$ | $0.876 \pm 0.000$ | $0.915 \pm 0.000$ | $0.915 \pm 0.000$ |
| | Sens. | $0.616 \pm 0.028$ | $0.606 \pm 0.004$ | $0.053 \pm 0.012$ | $0.077 \pm 0.006$ | $0.099 \pm 0.014$ | $0.064 \pm 0.010$ |
| | Spec. | $0.800 \pm 0.018$ | $0.814 \pm 0.003$ | $0.996 \pm 0.001$ | $0.994 \pm 0.001$ | $0.992 \pm 0.001$ | $0.995 \pm 0.001$ |
| | Prec. | $0.679 \pm 0.010$ | $0.690 \pm 0.002$ | $0.642 \pm 0.027$ | $0.640 \pm 0.019$ | $0.546 \pm 0.012$ | $0.564 \pm 0.016$ |
| | F1 | $0.645 \pm 0.011$ | $0.645 \pm 0.002$ | $0.097 \pm 0.020$ | $0.138 \pm 0.009$ | $0.167 \pm 0.020$ | $0.115 \pm 0.015$ |
| **CheXpert CLF** | AUC | $0.804 \pm 0.000$ | $0.801 \pm 0.000$ | $0.799 \pm 0.000$ | $0.794 \pm 0.001$ | $0.854 \pm 0.000$ | $0.844 \pm 0.000$ |
| | Acc. | $0.732 \pm 0.001$ | $0.731 \pm 0.001$ | $0.878 \pm 0.000$ | $0.878 \pm 0.000$ | $0.916 \pm 0.000$ | $0.915 \pm 0.000$ |
| | Sens. | $0.686 \pm 0.021$ | $0.670 \pm 0.006$ | $0.102 \pm 0.014$ | $0.122 \pm 0.007$ | $0.144 \pm 0.029$ | $0.161 \pm 0.017$ |
| | Spec. | $0.764 \pm 0.014$ | $0.772 \pm 0.003$ | $0.991 \pm 0.002$ | $0.989 \pm 0.001$ | $0.989 \pm 0.003$ | $0.987 \pm 0.002$ |
| | Prec. | $0.665 \pm 0.006$ | $0.668 \pm 0.001$ | $0.634 \pm 0.021$ | $0.619 \pm 0.007$ | $0.563 \pm 0.013$ | $0.541 \pm 0.010$ |
| | F1 | $0.675 \pm 0.007$ | $0.669 \pm 0.003$ | $0.175 \pm 0.020$ | $0.204 \pm 0.009$ | $0.228 \pm 0.034$ | $0.248 \pm 0.018$ |

Table T.4: Prediction performance original versus orthogonalized data on the MIMIC and CheXpert datasets. The table shows the mean and standard deviation over 10 randomly initialized runs for the labels *Pleural Effusion*, *Cardiomegaly* and *No Finding* and the metrics AUC, accuracy, sensitivity, specificity, precision, and F1-score.